



MEET3220C & MEET6480
Computational Meteorology



Exploratory Data Analysis
Empirical Distributions

Descriptive Statistics
Robustness of Statistics

<http://cam.psu.edu/bouassa@met.fsu.edu>



The Florida State University



Computational Statistics
Introduction 1

Exploratory Data Analysis

- Exploratory data analysis employs statistics to summarize characteristics of the data.
 - Converts data to information.
 - Graphical representations are often used to examine the data, and to infer additional qualities of the data set.
- We will discuss many graphical statistical/graphical techniques for examining data.
- One key question is 'how robust are the statistics?'
- If the summary statistics change greatly depending on the subset of the data that is being examined (assuming the subset is sufficiently large), then we should place little value in the statistics!
 - We will discuss which techniques are more reliable than others.
 - Sometimes the consequences of the data not meeting our assumptions (e.g., a bell curve) are quite serious.
- Key Point: ALWAYS LOOK AT THE DATA !!!!

<http://cam.psu.edu/bouassa@met.fsu.edu>



The Florida State University



Computational Statistics
Introduction 2

3M : Mean, Mode, and Median

- The 3 M's each give a measure of typical values.
 - The mean is the average. The notation for the mean of x is \bar{x} .
 - Mode is the most frequently occurring number.
 - If we order a data set, then the median is the value in the middle of the list.
- Consider the grades on a hypothetical homework assignment:
 - 21 values
 - 4, 5, 6, 6, 7, 7, 7, 8, 8, 8, 8, 9, 9, 9, 9, 10, 10, 10, 10, 10, 10
 - The mean is 8.1 (A-) $\left(\sum_{i=1}^n \text{grade}_i \right) / N$
 - The mode is 10 (A)
 - The median is 8 (A-) The 11th value, $\text{grade}_{(N+1)/2}$

<http://cam.psu.edu/bouassa@met.fsu.edu>



The Florida State University



Computational Statistics
Introduction 3

Quartiles and Percentiles

- Quartiles and percentiles are used to describe the spread or distribution of a data set.
- Consider an ordered set of data $\{x_1, x_2, x_3, \dots, x_{n-1}, x_n\}$
- The median is defined as
 - $Q_{0.5} = \begin{cases} x_{(n+1)/2} & \text{if } n \text{ is odd} \\ (x_{n/2} + x_{n/2+1})/2 & \text{if } n \text{ is even} \end{cases}$
 - Halfway through the series
- Quartiles are 25% through the series for the lower quartile $Q_{0.25}$, and 75% through the series for the upper quartile $Q_{0.75}$.
- Percentiles are the value that are greater than a percentage of the data set.
 - Example: the 50th percentile is the median.
 - In data set with 100 values, the 99th percentile is the greatest value.

<http://cam.psu.edu/bouassa@met.fsu.edu>



The Florida State University



Computational Statistics
Introduction 4

FORTRAN : Arrays

- We will discuss a new type of variable designed to hold a series of data.


```
REAL fake_obs(125) !A variable that can hold 125 real values.
DO index = 1, 125
    fake_obs(index) = (REAL(index) + 0.5) ** 2
! Converts index to a real number, then adds 0.5, the squares to total
ENDDO
PRINT *, fake_obs(10:20) ! Prints the 10th to 20th elements of the array
```
- Consider a sorted array $\{x_1, x_2, x_3, \dots, x_{n-1}, x_n\}$.
- If the change with index number is uniform, and n is large, then the mean can be approximated as $x_{(n+1)/2}$
 - This value is the median, even if the change with index number is non-uniform
- And the 30th percentile can be approximated as $x(0.3 * n + 1)$

<http://cam.psu.edu/bouassa@met.fsu.edu>



The Florida State University



Computational Statistics
Introduction 5

More Robust Estimates of Central Location

- The mean is sensitive to outliers
- Consider the data set $\{11, 12, 13, 14, 15, 16, 17, 18, 19\}$
 - The mean and median are 15.
- However, if the final value (19) is replaced with 91, then the mean becomes 23.
 - The robustness can be examined by examining the differences between the mean of the whole data set and the mean with a small fraction of the data removed.
 - This test should be done for many samples.
 - This example is a form of cross validation.
- A trim mean is a more robust measure of the central location.
 - Trim mean = $(Q_{0.25} + 2Q_{0.5} + Q_{0.75}) / 4$
- A mean could also be determined from a trimmed portion of the data set:

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

<http://cam.psu.edu/bouassa@met.fsu.edu>



The Florida State University



Computational Statistics
Introduction 6

Example Cross Validation Code
 Useful for testing the robustness of a mean of a small number of independent obs.

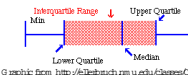
```

REAL sum
REAL, dimension(365) :: daily_rain !array of 365 daily temperatures
REAL test_means(365) !An alternative version of similar declaration
INTEGER index, n, skip
N = 365 !Assume that the values are read by the program
sum = 0.0
DO index = 1, n
    sum = sum + daily_rain(index)
ENDDO
DO index = 1, n
    test_means(index) = sum - daily_rain(index)
    test_means(index) = test_means(index) / REAL(n - 1)
ENDDO
sum = sum / n
  
```

<http://cm.psu.edu/boussas@mctfau.edu>  The Florida State University  Computational Statistics Introduction 7

Spread: Interquartile Range (IQR)

- One measure of spread is the interquartile range.
- Spread is an indication of departure from the mean.
- $IQR = Q_{0.75} - Q_{0.25}$
- The IQR is a very robust measure of the spread of values near the mean, but does not give any information on outliers.



Graphic from http://mathscholar.zsu.edu/Emanuel/CS560w/94/Students/Wiki/IQR/IQR_Keypoint/index.html

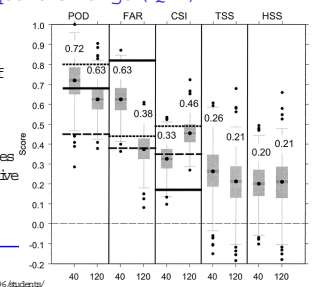


Figure from <http://cm.psu.edu/boussas@mctfau.edu> / requests-display-figures/run/e-8520-0434-18-6-953-04

<http://cm.psu.edu/boussas@mctfau.edu>  The Florida State University  Computational Statistics Introduction 8

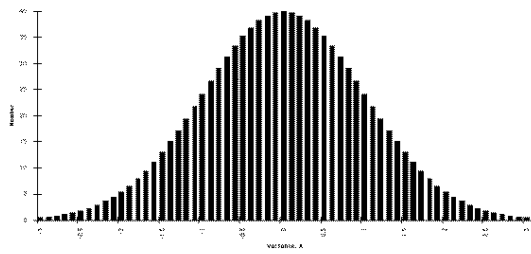
Standard Deviation

- The standard deviation is the most common measure of spread.
 - Unlike the IQR, it does consider outliers.
- The standard deviation is defined as

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$
 - Where s is an estimate of the standard deviation
 - Estimate because it is assumed (?) to be based on an incomplete sample of the population.
 - The true standard deviation is usually notated as σ
- The standard deviation is highly sensitive to outliers. Why?
- Because of the square of the difference from the mean.
- If the data has a Gaussian distribution, then
 - 68% of the data are within 1 standard deviation from the mean
 - 99% of the data are within 3 standard deviations from the mean

<http://cm.psu.edu/boussas@mctfau.edu>  The Florida State University  Computational Statistics Introduction 9

Gaussian Distribution



Graphic from www.cmu.edu/epfl/~dave/ai/1/normal_distribution.html

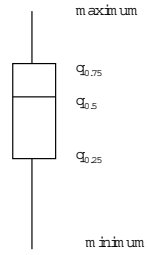
<http://cm.psu.edu/boussas@mctfau.edu>  The Florida State University  Computational Statistics Introduction 10

Median Absolute Deviation (MAD)

- The MAD does consider outliers, but unlike the standard deviation, the outliers have similar influence to the non-outliers.
- $MAD = \text{median}(|x_i - Q_{0.5}|)$
- Why is the influence of outliers reduced?
- Two reasons:
 - No square of the difference from the central location
 - The median (rather than the mean) is not influenced by outliers.

<http://cm.psu.edu/boussas@mctfau.edu>  The Florida State University  Computational Statistics Introduction 11

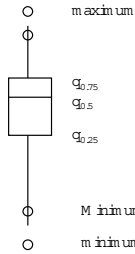
Box and Whiskers Plots



- The box plot is modified to show the extreme values of the data set.
- Alternatively, the whiskers can indicate the 5th and 95th percentiles
- Alternatively, the whiskers can indicate multiples of the IQR.

<http://cm.psu.edu/boussas@mctfau.edu>  The Florida State University  Computational Statistics Introduction 12

Box and Whiskers Variant



- The box plot is modified to show the additional measure of spread, and extreme values of the data set.
- Also shown are the most extreme values within the measure of spread.

<http://cam.psu.edu/boussess@met.fsu.edu>



Computational Statistics Introduction 13

Symmetry – or a lack thereof

- “Is the data symmetric?” is a key question for many assumptions
 - Example: are the ENSO impacts for ELNño equal and opposite those for La Niña? If so, forecasting is a lot easier!
- Skewness is one measure of asymmetry

$$\gamma = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$$

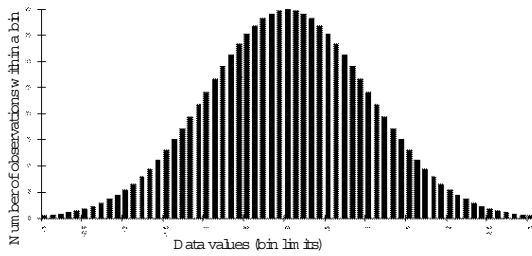
- Skewness is far from a robust statistic
 - There are several alternatives.
 - However, these are not commonly used in statistical analyses

<http://cam.psu.edu/boussess@met.fsu.edu>



Computational Statistics Introduction 14

Histograms, or Probability Distribution Functions



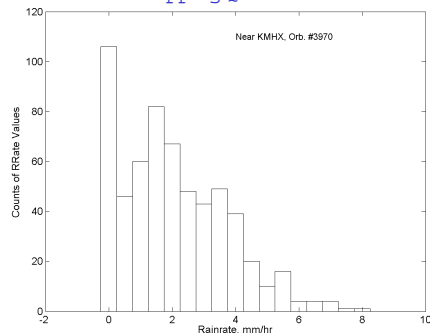
Graph from www.ctm.nyu.edu/ctmweb/CTM/ctmfrom.shtml

<http://cam.psu.edu/boussess@met.fsu.edu>



Computational Statistics Introduction 15

Hurr. Isabel Rainrate Value Distribution for Overlapping QSCAT & NEXRAD

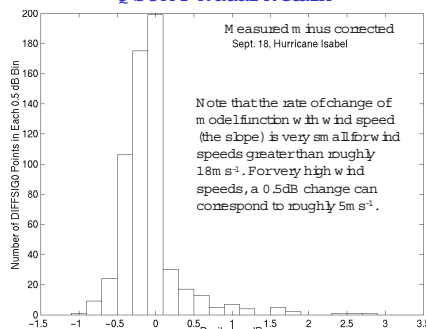


<http://cam.psu.edu/boussess@met.fsu.edu>



Computational Statistics Introduction 16

Physically Based Correction to QSCAT Radar Return

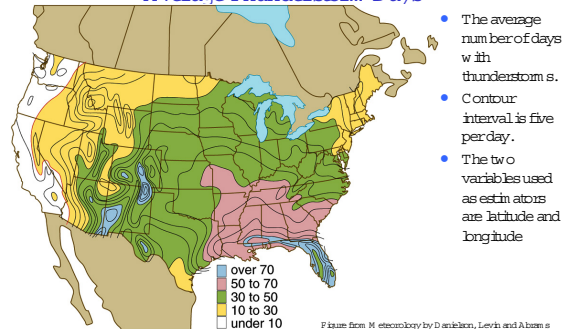


<http://cam.psu.edu/boussess@met.fsu.edu>



Computational Statistics Introduction 17

Examine of Histogram Bivariate Estimator Average Thunderstorm Days



<http://cam.psu.edu/boussess@met.fsu.edu>

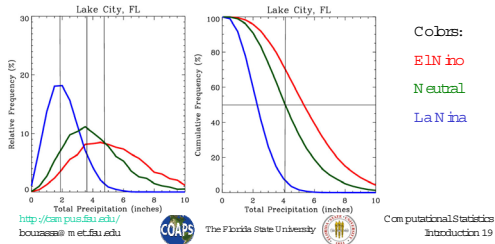


Computational Statistics Introduction 18

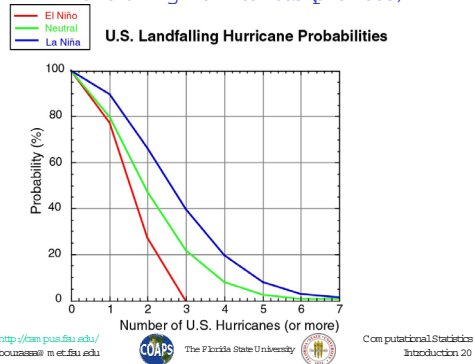
- The average number of days with thunderstorms.
- Contour interval is five per day.
- The two variables used as estimators are latitude and longitude

Cumulative Probability Distributions

- Cumulative probability distributions plot
 - X-axis: magnitude of events
 - Y-axis: cumulative probability of all events to the left of the point on the x-axis. $P(x \leq X)$
 - Probability of exceedence is $1 - \text{cumulative probability}$. $P(x \geq X)$



Probability of Exceedence Example: Landfalling Hurricanes (pre 2005)



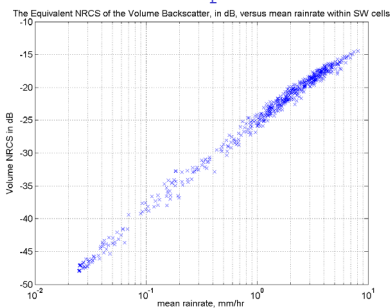
Extreme Value Distributions

- There are several types of extreme value distributions that can be used to describe the likelihood of an event (or an event of lesser magnitude).
- These methods are usually used on ordered (also called ranked) data, where i represents the rank (from lesser values to greater values).
- Most of these follow form given in the textbook.
- They are used to estimate the likelihood of events of certain magnitude, often for engineering or insurance purposes.
- Alternatively, they can be used to estimate that average time (with large margins of error) between events of given magnitudes.
 - Example: the average time between floods of a certain level.

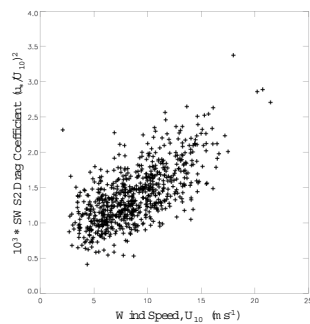
Standardized Anomalies

- We are often interested in departures from a mean value.
- If x_i is a value in a series, and $\langle x \rangle$ is the mean value, the departure associated with x_i is usually written as x'
 - $x_i = \langle x \rangle + x'_i$
- The standardized anomaly (z) is defined as
 - $z_i = (x_i - \langle x \rangle) / s_x$
 - Where s_x is the sample standard deviation

Exploration of Paired Data Scatterplots

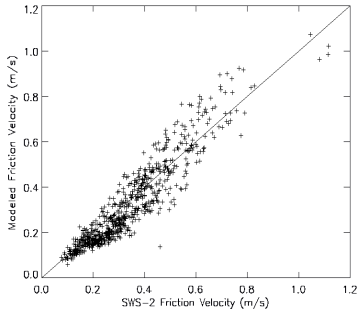


Example Drag Coefficients From Severe Wind Storms 2 Experiment



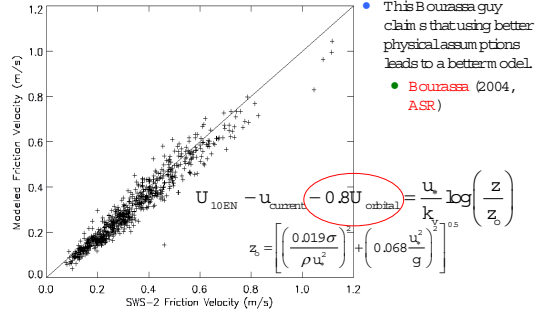
- Preliminary version of the data set provided by Peter K. Taylor.
- These drag coefficients are based on high quality observations.
- Observations that are mostly from rough seas.

Results of Taylor and Yelland's Parameterization on SW S2 data



<http://cam.psu.edu/boussas@m.tfsu.edu>
 COAPS The Florida State University
 Computational Statistics Introduction 25

Bourassa (2004) Comparison to Observations



<http://cam.psu.edu/boussas@m.tfsu.edu>
 COAPS The Florida State University
 Computational Statistics Introduction 26

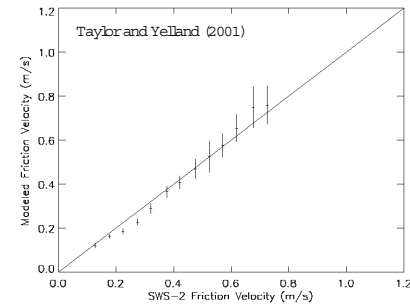
Uncertainty in a Mean

- A problem with using overall measurements of error is that many are very sensitive to outliers.
 - Comparison statistics are largely dependent on a small fraction of the data set.
 - The results are similar even for very different models.
- It is more useful to examine the statistics for small sub-samples of the data set.
 - This can also be misleading for some cases, which will be addressed in later lectures.
 - One useful diagnostic statistic is the uncertainty in the mean.
 - If the errors have a Gaussian distribution, then the uncertainty in the mean is

$$s_x = s_x / \sqrt{n}$$
 - Where n refers to the number of independent points in the sample (or sub-sample)

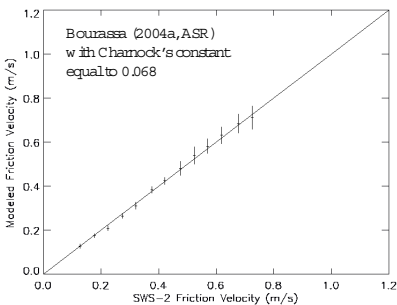
<http://cam.psu.edu/boussas@m.tfsu.edu>
 COAPS The Florida State University
 Computational Statistics Introduction 27

More Results: Means & Three Standard Deviations



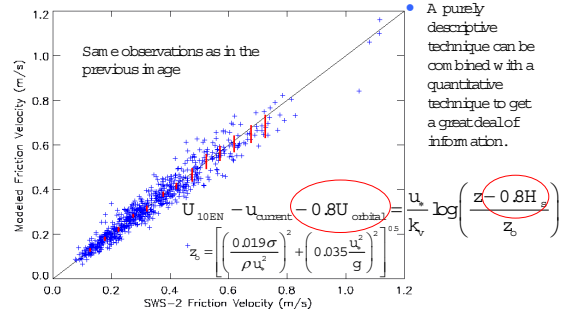
<http://cam.psu.edu/boussas@m.tfsu.edu>
 COAPS The Florida State University
 Computational Statistics Introduction 28

More Results: Means & Three Standard Deviations



<http://cam.psu.edu/boussas@m.tfsu.edu>
 COAPS The Florida State University
 Computational Statistics Introduction 29

Results of Bourassa (2005) Compared to SW S2 Observations



<http://cam.psu.edu/boussas@m.tfsu.edu>
 COAPS The Florida State University
 Computational Statistics Introduction 30