



COAPS MET3220C & MET6480
Computational Statistics

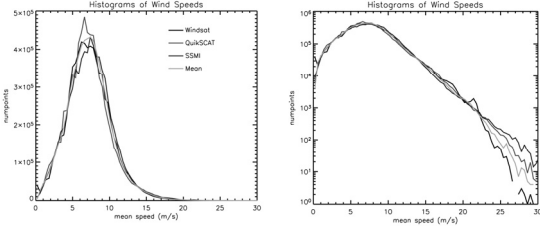
Exploratory Data Analysis
 For Paired Data

Scatterplots
 Correlation (several types)
 Star Plot
 Glyph Scatterplots

Key Point: ALWAYS LOOK AT THE DATA!!!!

<http://campus.fsu.edu/bourassa@met.fsu.edu>  The Florida State University  Data Exploration: Paired Data Sets 1



Wind Speed Histograms Based on Co-located Observations From 3 Satellites



Better for seeing differences in frequently occurring observations

Better for seeing differences in infrequently occurring observations

Graphics from talk by Mike Freilich and Barry Vanhoff

<http://campus.fsu.edu/bourassa@met.fsu.edu>  The Florida State University  Data Exploration: Paired Data Sets 2

Standard Deviation and Variance

- Recall that the standard deviation is defined as



$$s_x = \left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2}$$

- The variance is the square of the standard deviation.

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- The variance is a particularly useful quantity because it is additive in many applications, and the total variance is often preserved.
- For example, if a variable f is dependent on three independent variables x , y , and z , then

$$s_f^2 = s_x^2 + s_y^2 + s_z^2$$

<http://campus.fsu.edu/bourassa@met.fsu.edu>  The Florida State University  Data Exploration: Paired Data Sets 3

Covariance

- Covariance is a measure of sort of like variance.
- However, covariance (cov) examines how one variable changes in proportion to another.



$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- If x' is proportional to y' , then

$$\text{cov}(x, y) = s_x s_y$$

- If x' is independent of y' , then

$$\text{cov}(x, y) = 0$$



<http://campus.fsu.edu/bourassa@met.fsu.edu>  The Florida State University  Data Exploration: Paired Data Sets 4

Pearson (Ordinary) Correlation (AKA Linear Correlation)

- The Pearson correlation assumes that there is (or more accurately could be) a linear relationship between the two variables being considered: $x \propto y$.
- This correlation coefficient is defined as



$$r_{xy} = \frac{\text{covariance}(x, y)}{s_x s_y} = \frac{\frac{1}{n-1} \sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} \left[\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{1/2}}$$

$$r_{xy} = \frac{\sum_{i=1}^n (x'_i y'_i)}{\left[\sum_{i=1}^n (x'^i)^2 \right]^{1/2} \left[\sum_{i=1}^n (y'^i)^2 \right]^{1/2}}$$

<http://campus.fsu.edu/bourassa@met.fsu.edu>  The Florida State University  Data Exploration: Paired Data Sets 5

Properties of the Correlation Coefficient

- $-1 < r_{xy} < 1$
- If $r_{xy} = 1$, then $x' \propto y'$
 - Indicating a positive slope of the best fit line
- If $r_{xy} = -1$, then $x' \propto -y'$
 - Indicating a negative slope of the best fit line
- If $r_{xy} = 0$, then x' is independent of y'
 - The slope of the best fit line is meaningless
- It is often said that r^2 is the fraction of the variance explained by a linear relationship. This is true provided that the uncertainty in both sets of observations is negligible.
 - Another key consideration is that both variables should not be calculated from the same variable or variables.
 - The above problem is called cross correlation, and it results in a much larger correlation than would be otherwise determined.
- Correlation does NOT imply cause and effect.

<http://campus.fsu.edu/bourassa@met.fsu.edu>  The Florida State University  Data Exploration: Paired Data Sets 6

Example Problems

Set I
 $r = 0.88$
 $r^2 = 0.74$

Set II
 $r = 0.61$
 $r^2 = 0.37$
Outlier

- Pearson linear correlation does not work well in either of these examples.
 - Why are there problems?
- Set I: the relationship is substantially non-linear.
 - An engineering solution might be linear fits over several ranges.
- Set II: The outlier leads to a large covariance, resulting in a questionable value for the correlation.

Graphics from Wilk's Statistical Methods in the Atmospheric Sciences
<http://campus.fsu.edu/> Data Exploration: Paired Data Sets 7
 bourassa@met.fsu.edu COAPS The Florida State University

Example of Cross Correlation Non-Dimensional Fetch Relationship

Figure from Hasseleman et al. (1973) and Canady's (2001) Fig. 2.7.

- Kitigorodski has developed a non-dimensional relationship that applies over a wide range of conditions.
- The u_* in both the x and y terms is a serious problem (cross correlation)!

<http://campus.fsu.edu/> Data Exploration: Paired Data Sets 8
 bourassa@met.fsu.edu COAPS The Florida State University

Comparison of PAC3 Stress to Modeled Stress

PACS South stress

Total	R	R ²	regression	RMSE
Taylor and Yelland	0.99	0.97	$y=0.70x-0.00$	0.006
Wu	0.99	0.97	$y=0.79x+0.00$	0.007
Smith et al.	0.99	0.97	$y=0.70x-0.00$	0.007

- All three models are very well correlated to the data.
- Which model is better?
 - T&Y model has lower RMS errors, but Wu's model seems to have a much better slope.

Graphics from talk by Yoshi Goto
<http://campus.fsu.edu/> Data Exploration: Paired Data Sets 9
 bourassa@met.fsu.edu COAPS The Florida State University

Computationally Efficient Correlation Step 1: The Covariance

- Computational efficiency can be ignored for small data sets.
- However, for every large data sets it can be very important.
 - Example: one pass through a data set is used to determine the mean
 - A second pass is used to determine the standard deviation.
 - If the data set is read in each time, then the process is miserably slow!
- Consider the covariance squared:

$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(x, y) = \frac{1}{n-1} \left[\sum_{i=1}^n (x_i y_i) - \bar{y} \sum_{i=1}^n (x_i) - \bar{x} \sum_{i=1}^n (y_i) + \bar{x} \bar{y} \sum_{i=1}^n (1) \right]$$

$$\text{cov}(x, y) = \frac{1}{n-1} \left[\sum_{i=1}^n (x_i y_i) - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y} \right]$$

$$\text{cov}(x, y) = \frac{1}{n-1} \left[\sum_{i=1}^n (x_i y_i) - \frac{1}{n} \sum_{i=1}^n (x_i) \sum_{i=1}^n (y_i) \right]$$

<http://campus.fsu.edu/> Data Exploration: Paired Data Sets 10
 bourassa@met.fsu.edu COAPS The Florida State University

Computationally Efficient Correlation Step 2: The Standard Deviation

- Consider the standard deviation:

$$s_x = \left[\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \right]^{1/2}$$

$$s_x = \left[\frac{\sum_{i=1}^n (x_i^2) - n \bar{x}^2}{n-1} \right]^{1/2}$$

$$s_x = \left[\frac{\sum_{i=1}^n (x_i^2) - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2}{n-1} \right]^{1/2}$$

<http://campus.fsu.edu/> Data Exploration: Paired Data Sets 11
 bourassa@met.fsu.edu COAPS The Florida State University

Computationally Efficient Correlation Step 3: The Correlation

- Combine covariance and standard deviation to get correlation:

$$r_{xy} = \frac{\text{covariance}(x, y)}{s_x s_y}$$

$$r_{xy} = \frac{\frac{1}{n-1} \left[\sum_{i=1}^n (x_i y_i) - \frac{1}{n} \sum_{i=1}^n (x_i) \sum_{i=1}^n (y_i) \right]}{\left[\frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right] \right]^{1/2} \left[\frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2 \right] \right]^{1/2}}$$
- All the terms in the above equation can be calculated in one pass of a data set.
- For example: reading satellite data can be extremely time intensive, and often the data are too massive to store.
- Calculating standard deviations or correlations in one pass is great.

<http://campus.fsu.edu/> Data Exploration: Paired Data Sets 12
 bourassa@met.fsu.edu COAPS The Florida State University

FORTRAN Tidbit: Procedures AKA Subprograms

- There are two types of procedures: functions and subroutines.
- In general, procedures have zero or more arguments
 - E.g., subprogram1(x1, x2, x3, x4)
 - The variables x1, x2, x3, and x4 are arguments
 - Each procedure must end with the END command
 - Each procedure will cause the program to stop if the program reaches any END command
 - If procedure is not suppose to cause the program to stop, then the program must reach a RETURN command prior to the END.
- Subroutines change the value of one or more of the arguments.
 - Executed with a CALL command. E.g., CALL MEAN(x, ave)
- Functions do not alter any arguments, but return a value.
 - Example: y = mean(x)
- There are several ways to declare procedures.
 - Procedures must be declared in any program that uses them.

http://campus.fsu.edu/bourassa@met.fsu.edu  The Florida State University  Data Exploration: Paired Data Sets 13

FORTRAN90 Example Function

```

• Consider a subroutine to calculate standard deviation.
FUNCTION STANDEV( x, n )
! n the number of values in array x
! x array of values for which the standard deviation will be determined
INTEGER :: i_data, n
REAL :: standev, sum_x, sum_x_sqd
REAL, dimension( n ) :: x
sum_x = 0.0
sum_x_sqd = 0.0
DO i_data = 1, n
    sum_x = sum_x + x(i_data)
    sum_x_sqd = sum_x_sqd + x(i_data)**2
ENDDO
standev = SQRT( ( sum_x_sqd - ( sum_x ** 2 ) / REAL(n) ) / REAL(n-1) )
RETURN
END FUNCTION STANDEV
    
```

http://campus.fsu.edu/bourassa@met.fsu.edu  The Florida State University  Data Exploration: Paired Data Sets 14

Alternative Correlation Methods

- Spearman Rank Correlation
 - More robust than Pearson correlation
 - Computes a Pearson correlation using the ranks of the data, rather than the actual data values.
 - Reduces the influence of outliers
 - Still hampered by noisy data
- Can be simplified to

$$r_{\text{rank}} = 1 - \frac{6 \sum_{i=1}^n D_i^2}{n(n^2 - 1)}$$
 - Where D_i is the difference in ranks between the i^{th} pair of ranks.

http://campus.fsu.edu/bourassa@met.fsu.edu  The Florida State University  Data Exploration: Paired Data Sets 15

Autocorrelation

- Autocorrelation is used to investigate how information at one time is related to information at other time.
- It is useful for examinations of:
 - Persistence
 - Repeating cycles
- Autocorrelation is a correlation of a data set with itself, but with one of the series lagged
- For example: With a series of daily temperatures $\{T_0, T_1, T_2, \dots, T_{31}\}$ could be correlated with a series lagged by one day $\{T_{-1}, T_0, T_1, \dots, T_{30}\}$
 - Correlations for days of two or more days could also be calculated.
- The autocorrelation for the k^{th} lag could be written as

$$r_k = \frac{\sum_{i=1}^{n-k} [(x_i - \bar{x}_0)(x_{i+k} - \bar{x}_k)]}{\left[\sum_{i=1}^{n-k} (x_i - \bar{x}_0)^2 \sum_{i=k+1}^n (x_i - \bar{x}_k)^2 \right]^{0.5}}$$

Where the k^{th} mean is based on values from x_{i+k} to x_n

http://campus.fsu.edu/bourassa@met.fsu.edu  The Florida State University  Data Exploration: Paired Data Sets 16

Alternative Autocorrelation

- The previous version of autocorrelation is useful when the overlapping portions of the data set are too small.
 - However, it is odd to work with means from different periods
- If there is a large number of overlapping points, then an alternative version can be applied.
 - Examines only the overlapping period.
 - The means are identical
- For example: With a series of daily temperatures $\{T_0, T_1, T_2, \dots, T_{31}\}$ could be correlated with a series lagged by one day using the values $\{T_0, T_1, T_2, \dots, T_{30}\}$
- The autocorrelation for the k^{th} lag could be written as

$$r_k = \frac{\sum_{i=1}^{n-k} [(x_i - \bar{x})(x_{i+k} - \bar{x})]}{\sum_{i=1}^{n-k} (x_i - \bar{x})^2}$$

Where all means are based on values from x_1 to x_n

http://campus.fsu.edu/bourassa@met.fsu.edu  The Florida State University  Data Exploration: Paired Data Sets 17

Example Autocorrelation

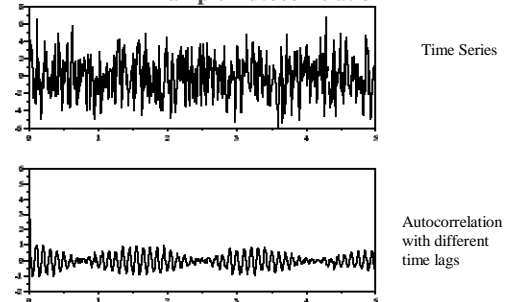


Figure from www.neurotraces.com/scilab/scilab2/node39.html
 http://campus.fsu.edu/bourassa@met.fsu.edu  The Florida State University  Data Exploration: Paired Data Sets 18

