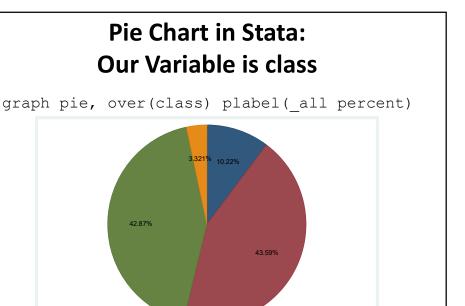
Describing Distributions with Graphs





For all: frequency distribution, percentages

		· ·	•
Variable Type	Central Tendency	Dispersion	Graphs
Nominal	mode	variation ratio	bar, pie
Ordinal	median, mode, (sometimes: mean)	range, interquartile range	bar, pie
Interval/ratio	mean, median (preferred if outliers), mode (rarely useful)	range, interquartile range, standard deviation, variance	histogram, frequency polygon, ogive, boxplot, stem- and-leaf



WORKING CLASS

UPPER CLASS

Exporting Stata Graphs

• Export into png format

graph export class_pie.png

LOWER CLASS

MIDDLE CLASS

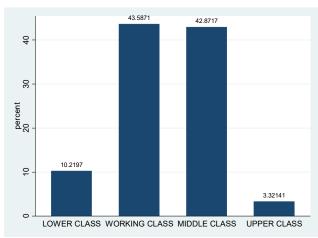
- When in Word, use Insert -> Picture to add graphs
- If file already exists and want to overwrite: use replace option

graph export class pie.png, replace



Bar Graph in Stata: Our Variable is class

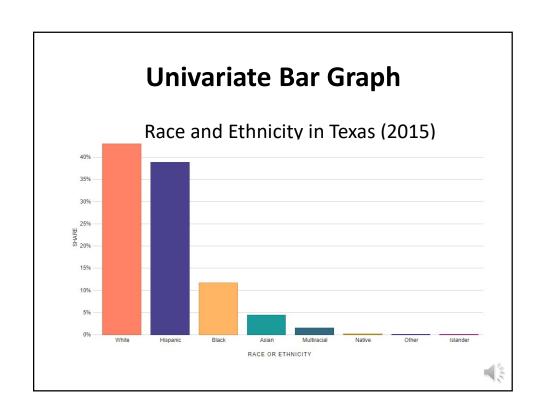
graph bar, over(class) blabel(bar)

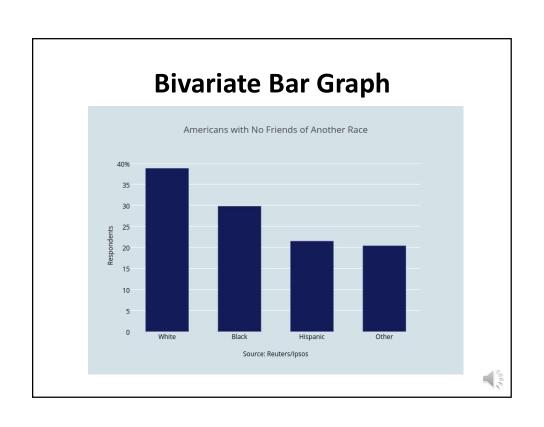


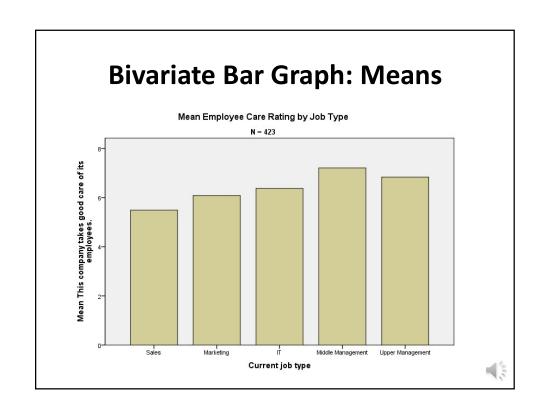
Univariate vs Bivariate Bar Graph

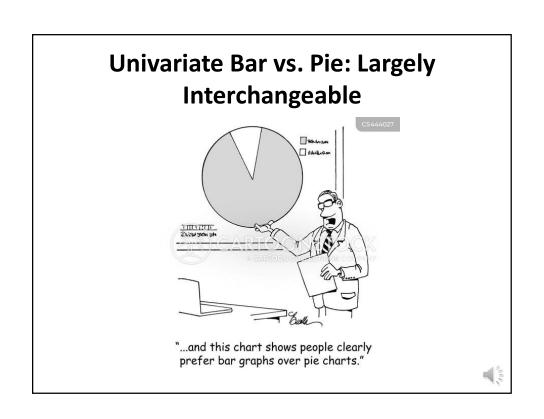
- Univariate Bar Graph = percentages for categories of one nominal/ordinal variable → add up to 100%
- Bivariate Bar Graph = represents the mean or percentage of something calculated separately for each group (2 variables – one main variable and one group variable) → does not add up to 100%

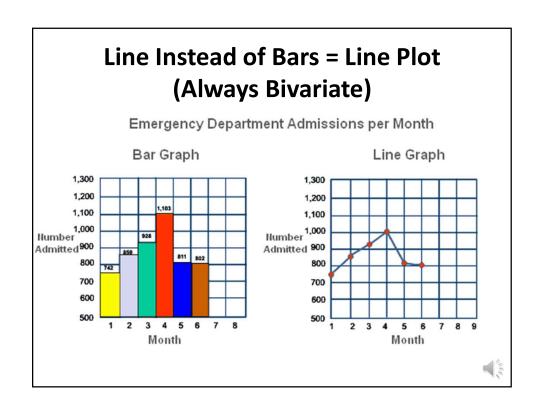


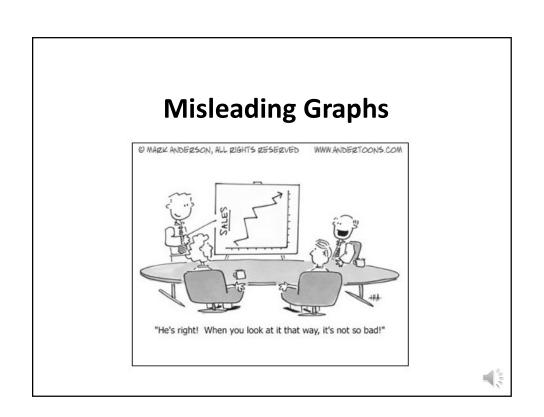










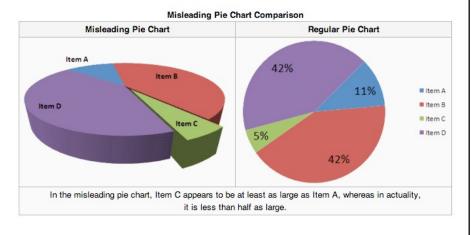


Typical Problems with Graphs

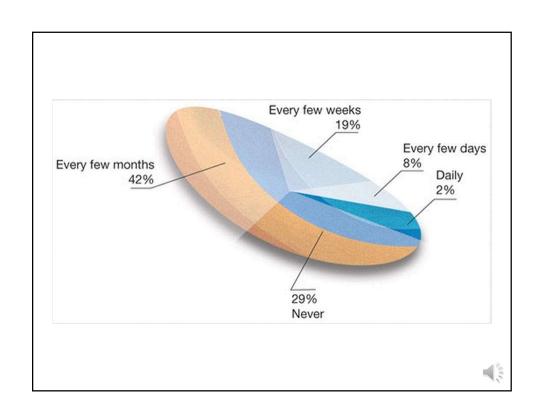
- The proportional ink principle is violated
- The graph isn't labeled properly
- The vertical scale is too big or too small, or skips numbers, or doesn't start at zero
- The horizontal scale is unevenly spaced
- · Data are left out

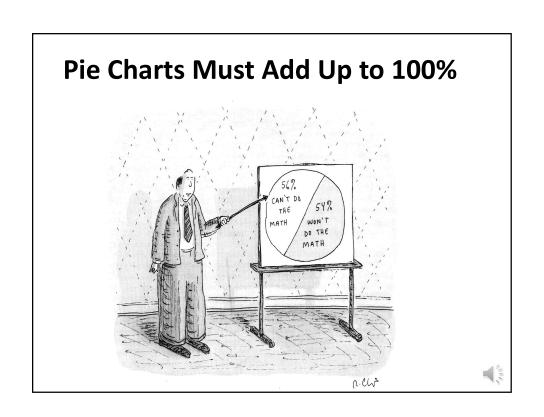


3D Pie Charts: Violating Proportional Ink Principle

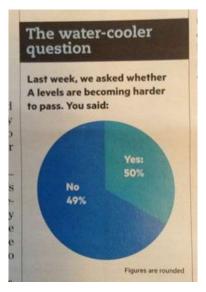






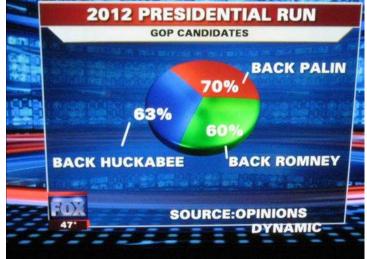




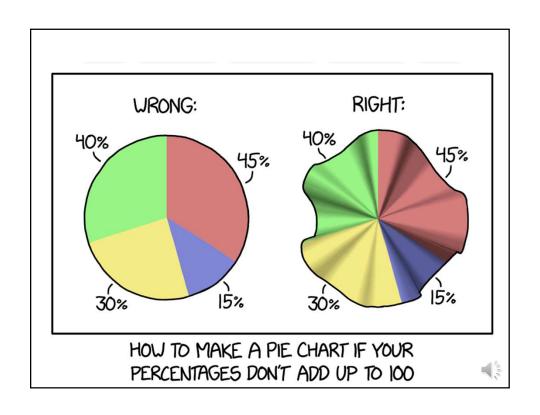


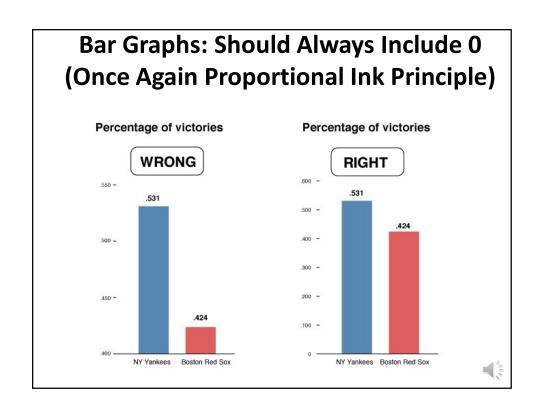
Source: https://www.shinobicontrols.com/blog/6-common-mistakes-with-data-visualization

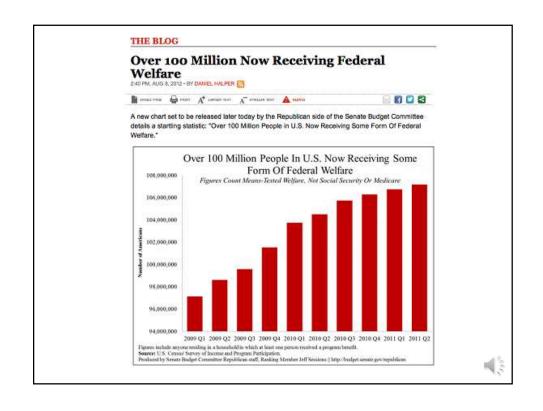
Don't Use Pie Charts With Non-Mutually Exclusive Categories



Source: https://www.shinobicontrols.com/blog/6-common-mistakes-with-data-visualization

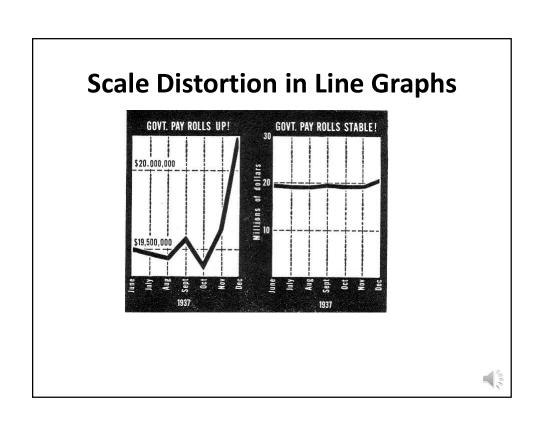


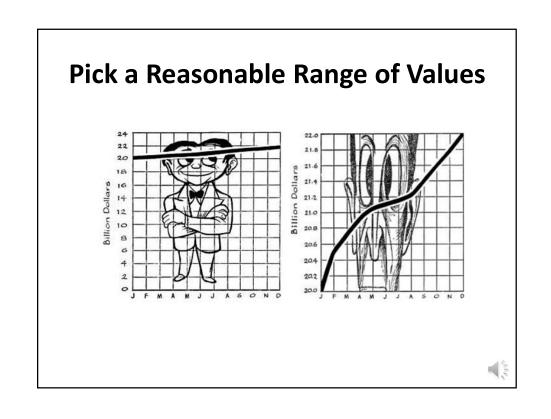


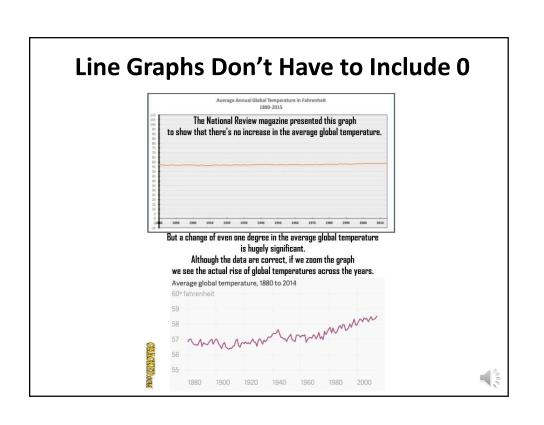


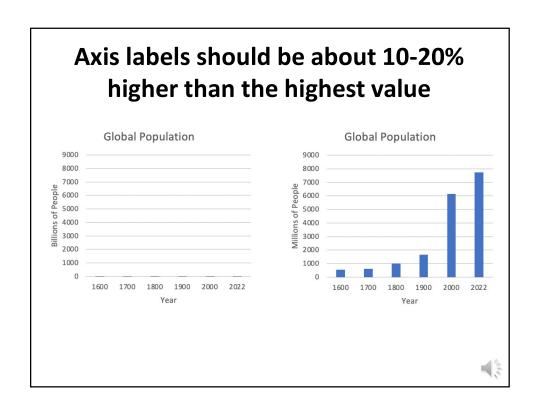






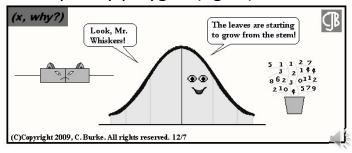






Now For the More Technical Plots that Rarely Make News

- Histogram
- Stem-and-leaf plot
- Frequency polygon
- Cumulative frequency polygon (ogive)
- Boxplot



Frequency Distribution Is
Too Long for a Bar Graph?

tab tvhours

HOURS PER DAY			
WATCHING TV	Freq.	Percent	Cum.
0	90	6.93	6.93
1	255	19.65	26.58
2	325	25.04	51.62
3	238	18.34	69.95
4	171	13.17	83.13
5	61	4.70	87.83
6	58	4.47	92.30
7	19	1.46	93.76
8	31	2.39	96.15
9	3	0.23	96.38
10	17	1.31	97.69
12	11	0.85	98.54
13	2	0.15	98.69
14	4	0.31	99.00
15	3	0.23	99.23
16	1	0.08	99.31
18	1	0.08	99.38
20	2	0.15	99.54
22	1	0.08	99.61
24	5	0.39	100.00
Total	1,298	100.00	



tab tvhours

HOURS PER DAY			
WATCHING TV	Freq.	Percent	Cum.
Γ 0 Ι	90	6.93	6.93
1	255	19.65	26.58
2	325	25.04	51.62
3	238	18.34	69.95
4	171	13.17	83.13
ſ 5 I	61	4.70	87.83
6	58	4.47	92.30
7	19	1.46	93.76
8	31	2.39	96.15
ا 9 ا	3	0.23	96.38
[10	17	1.31	97.69
12	11	0.85	98.54
13	2	0.15	98.69
14	4	0.31	99.00
ſ15 I	3	0.23	99.23
- 16	1	0.08	99.31
18	1	0.08	99.38
20	2	0.15	99.54
- 22	1	0.08	99.61
24	5	0.39	100.00
Total	1,298	100.00	



Resulting Distribution

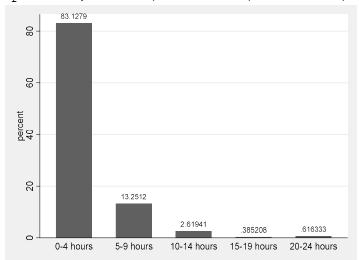
tab tvhours5

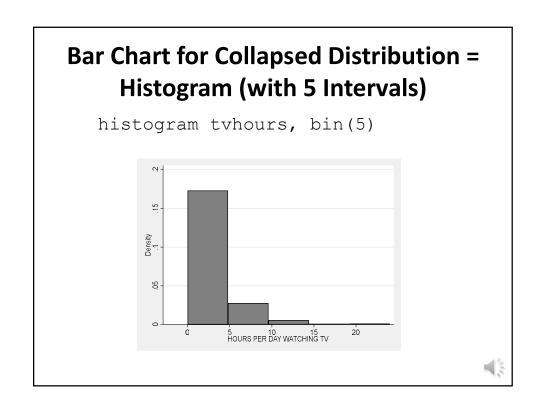
TV hours in intervals of 5 hours	Freq.	Percent	Cum.
0-4 hours 5-9 hours 10-14 hours 15-19 hours 20-24 hours	1,079 172 34 5	83.13 13.25 2.62 0.39 0.62	83.13 96.38 99.00 99.38 100.00
Total	1,298	100.00	

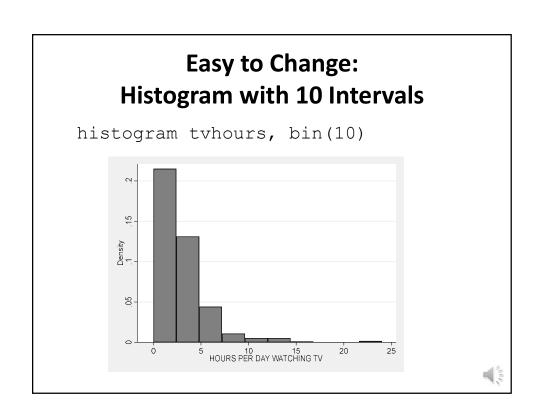


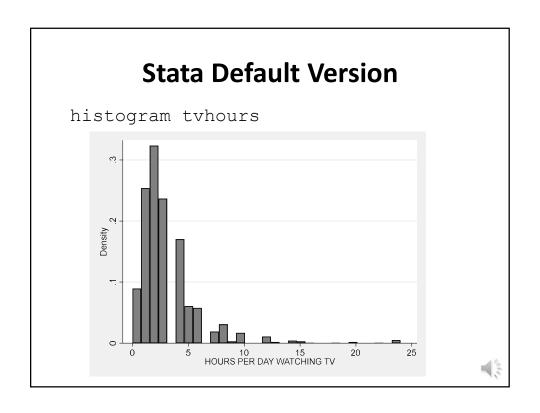
Bar Chart for Collapsed Distribution

graph bar, over(tvhours5) blabel(bar)



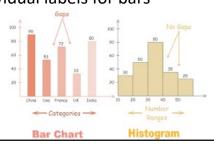


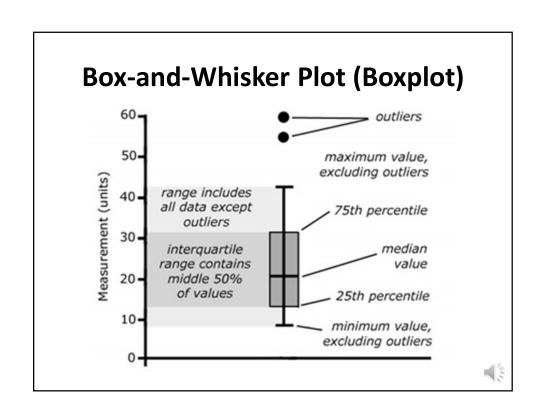


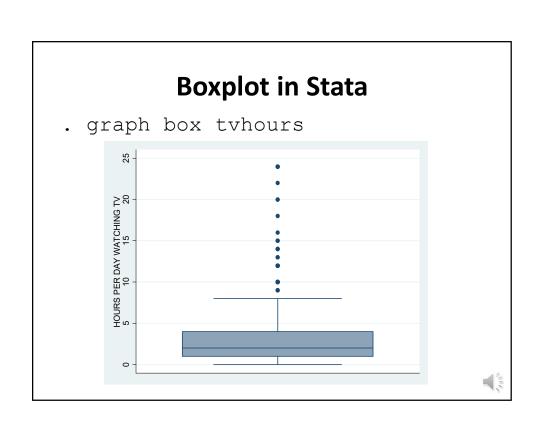


Histogram or Univariate Bar Graph?

- Bar graph (nominal/ordinal variable or 2 vars):
 - Spaces between bars
 - Bars separately labeled
- Histogram (one interval/ratio variable):
 - No spaces between bars
 - No individual labels for bars



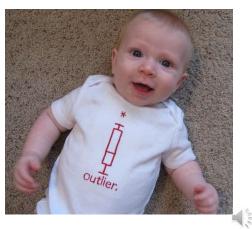


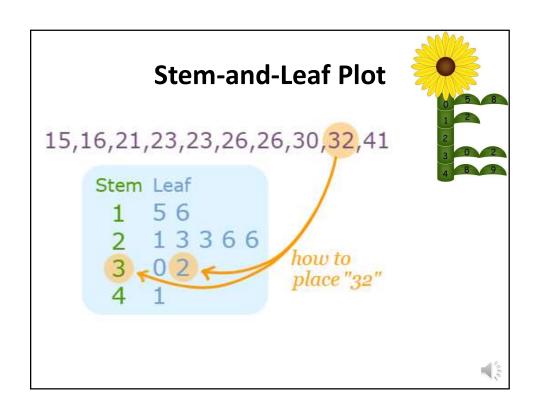


Comparing Distributions graph box emailhr wwwhr EMAIL HOURS PER WEEK WWW HOURS PER WEEK

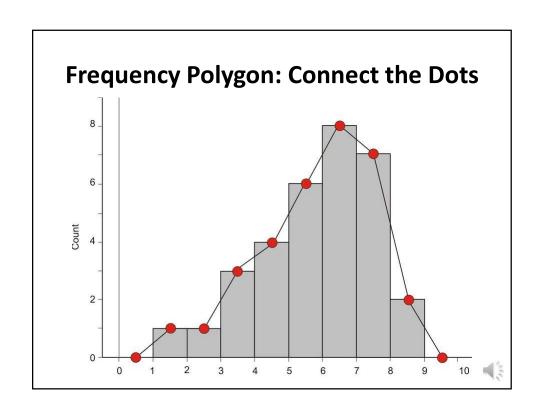
Outliers in Boxplot

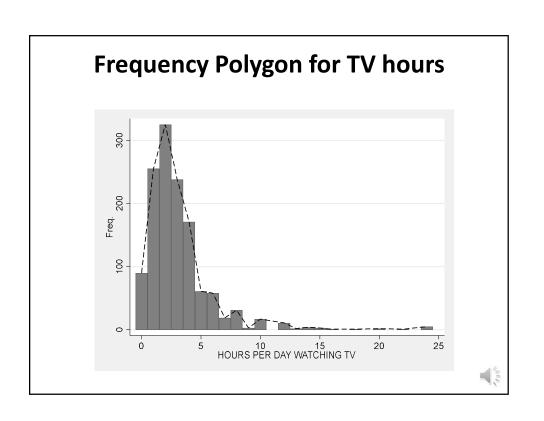
- Below Q1 1.5*IQR
- Above Q3 + 1.5*IQR





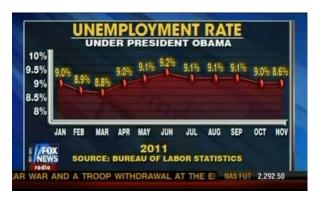
Stem-and-Leaf Plot in Stata . stem tvhours 0* | 7777777777777777777 1* | 00000000000000000 1* | 1* | 2222222222 1* | 4444 1* | 555 1* | 6 1* | 8 1.* | 2* | 00 2* | 2 2* | 2* | 44444

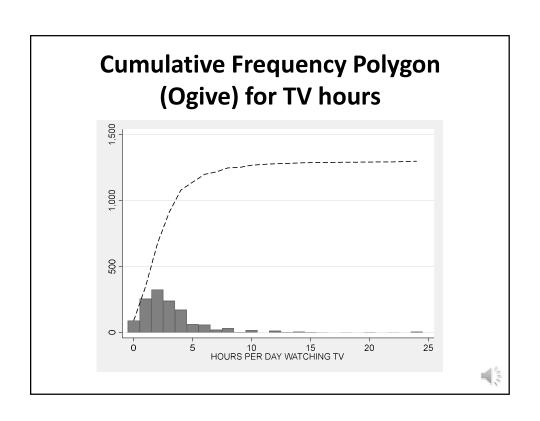


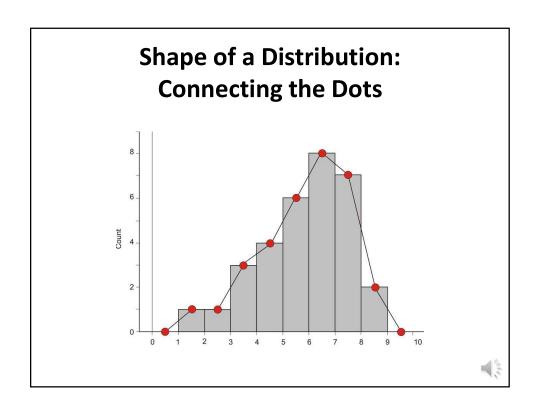


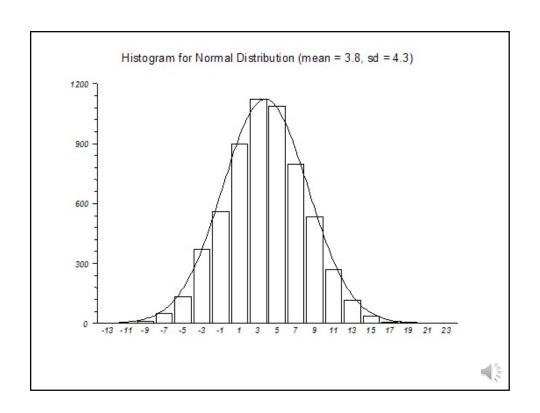
Line Plot vs Frequency Polygon

- Line plot is bivariate (two variables)
- Frequency polygon is univariate (one variable)



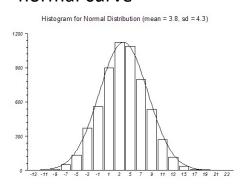




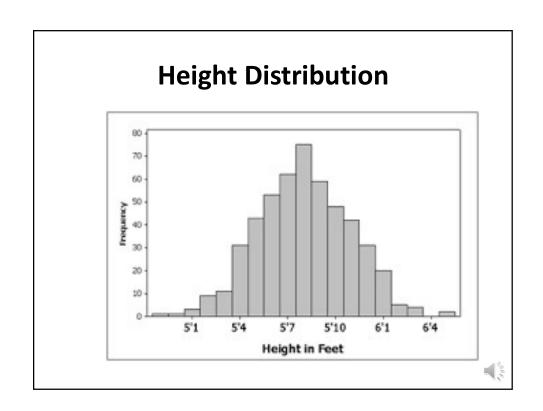


Normal Curve (Bell Curve)

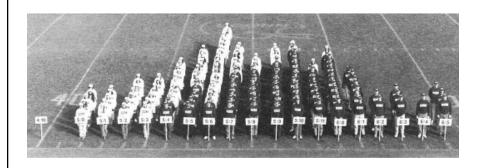
- The most well known "shape" is the normal curve
- Random processes (chance) often result in a normal curve



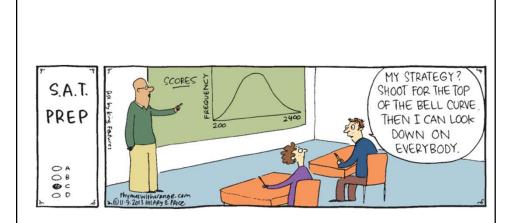




Height Distribution

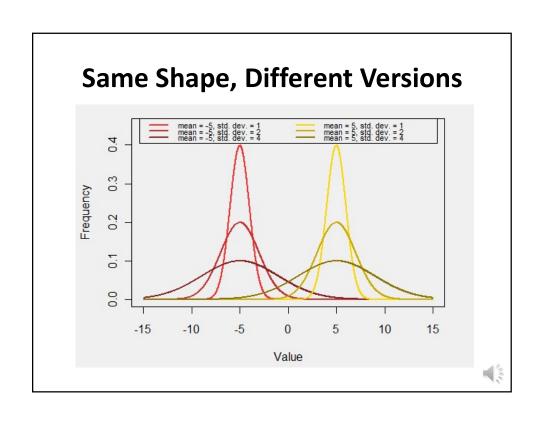


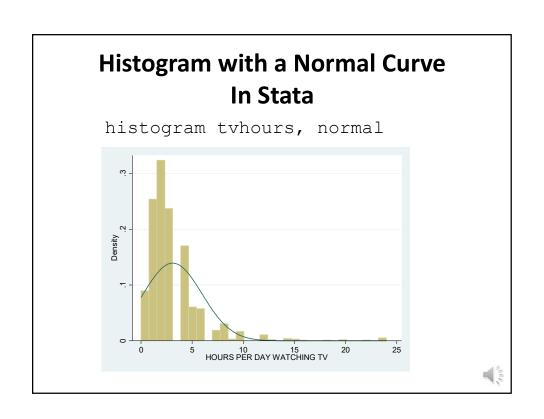


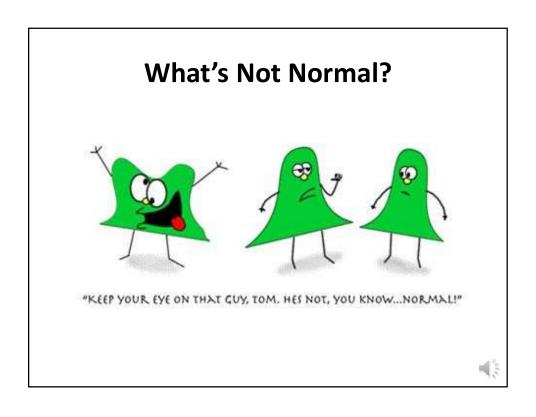


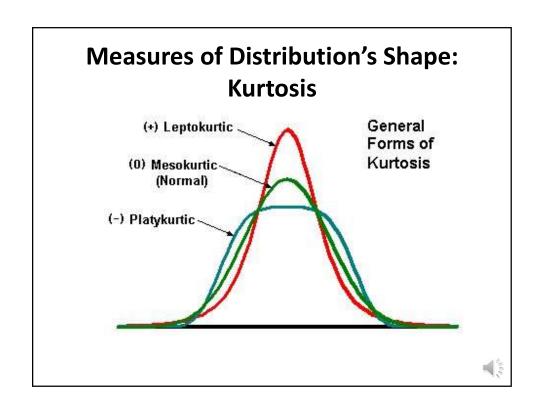
• Normal curve: Mean = median = mode

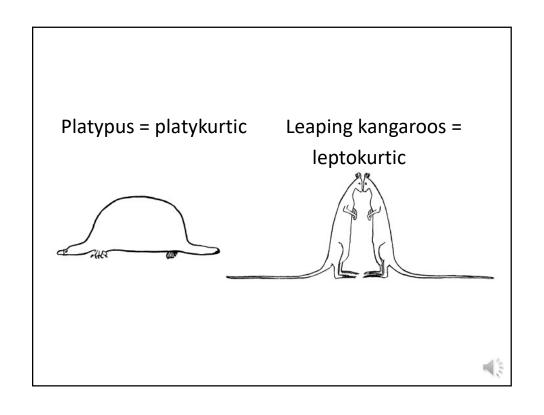


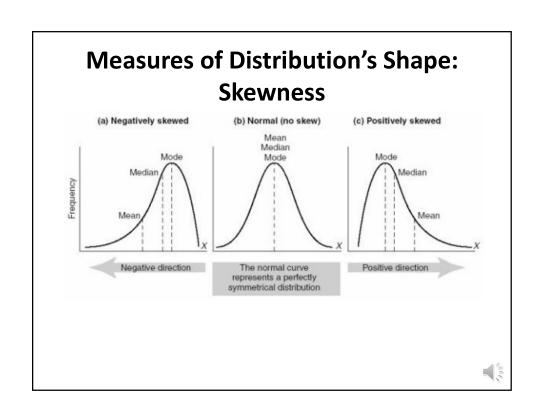












Interpretation of Skewness and Kurtosis Statistics

- Skewness:
 - Near zero = symmetric (close to normal)
 - Positive = right, positive skew (especially >2)
 - Negative = left, negative skew (especially <-2)</p>
- Kurtosis:
 - In Stata: subtract 3 before interpreting!
 - Near zero = close to normal
 - Positive = leptokurtic (especially >2)
 - Negative = platykurtic (especially <-2)</p>



Skewness and Kurtosis in Stata

tabstat tvhours, stats(mean skewness
kurtosis)

