

# Exploratory Data Analysis

2906 Multivariate Statistical Techniques

Session 1

January 24, 2019

## Data Entry + Treatment

### ➤ Take control

- Think through your output in advance
- Study experimental software output files
- Stay flexible (e.g., per trial vs. per person; flipped)

### ➤ Minimize error

- Double check all steps from measurement to analysis
- Stay raw (keep multiple copies when you recode)
- Record value labels (and especially any changes)

## Exploratory Data Analysis

- Understand what you have in front of you
- Active and future-directed; an attempt to understand the process that generates the data and improve future data collections.
- Sometimes meant to generate hypotheses that can be tested preliminarily in the same data set;
- Then need replication in a different data set.

## Data Screening

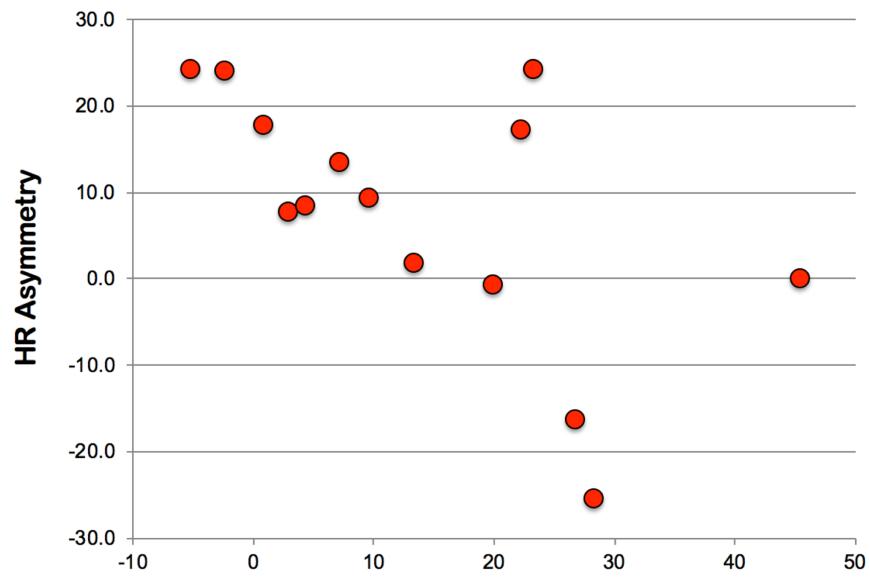
- Accuracy of data entry
  - Order manipulations?
  - Variables within their logical value range?
  - Label values (e.g., for GENDER: 1 = ‘female’ 2 = ‘male’ )
- Outliers
  - Within an acceptable range + frequency? (e.g.,  $z = 2.56$ )
  - Clarify *why* you have outliers
  - In a particular subgroup? (e.g., experimental condition)
  - w/s vs. b/s outliers
- Multivariate outliers?

# Detecting and Handling Outliers

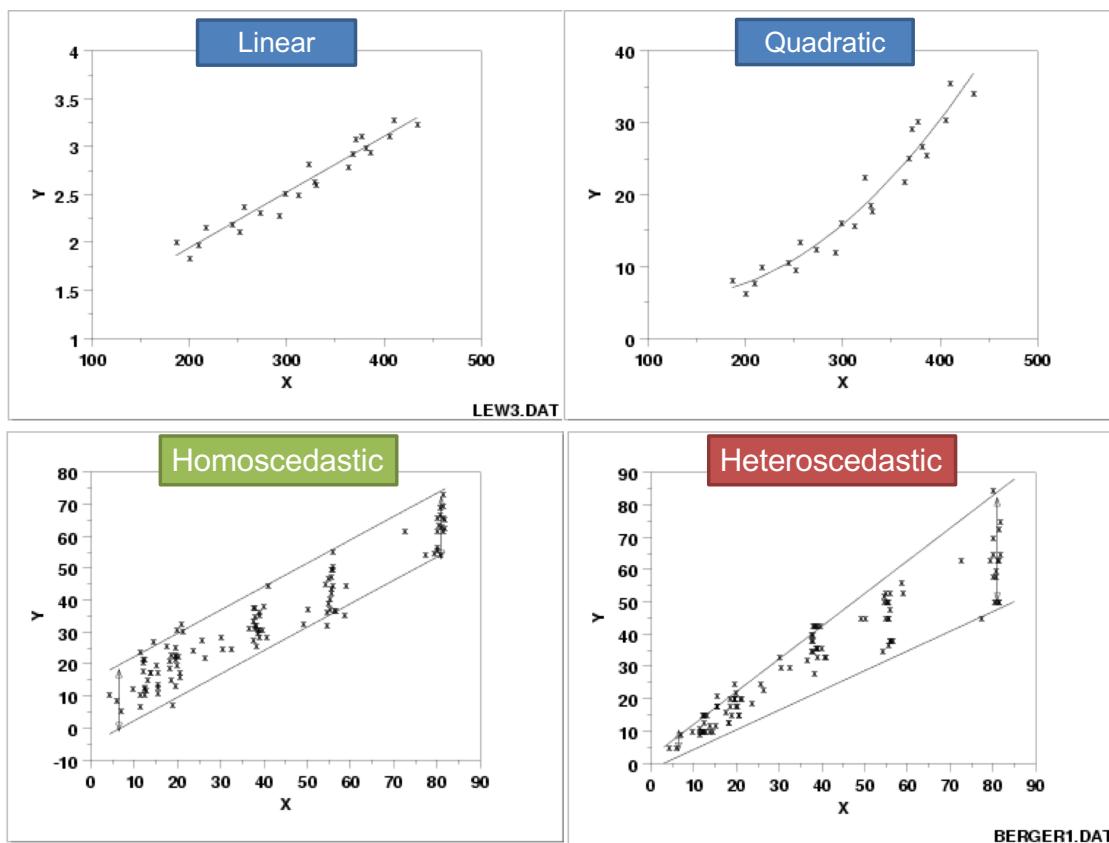
Example from Malle & Holbrook (2012):

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
				<b>z</b>	<b>QQR</b>										
1				-1.26	-2.92										
2		1	1003	-1.36	-2.92										
3		1	1020	-1.31	-2.44										
4		1	1037	-1.11	-2.05										
5		1	1098	-1.10	-2.05										
6		1	1137	-1.00	-1.86										
7		1	1147	-0.97	-1.73										
8		1	1190	-0.86	-1.59										
9		1	1283	-0.62	-1.15										
10		1	1302	-0.44	-0.87										
11		1	1341	-0.46	-0.85										
12		1	1373	-0.73	-0.70										
13		1	1392	-0.33	-0.59										
14		1	1512	-0.01	0.00										
15		1	1519	0.00	0.00										
16		1	1548	0.00	0.18										
17		1	1560	0.17	0.34										
18		1	1565	0.20	0.41										
19		1	1597	0.21	0.42										
20		1	1613	0.26	0.50										
21		1	1615	0.32	0.62										
22		1	1695	0.47	0.91										
23		1	1700	1.28	2.94										
24		1	2272	1.99	3.76										
25		1	2451	2.46	4.65										
26		1	2468	2.48	4.62										
27		1	2488	4.00	7.49										
28		1	2490	4.00	7.49										
29		1	2491	4.00	7.49										
30		1	2492	4.00	7.49										
31		1	2493	4.00	7.49										
32		1	2494	4.00	7.49										
33		1	2495	4.00	7.49										
34		1	2496	4.00	7.49										
35		1	2497	4.00	7.49										
36		1	2498	4.00	7.49										
37		1	2499	4.00	7.49										
38		1	2500	4.00	7.49										
39		1	2501	4.00	7.49										
40		1	2502	4.00	7.49										
41															
42															
43															
44															
45															
46															
47															
48															
49															
50															
51															
52															
53															
54															
55															
56															
57															
58															
59															
60															
61															
62															
63															
64															
65															
66															
67															
68															
69															
70															
71															
72															
73															
74															
75															
76															
77															
78															
79															
80															
81															
82															
83															
84															
85															
86															
87															
88															
89															
90															
91															
92															
93															
94															
95															
96															
97															
98															
99															
100															
101															
102															
103															
104															
105															
106															
107															
108															
109															
110															
111															
112															
113															
114															
115															
116															
117															
118															
119															
120															
121															
122															
123															
124															
125															
126															
127															
128															
129															
130															
131															
132															
133															
134															
135															
136															
137															
138															
139															
140															
141															
142															

Excel



Overall Action-Inaction Discrepancy



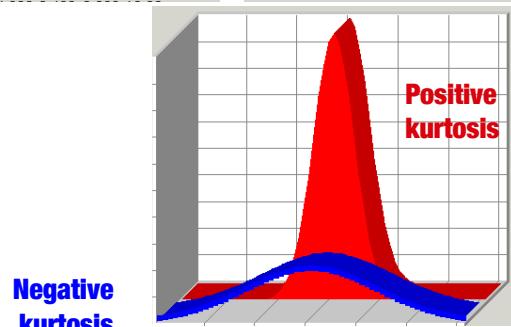
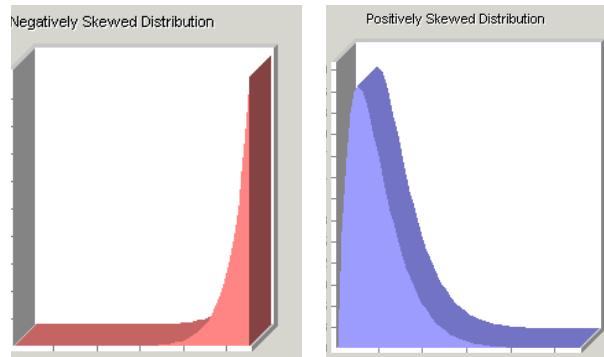
# Distributional Properties

Mean:  $\frac{1}{N} \sum_1^N (X_i)^1$

Variance:  $\frac{1}{N-1} \sum_1^N (X_i - \bar{X})^2$

Skewness:  $\frac{1}{N} \sum_1^N \left( \frac{X_i - \bar{X}}{s} \right)^3$

Kurtosis:  $\frac{1}{N} \sum_1^N \left( \frac{X_i - \bar{X}}{s} \right)^4 - 3$



## Distributional Properties

### EXAMINE

```
VARIABLES=G_GOA_RTY
/PLOT STEMLEAF NPLOT
/STATISTICS EXTREME (5)
/CINTERVAL 95.
```

IQR (Inter-Quartile Range) =<sub>df</sub> value at 75<sup>th</sup> – value at 25<sup>th</sup> percentile

Extremes =<sub>df</sub> values that lie 2\*IQR outside the 25<sup>th</sup>/75<sup>th</sup> percentile

Outliers =<sub>df</sub> values that lie 1.5\*IQR outside the 25<sup>th</sup>/75<sup>th</sup> percentile, but less than 2\*IQR

### Descriptives

	Statistic	Std. Error
G_GOA_RTY Mean	1573.0429	68.75588
95% Confidence Lower Bound	1433.3141	
Interval for Mean Upper Bound	1712.7716	
5% Trimmed Mean	1553.9841	
Median	1528.6667	
Variance	165457.992	
Std. Deviation	406.76528	
Minimum	973.00	
Maximum	2516.00	
Range (IQR)	1543.00	
Interquartile Range	497.33	
Skewness	.699	.398
Kurtosis	.202	.778

### Extreme Values

	Case Number	Value
G_GOA_RTY Highest	20	2516.00
	2	2513.33
	3	2325.33
	4	2201.50
	5	2124.67
Lowest	1	973.00
	2	976.00
	3	988.33
	4	1099.33
	5	1128.00

### G\_GOA\_RTY Stem-and-Leaf Plot

Frequency	Stem & Leaf
3.00	0 . 999
13.00	1 . 0111223333444
14.00	1 . 5555666777788
3.00	2 . 123
2.00	Extremes (>=2513)

Stem width: 1000.00  
Each leaf: 1 case(s)

# Dotplot

XGRAPH  
CHART=[POINT] BY rt9[s]  
/DISPLAY DOT=ASYMMETRIC.

```

XGRAPH CHART=yyvars BY xvars BY zvars

/BIN START={AUTO**} SIZE={AUTO**}
{x } {WIDTH (x)}
{ } {COUNT (n)}

/DISPLAY DOT={ASYMMETRIC**}
{SYMMETRIC }
{FLAT }

/DISTRIBUTION TYPE=NORMAL

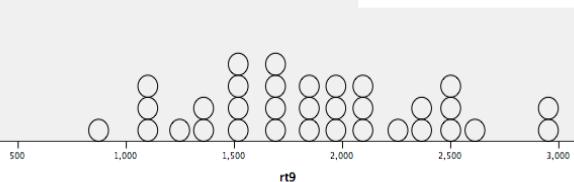
/COORDINATE SPLIT={NO**}
{YES }

/ERRORBAR {CI {(.95)}}
{(n )}
{STDDEV {(2 )}}
{(n )}
{SE {(.2 )}}
{(n )}

/MISSING USE={LISTWISE** } REPORT={NO**}
{VARIABLEWISE } {YES }

/PANEL COLVAR=varlist COLOP={CROSS**} ROWVAR=varlist ROWOP={CROSS**}
{NEST }
{ } {NEST }

```



# Stem-and-Leaf Plot

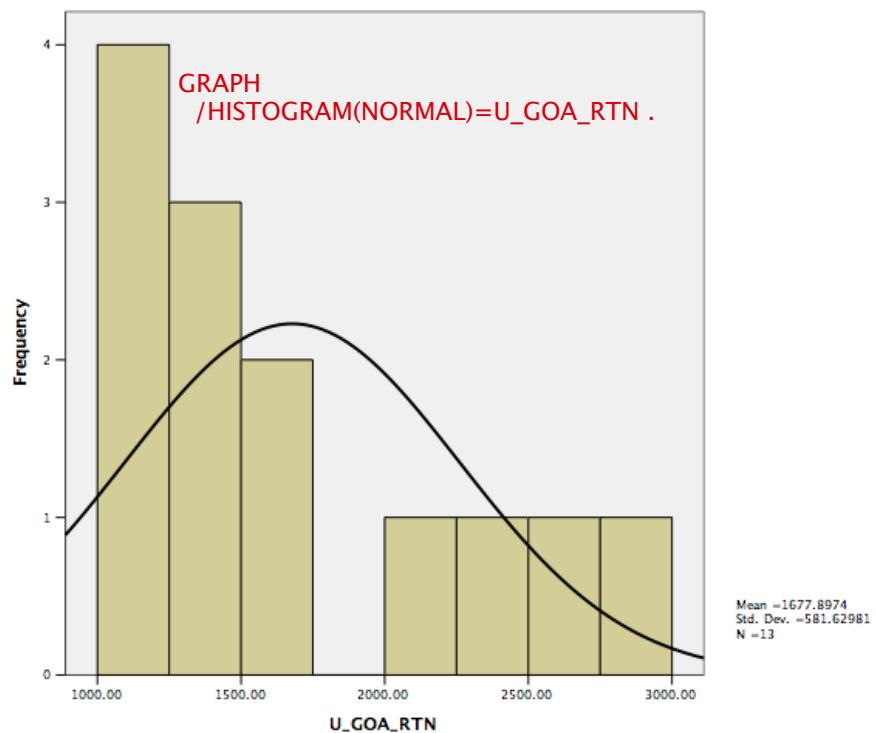
### Frequency Stem & Leaf

Stem width: 1  
Each leaf: 1 case(s)

# Histogram

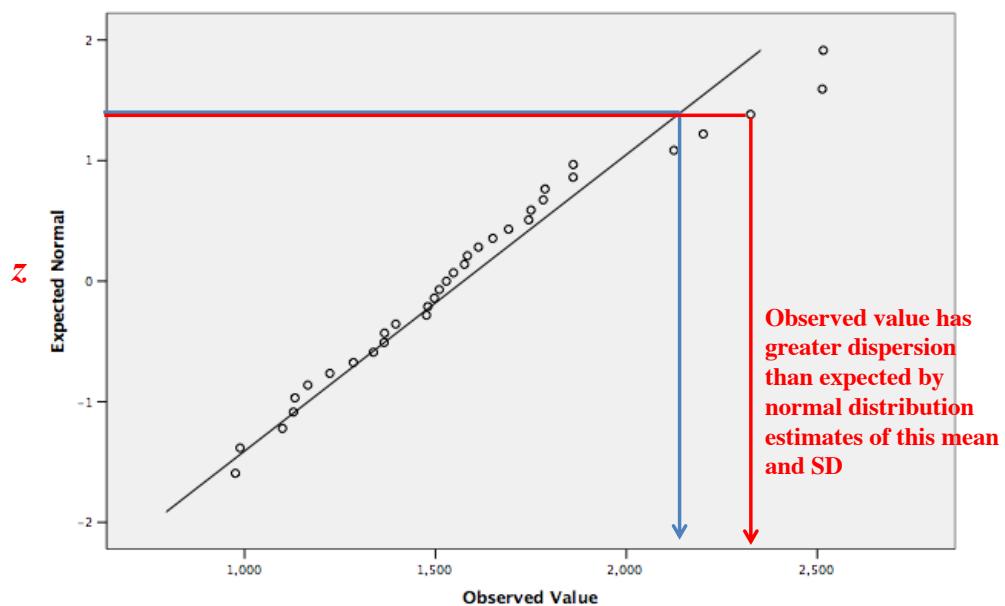
FREQUENCIES VARIABLES=U\_GOA\_RTN  
/HISTOGRAM NORMAL .

U_GOA_RTN		Freq	Percent	Valid Percent	Cumulative
Percent					
Valid	1087.00	1	2.8	7.7	7.7
	1141.00	1	2.8	7.7	15.4
	1180.67	1	2.8	7.7	23.1
	1228.00	1	2.8	7.7	30.8
	1275.00	1	2.8	7.7	38.5
	1343.00	1	2.8	7.7	46.2
	1449.00	1	2.8	7.7	53.8
	1552.50	1	2.8	7.7	61.5
	1703.50	1	2.8	7.7	69.2
	2238.00	1	2.8	7.7	76.9
	2336.00	1	2.8	7.7	84.6
	2511.00	1	2.8	7.7	92.3
	2768.00	1	2.8	7.7	
	Total	13	36.1	100.0	100.0
Missing	System	23	63.9		
	Total	36	100.0		

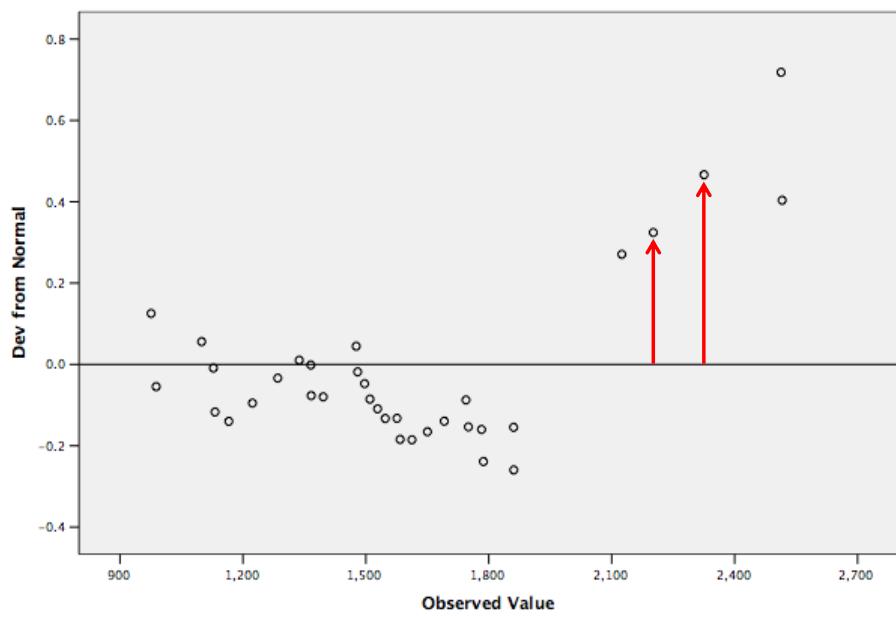


EXAMINE  
 VARIABLES=G\_GOA\_RTY  
 /PLOT STEMLEAF NPLOT  
 /STATISTICS EXTREME (5)  
 /CINTERVAL 95.

Normal P-P Plot of G\_GOA\_RTY



Detrended Normal Q-Q Plot of G\_GOA\_RTY



# Dealing with Skewness and Kurtosis

## ➤ Deviations from normality assumption

- Significance test:  $parameter/s_e = +/- 2 \Rightarrow p < .05$
- If significant, think about its meaning
- Consider aggregation, transformation, non-parametrics
  - aggregation ~ averaging of variables
  - helpful because skew and kurtosis of individual variables can cancel out
  - aggregated variables must be on the same scale

## ➤ In SPSS, aggregating can be done this way:

`COMPUTE newvar = mean.2(oldvar1, oldvar2, oldvar3).`

# Transformations?

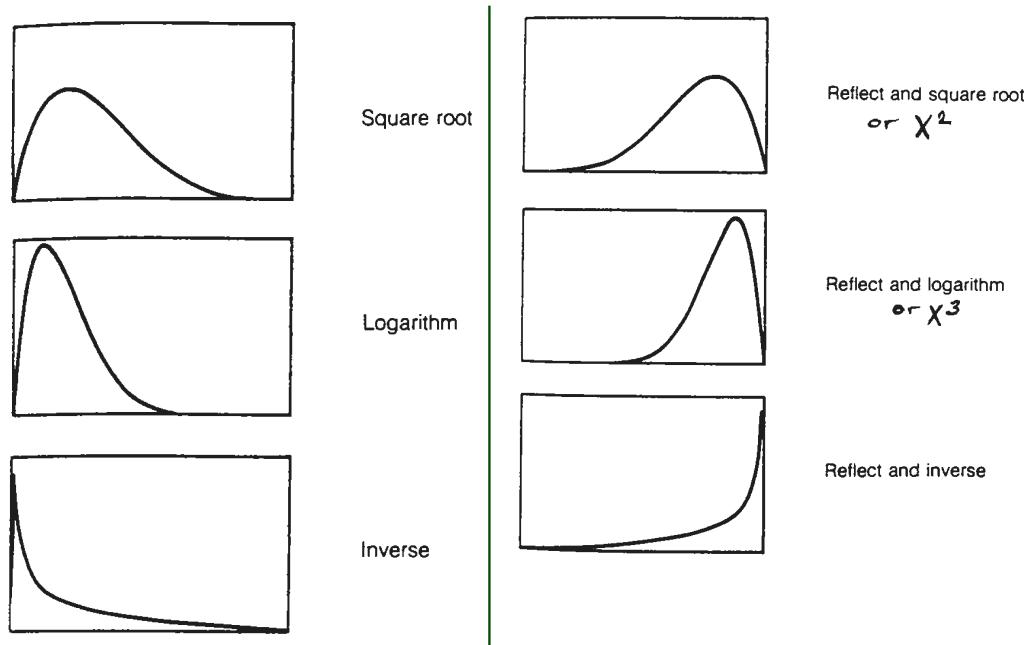
TABLE 4.2 CONTROL LANGUAGE FOR COMMON DATA TRANSFORMATIONS

	BMDP /TRANSFORM	SPSS* COMPUTE	SAS DATA Procedure	SYSTAT Data Module
Moderate positive skewness	NEWX = SQRT(X).	NEWX = SQRT(X)	NEWX = SQRT(X)	LET NEWX = SQR(X)
Substantial positive skewness	NEWX = LOG(X).	NEWX = LG10(X)	NEWX = LOG10(X)	LET NEWX = LOG(X)/LOG(10)
With zero	NEWX = LOG(X + C).	NEWX = LG10(X + C)	NEWX = LOG10(X + C)	LET NEWX = LOG(X + C)/LOG(10)
Severe positive skewness, L-shaped	NEWX = 1/X.	NEWX = 1/X	NEWX = 1/X	LET NEWX = 1/X
With zero	NEWX = 1/(X + C).	NEWX = 1/(X + C)	NEWX = 1/(X + C)	LET NEWX = 1/(X + C)
Moderate negative skewness	NEWX = SQRT(K - X).	NEWX = SQRT(K - X)	NEWX = SQRT(K - X)	LET NEWX = SQR(K - X)
Substantial negative skewness	NEWX = LOG(K - X).	NEWX = LG10(K - X)	NEWX = LOG10(K - X)	LET NEWX = LOG(K - X)/LOG(10)
Severe negative skewness, J-shaped	NEWX = 1/(K - X).	NEWX = 1/(K - X)	NEWX = 1/(K - X)	LET NEWX = 1/(K - X)

C = a constant added to each score so that the smallest score is 1.

K = a constant from which each score is subtracted so that the smallest score is 1; usually equal to the largest score + 1.

From Tabachnick & Fidell and Handout 1

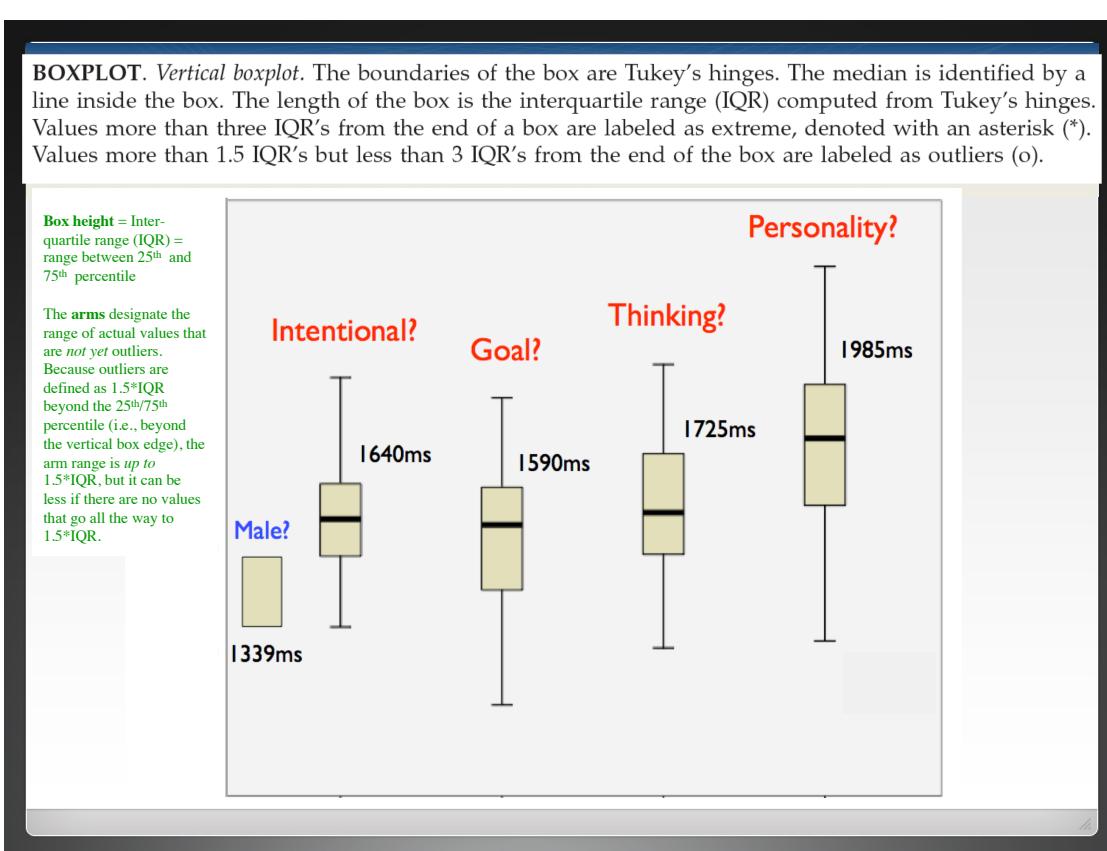


## Missing Data

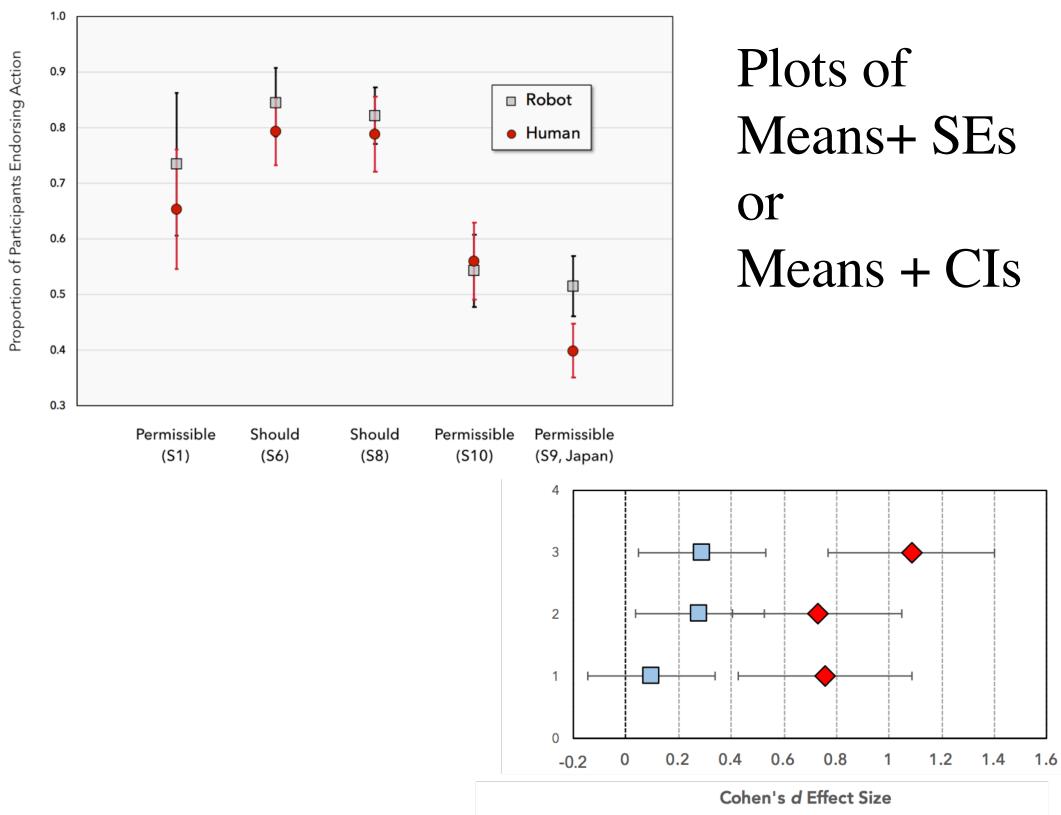
- Are they random or meaningful? (e.g., part of experimental group, indicating misunderstanding)
- Ways of dealing with missing data
  - Exclude the entire case or a specific value.
  - Substitution (e.g., by M/Md of the entire sample; by M/Md of the subgroup of which the case is a member; by regressing).
- Golden Rule: Compare analyses with or without outliers, before and after transformation, etc.
  - Create dummy variable for cases with substituted values.

Example from Malle & Holbrook (2012):

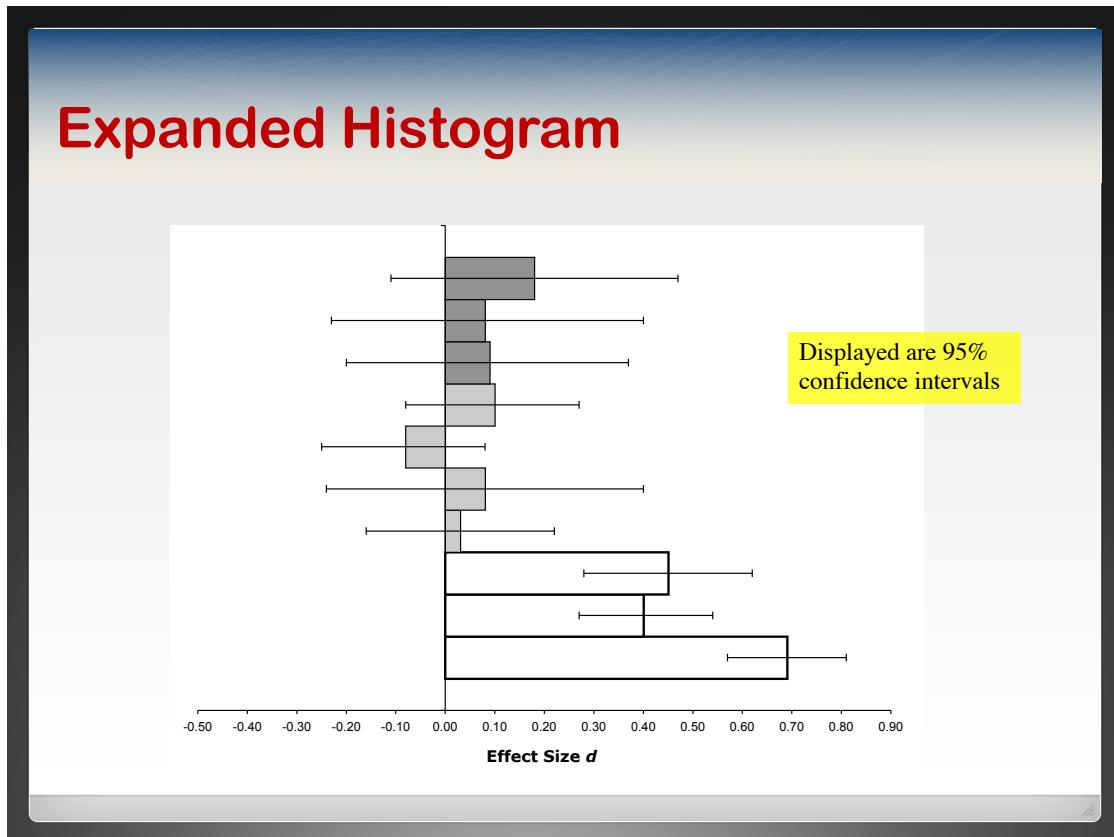
**Missing value replacements.** Some replacements for missing reaction times were necessary when individuals who had valid reaction times for most of the 12 design cells would have been omitted from the within-subject ANOVA because of one or more cells in which they had not provided any Yes responses and, therefore, had no score of inference speed. To retain 36 participants, we used cell-based sample means to replace 57 missing values within the matrix of 432 potential values (36 participants  $\times$  12 cells). To retain the same 36 participants in the analysis of inference likelihood, four entries were mean replaced. It should be noted that these replacements do not change the cell means of the design but somewhat lower their standard deviations.<sup>5</sup> (See supplementary material for all studies' means and standard deviations.)



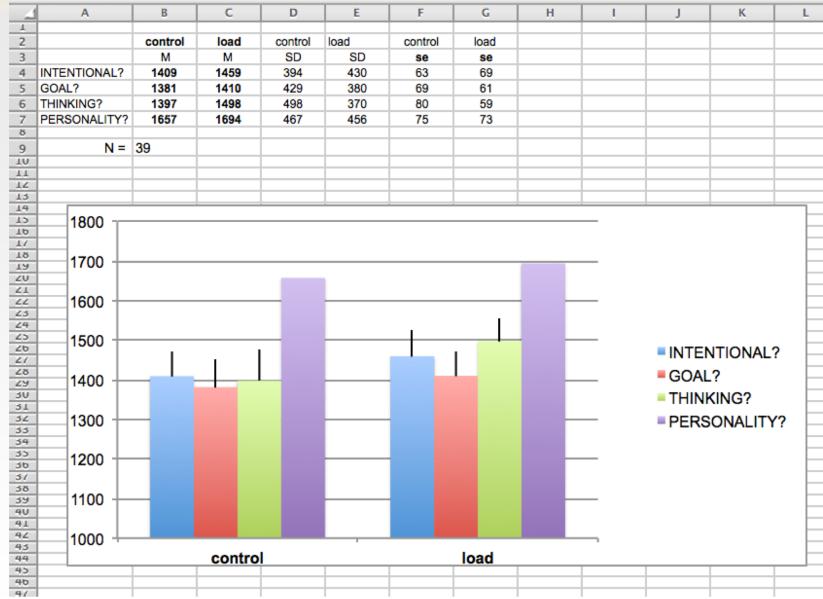
# Plots of Means+ SEs or Means + CIs



## Expanded Histogram



# Getting Those Error Bars



## Other Helpful Commands for EDA

- CROSSTABS
  - CORRELATION (PARTIAL CORR)
  - MEANS TABLES

## Graphing software options:

- SPSS, with editing
  - Excel, with patience (and Keynote or PPT on top of it)
  - JMP, R, many others (check what's available on Brown's software server)