A decorative graphic consisting of a large number of small, multi-colored dots (red, orange, yellow, green, blue) arranged in a curved, downward-sloping pattern that tapers to the right.

Tips, Tricks and Time-Savers: Using Statgraphics Operators

Presented by Dr. Neil W. Polhemus

Statgraphics Operators

- Operators are predefined functions that perform specific tasks.
- They are used in:
 - Data fields on analysis input dialog boxes.
 - “Select” field on analysis input dialog boxes.
 - Data editor when using “Generate Data”.
 - Procedures such as *Nonlinear Regression* to specify statistical models.

Types of Operators

- Algebraic operators (+, -, /, *, ^)
- Mathematical transformations (LOG, SQRT, ABS, ROUNDT0, ...)
- Sequential operators (RUNTOT, DIFF, SDIFF, ...)
- Random number generators (RNORMAL, RUNIFORM, ...)
- Statistical summaries (MEAN, SD, MEDIAN, ...)
- Distribution functions (NORMAL, INVNORMAL, ...)
- Boolean operators (<, >, =, <=, >=, <>, |, &)
- Data selectors (TAKE, TAKELAST, DROP, ...)
- Pattern generators (COUNT, REP, RESHAPE, ...)
- Utility functions (REPLACE, JOIN, JUXTAPOSE, ENDSWITH, ...)

Reference

- Select *Help – Procedure Documentation* from the main menu. Click on the PDF file titled “STATGRAPHICS Operators”.

PROPER(x)

Purpose: converts each string in a character column to a proper name by capitalizing the first letter of each word in the string.

Type: utility function

Argument: data column

Example: `PROPER(make)`

Result: column of strings

RNORMAL(n,mu,sigma)

Purpose: generates random numbers from a normal distribution

Type: random number generator

Argument: sample size, mean, standard deviation

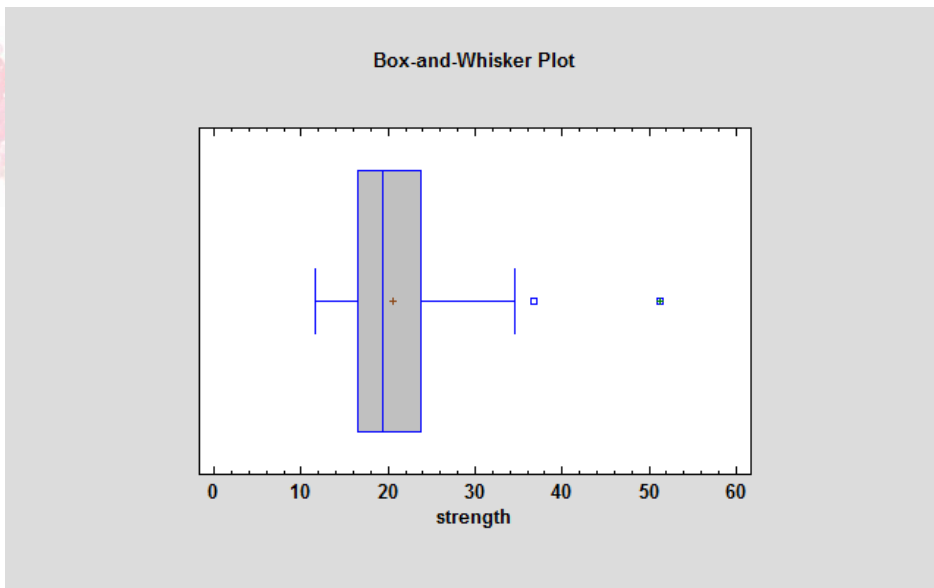
Example: `RNORMAL(3,10,3)`

Result: 13.4892 9.85616 11.9911

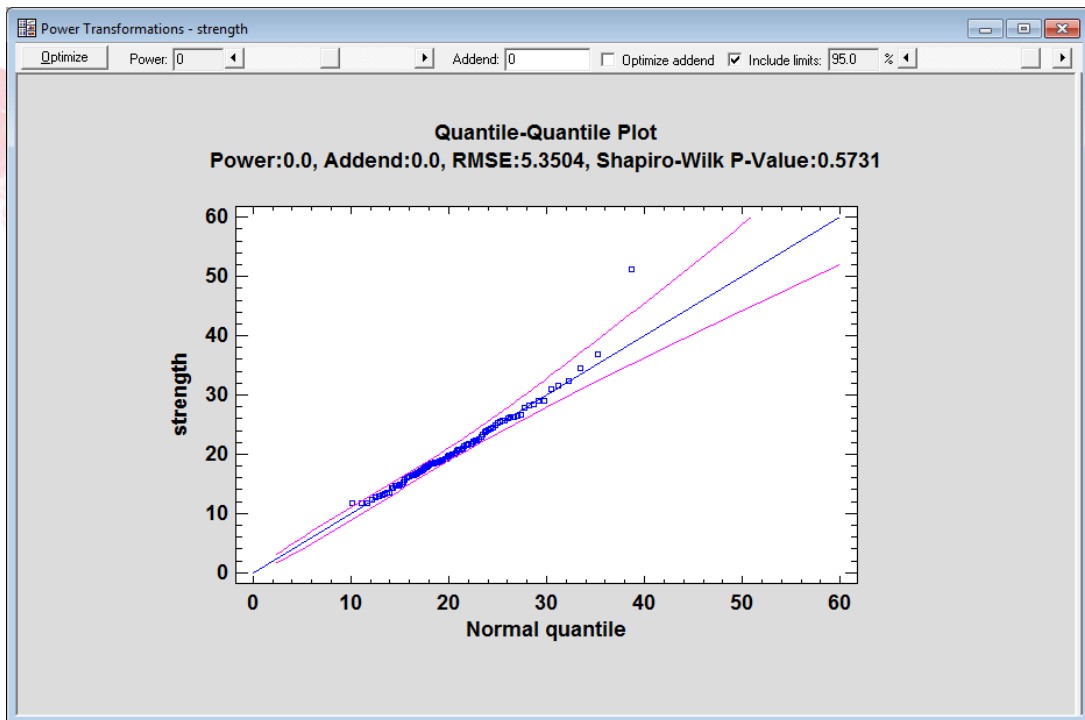
Application: Transform data to satisfy assumptions of method to be used

- Example #1: I have a set of data that may contain outliers. I'd like to use Grubbs' test, which assumes the data come from a normal distribution. My data appear to be skewed.
- Example #2: I've created a designed experiment. The data I wish to analyze are counts, which are likely to follow a binomial distribution in which the variance is a function of the mean.

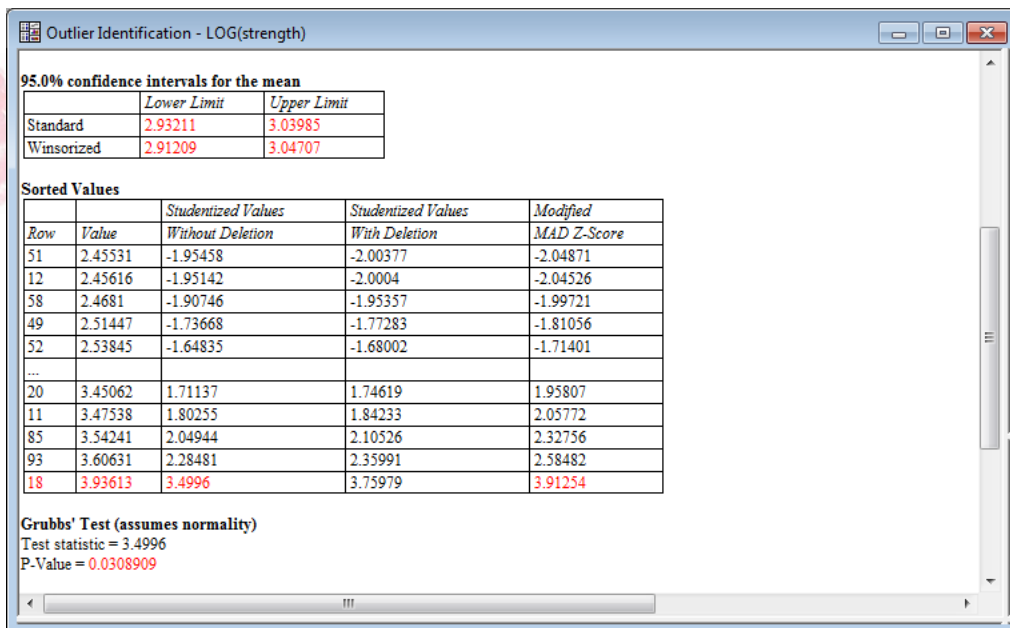
Example 1: Tensile Strength



Power Transformations Statlet



Outlier Identification



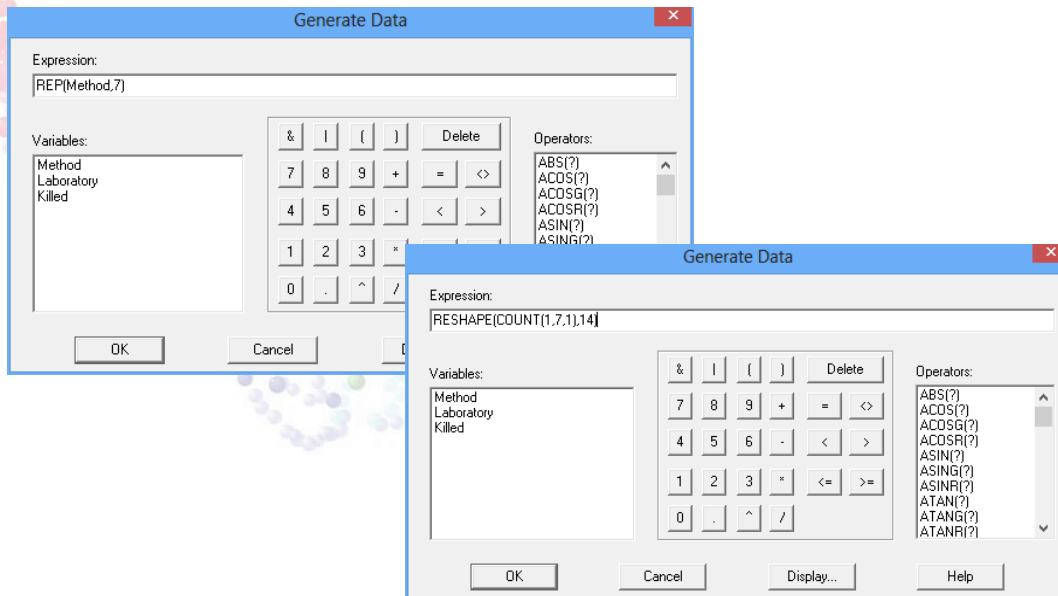
Example 2: Mothproofing

- Data: percentage of 20 moth larvae killed

	Lab 1	Lab 2	Lab 3	Lab 4	Lab 5	Lab 6	Lab 7
Method A	40	35	5	80	50	95	45
Method B	60	30	15	95	75	100	55

- Goal: compare the 2 methods
- Source: Statistics for Experimenters by Box, Hunter and Hunter

Data Setup



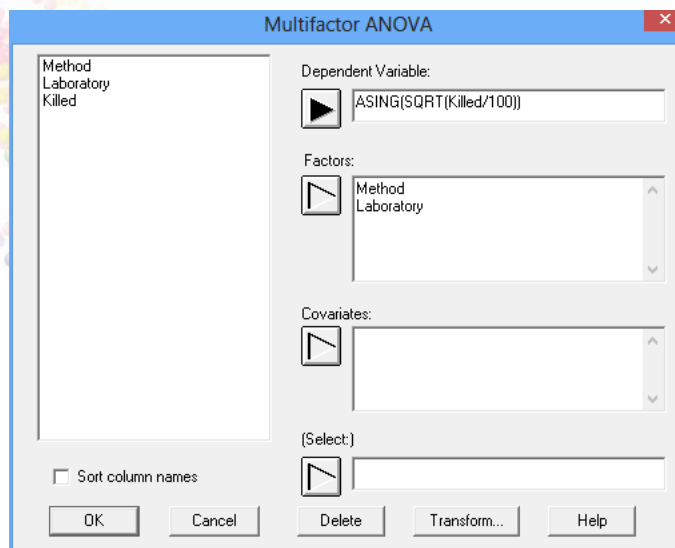
Data

	Method	Laboratory	Killed	Col_4	Col_5	Col_6
1	A	1	40			
2	A	2	35			
3	A	3	5			
4	A	4	80			
5	A	5	50			
6	A	6	95			
7	A	7	45			
8	B	1	60			
9	B	2	30			
10	B	3	15			
11	B	4	95			
12	B	5	75			
13	B	6	100			
14	B	7	55			
15						
16						
17						
18						
19						
20						

Multifactor ANOVA

- Fisher's variance stabilizing transformation:

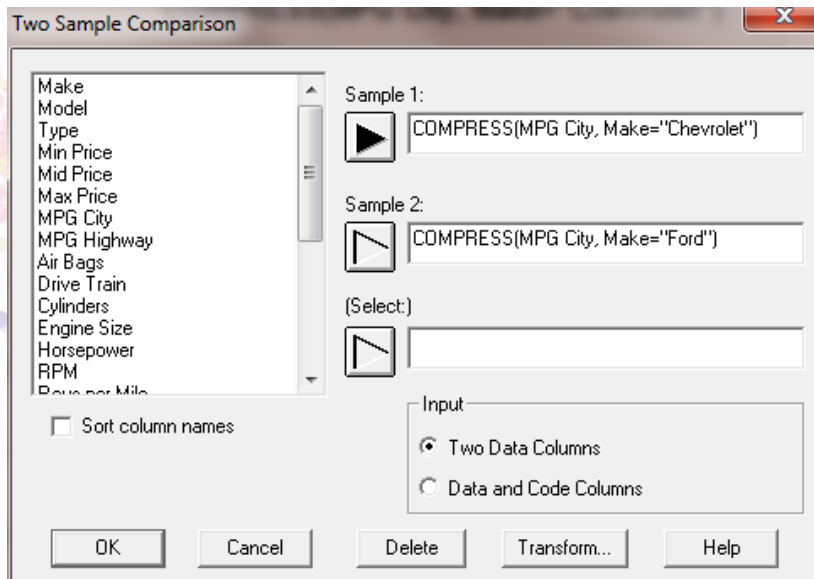
$$x = \sin^{-1} \sqrt{\hat{p}}$$



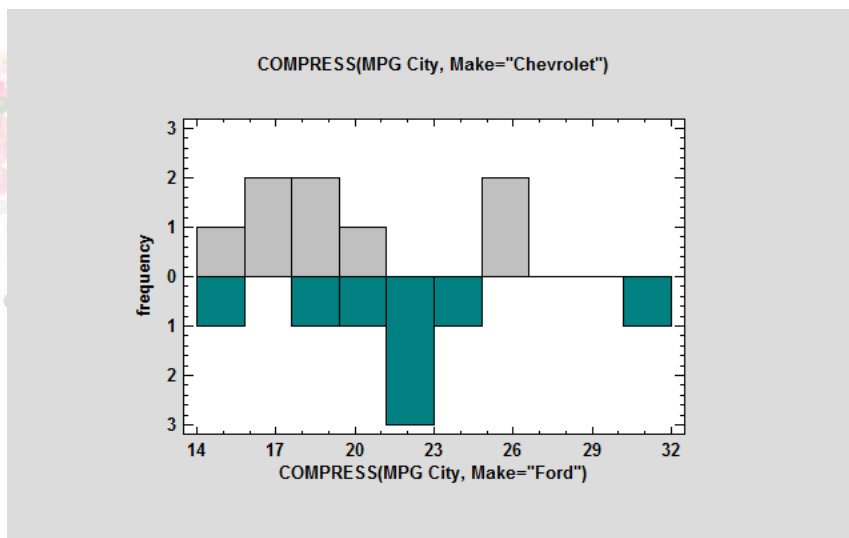
Application: Boolean operators

- Example #1: I have a data set containing information on 93 makes and models of automobiles. How can I compare Fords and Chevrolets?
- Example #2: I've collected data for a stability study regarding the loss of chlorine over time in a product I produce. How can I fit a piecewise linear model to the data?

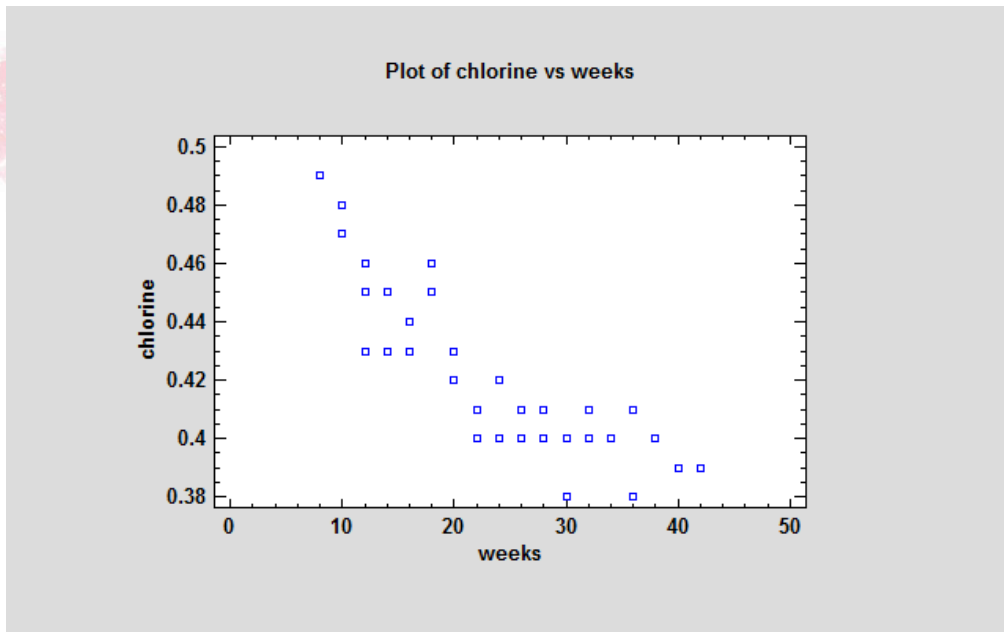
Two Sample Comparison



Histogram

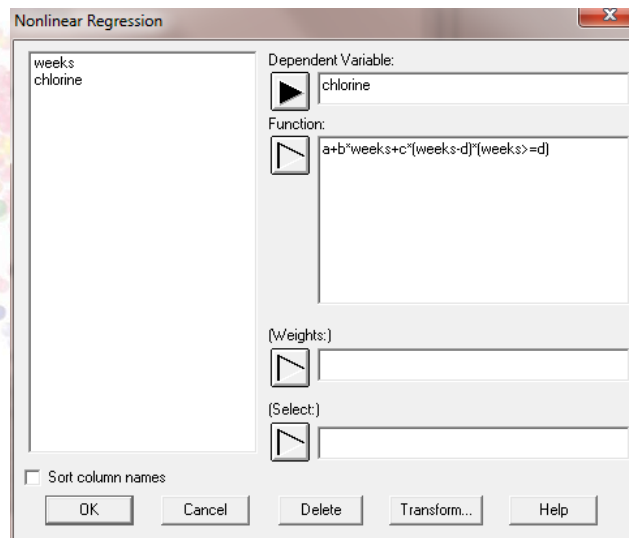


Stability Study

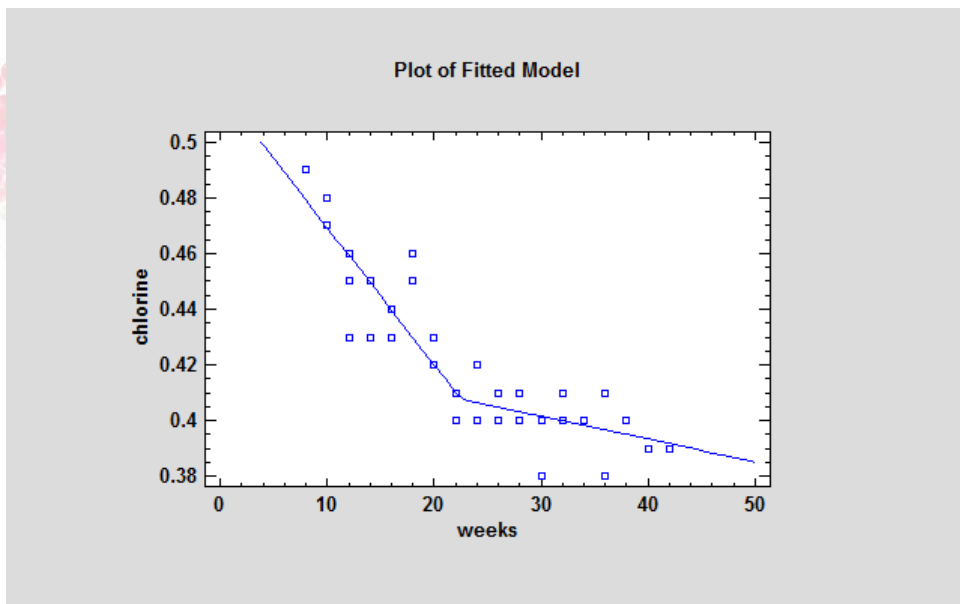


Nonlinear Regression Model

$$\text{chlorine} = a + b * \text{weeks} + c * (\text{weeks} - d) * (\text{weeks} \geq d)$$



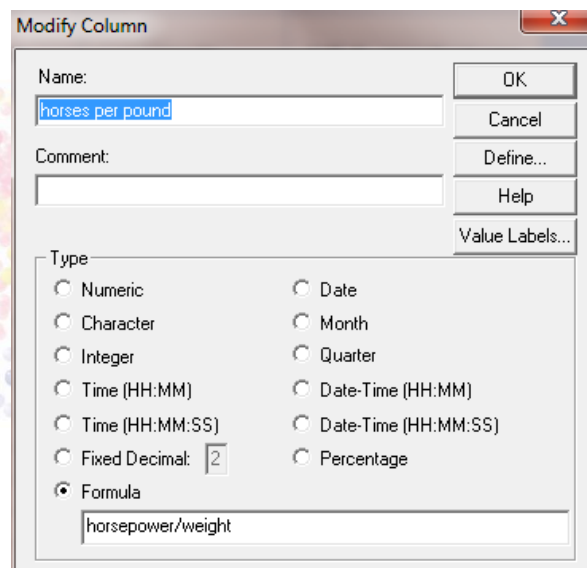
Piecewise Linear Fit



Application: Combining Columns

- Example #1: I will take my automobile data and construct a column consisting of the ratio of horsepower to weight.
- Example #2: I will join end to end my 3 columns with price.
- Example #3: I will combine “make” and “model” into a single column.

Defining a Function



Modify Column

Name:

Comment:

Type

Numeric Date

Character Month

Integer Quarter

Time (HH:MM) Date-Time (HH:MM)

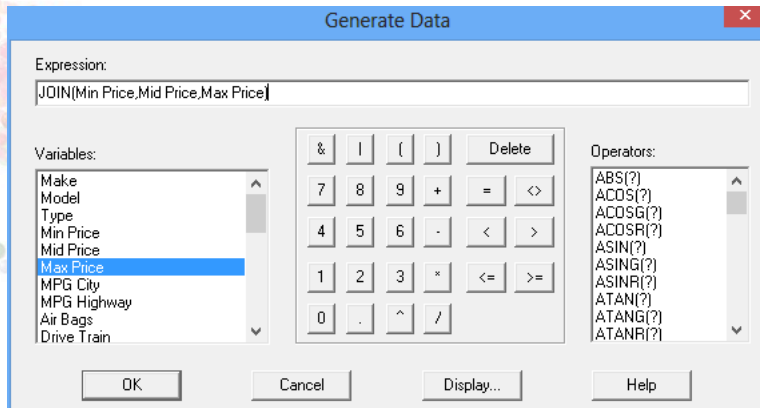
Time (HH:MM:SS) Date-Time (HH:MM:SS)

Fixed Decimal: Percentage

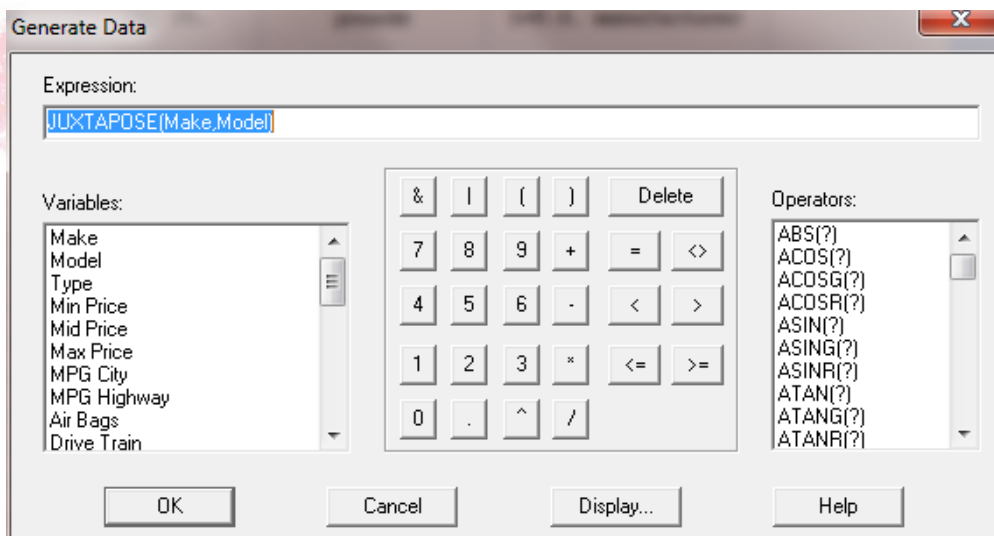
Formula

OK
Cancel
Define...
Help
Value Labels...

Generate Data



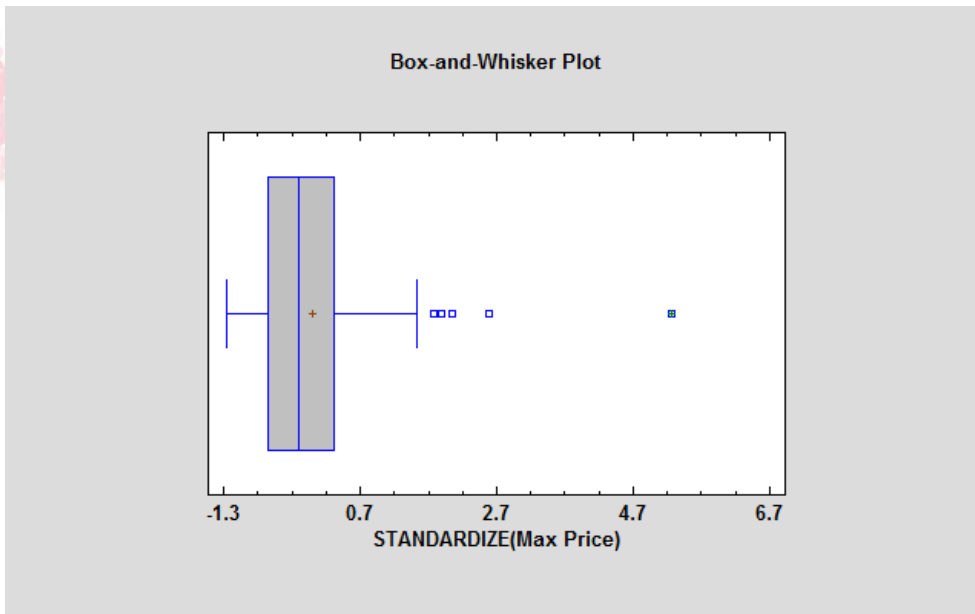
Generating Data



Application: Converting data to Z-scores

- Example: I've created a box-and-whisker plot for a sample of n observations and wish to determine how many standard deviations the outside points are from their mean.

Box-and-Whisker Plot



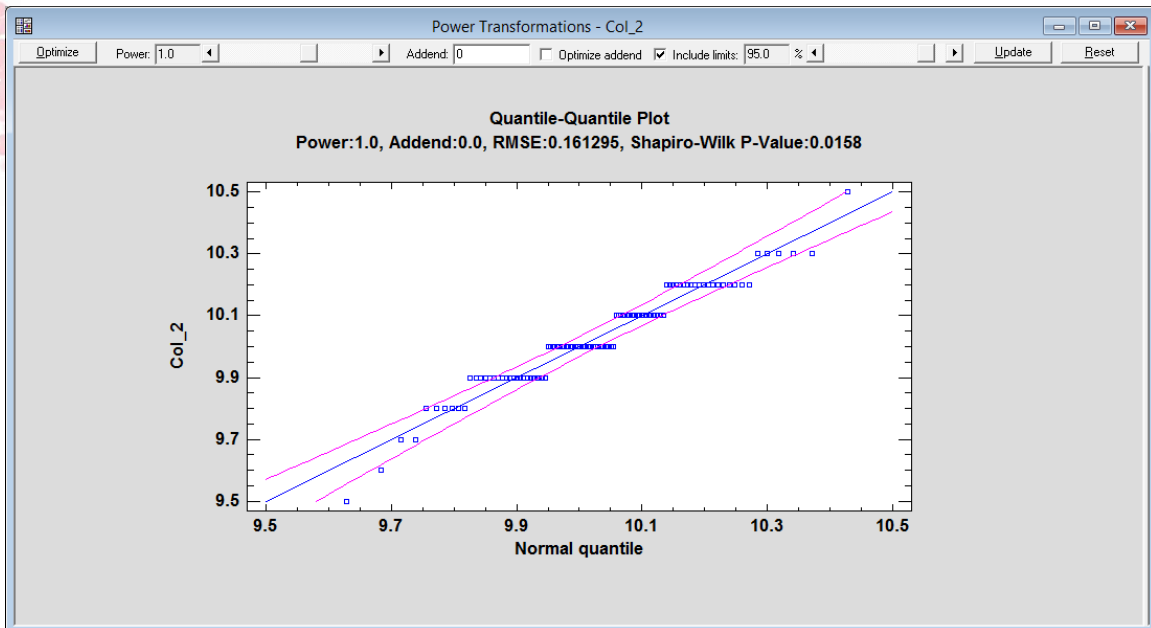
Application: Adding decimal places

- Example: Suppose I have a set of data that I wish to test for normality. It fails the test, but I suspect the failure is due to the fact that the measurements were only recorded to one decimal place. How can I test that suspicion?
- I'll create 3 columns:
 - Col_1: `RNORMAL(100,10,0.15)`
 - Col_2: `ROUNDTO(Col_1,1)`
 - Col_3: `Col_2+RUNIFORM(100,-.05,.05)`

Data

	Col_1	Col_2	Col_3	Col_4	Col_5	Col_6
1	9.87831789002	9.9	9.88			
2	9.86385795257	9.9	9.86			
3	9.91452151717	9.9	9.88			
4	10.0433132177	10	9.95			
5	10.1063153527	10.1	10.07			
6	10.1732957587	10.2	10.23			
7	10.0195983354	10	10.02			
8	10.3104340952	10.3	10.32			
9	9.9270451977	9.9	9.94			
10	10.273711203	10.3	10.27			
11	9.7235226067	9.7	9.71			
12	10.007261109	10	10.00			
13	10.115471377	10.1	10.07			
14	10.179958398	10.2	10.23			
15	9.9127099976	9.9	9.91			
16	10.122309318	10.1	10.15			
17	9.9068657219	9.9	9.89			
18	10.1762600494	10.2	10.24			
19	10.061514864	10.1	10.12			
20	10.221201732	10.2	10.24			

Q-Q Plot



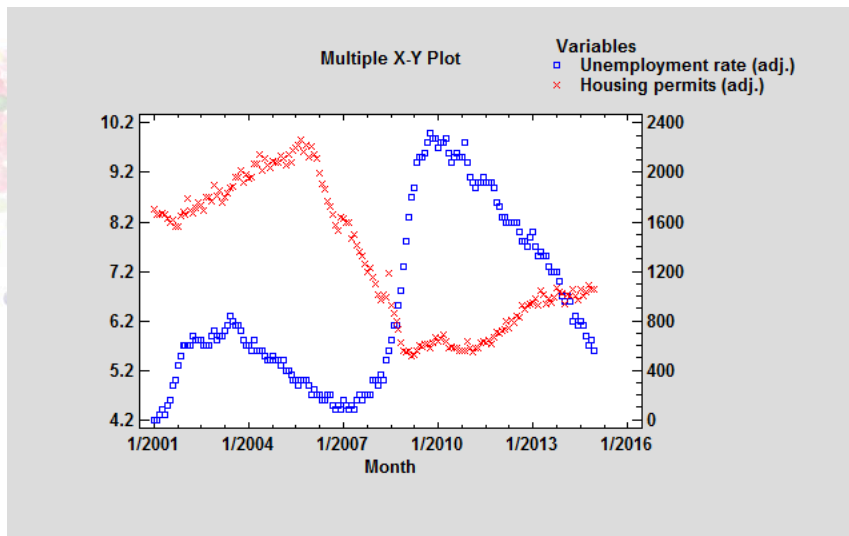
Application: Leading Indicators

- Example: I'd like to forecast the monthly U.S. unemployment rate. I've heard that the number of permits issued for new housing construction is a leading indicator for the economy. Can I use housing permits to improve my forecasts for unemployment?

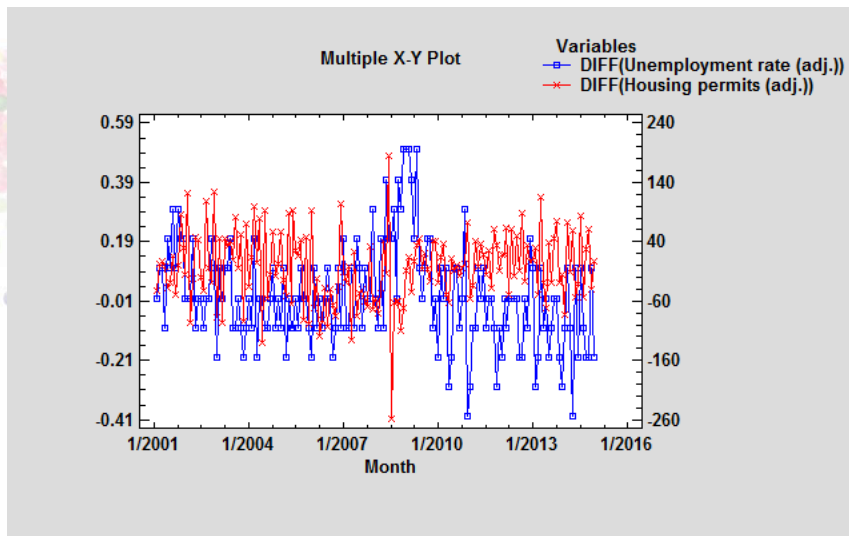
Housing Permits Data

	Month	Unemployment rate (unadj.)	Unemployment rate (adj.)	Housing permits (adj.)	Col_5	Col_6
1	1/2001	4.7	4.2	1699		
2	2/2001	4.6	4.2	1656		
3	3/2001	4.5	4.3	1659		
4	4/2001	4.2	4.4	1666		
5	5/2001	4.1	4.3	1665		
6	6/2001	4.7	4.5	1626		
7	7/2001	4.7	4.6	1598		
8	8/2001	4.9	4.9	1615		
9	9/2001	4.7	5.0	1565		
10	10/2001	5	5.3	1566		
11	11/2001	5.3	5.5	1651		
12	12/2001	5.4	5.7	1680		
13	1/2002	6.3	5.7	1665		
14	2/2002	6.1	5.7	1787		
15	3/2002	6.1	5.7	1691		
16	4/2002	5.7	5.9	1669		
17	5/2002	5.5	5.8	1716		
18	6/2002	6	5.8	1758		

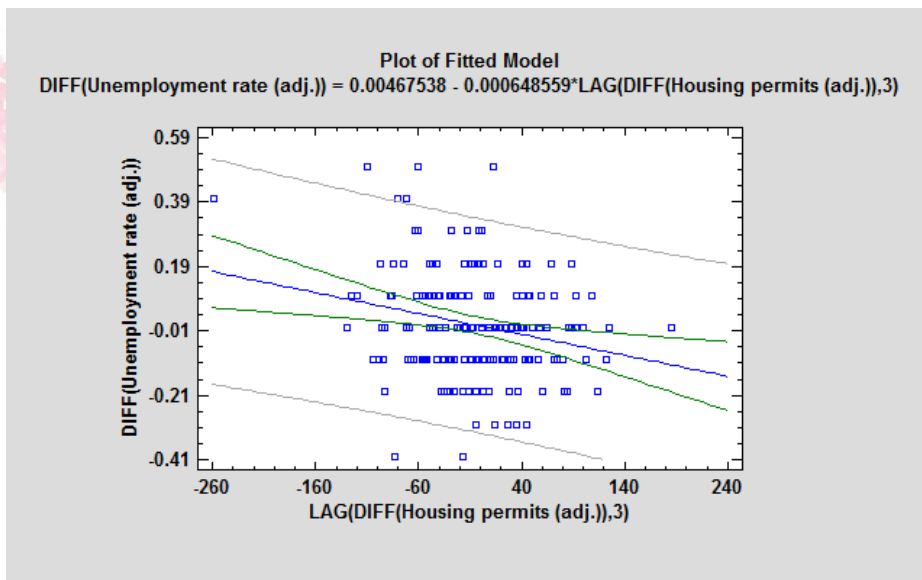
Time Sequence Plot



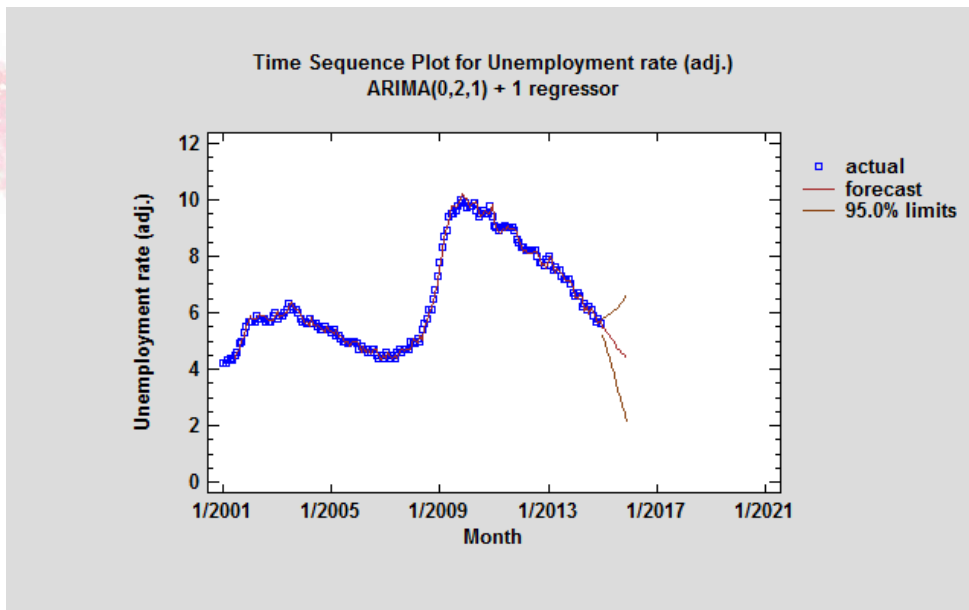
First Differences



Simple Regression



Forecasted Unemployment



Review

EXCLUDE

LOG

REP

RESHAPE

COUNT

ASING

SQRT

COMPRESS

JOIN3

JUXTAPOSE

STANDARDIZE

RNORMAL

ROUNDO

RUNIFORM

DIFF

LAG

Boolean

Algebraic