

# **A Biomed Data Analyst Training Program**

---

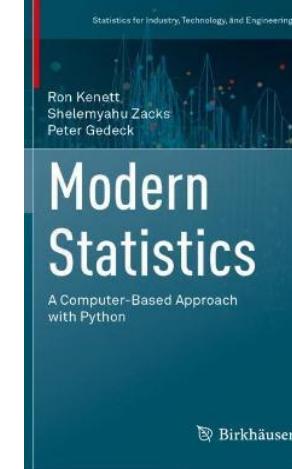
## **Data visualization**

---

### **Professor Ron S. Kenett**

# Chapter 1

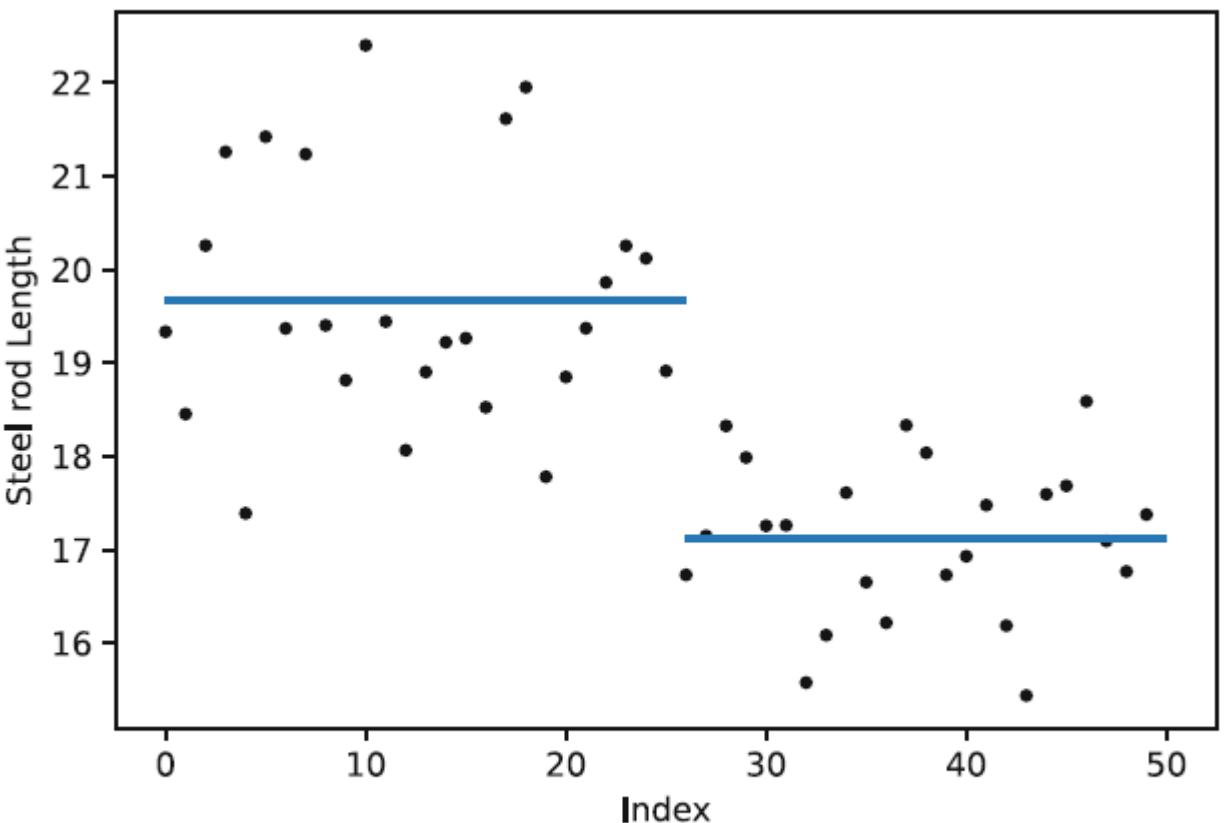
## Analyzing Variability: Descriptive Statistics



**Preview** The chapter focuses on statistical variability and various methods of analyzing random data. Random results of experiments are illustrated with distinction between deterministic and random components of variability. The difference between accuracy and precision is explained. Frequency distributions are defined to represent random phenomena. Various characteristics of location and dispersion of frequency distributions are defined. The elements of exploratory data analysis are presented.

```
steelrod[26:] = steelrod[26:] - 3
```

```
ax = steelrod.plot(y='STEELROD', style='.', color='black')
ax.set_xlabel('Index')
ax.set_ylabel('Steel rod Length')
ax.hlines(y=steelrod[:26].mean(), xmin=0, xmax=26)
ax.hlines(y=steelrod[26:].mean(), xmin=26, xmax=len(steelrod))
plt.show()
```



**Fig. 1.2** Level shift after the first 25 observations

```
from scipy.stats import beta, norm

x = np.linspace(-3, 3, 200)
df = pd.DataFrame({'x': x,
                    'steep': beta(8, 8, loc=-3, scale=6).pdf(x),
                    'flat': beta(2.5, 2.5, loc=-3, scale=6).pdf(x),
                    'normal': norm().pdf(x),
                    })
ax = df.plot.line(x='x', y='steep', legend=False, color='black')
df.plot.line(x='x', y='normal', legend=False, color='black',
             linestyle='--', ax=ax)
df.plot.line(x='x', y='flat', legend=False, color='black',
             linestyle='-.', ax=ax)
ax.set_ylabel('y')
ax.text(0.5, 0.5, 'Steep')
ax.text(1.0, 0.35, 'Normal')
ax.text(2.0, 0.2, 'Flat')
plt.show()
```

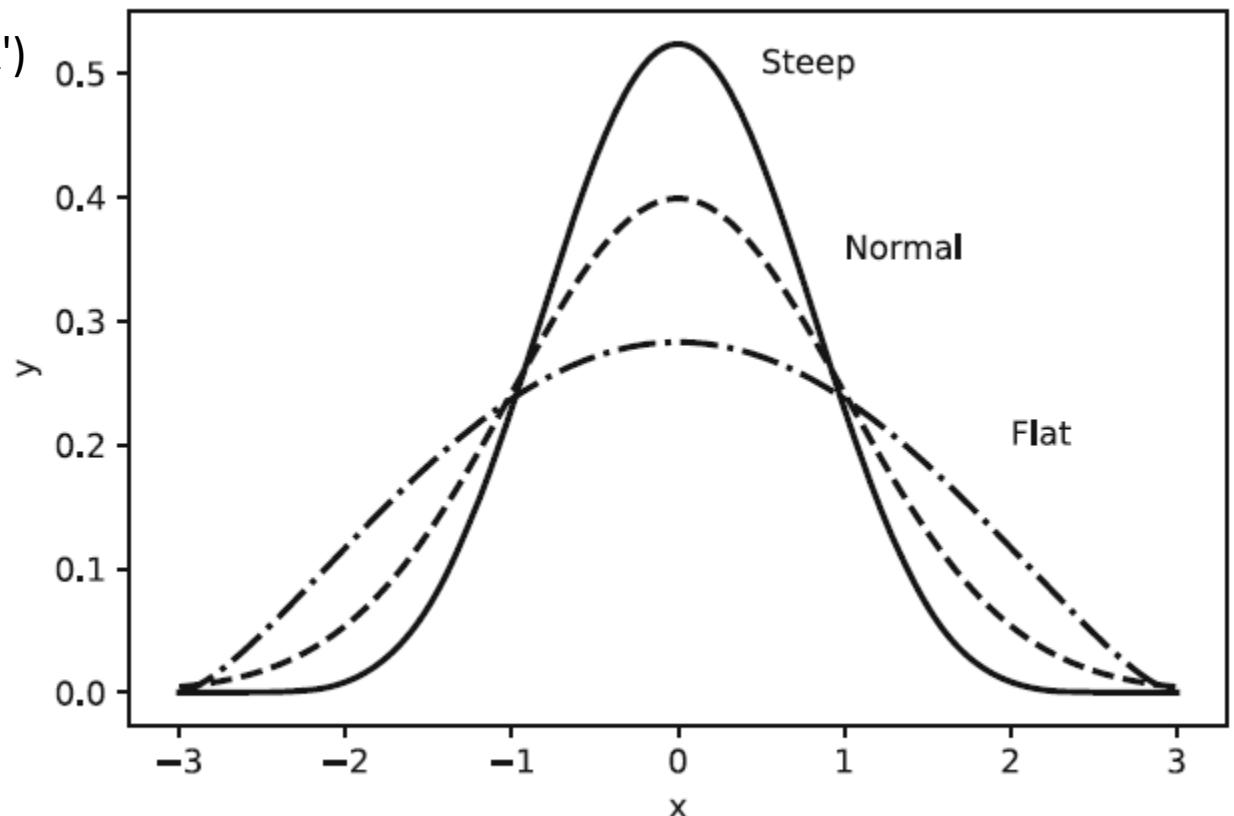
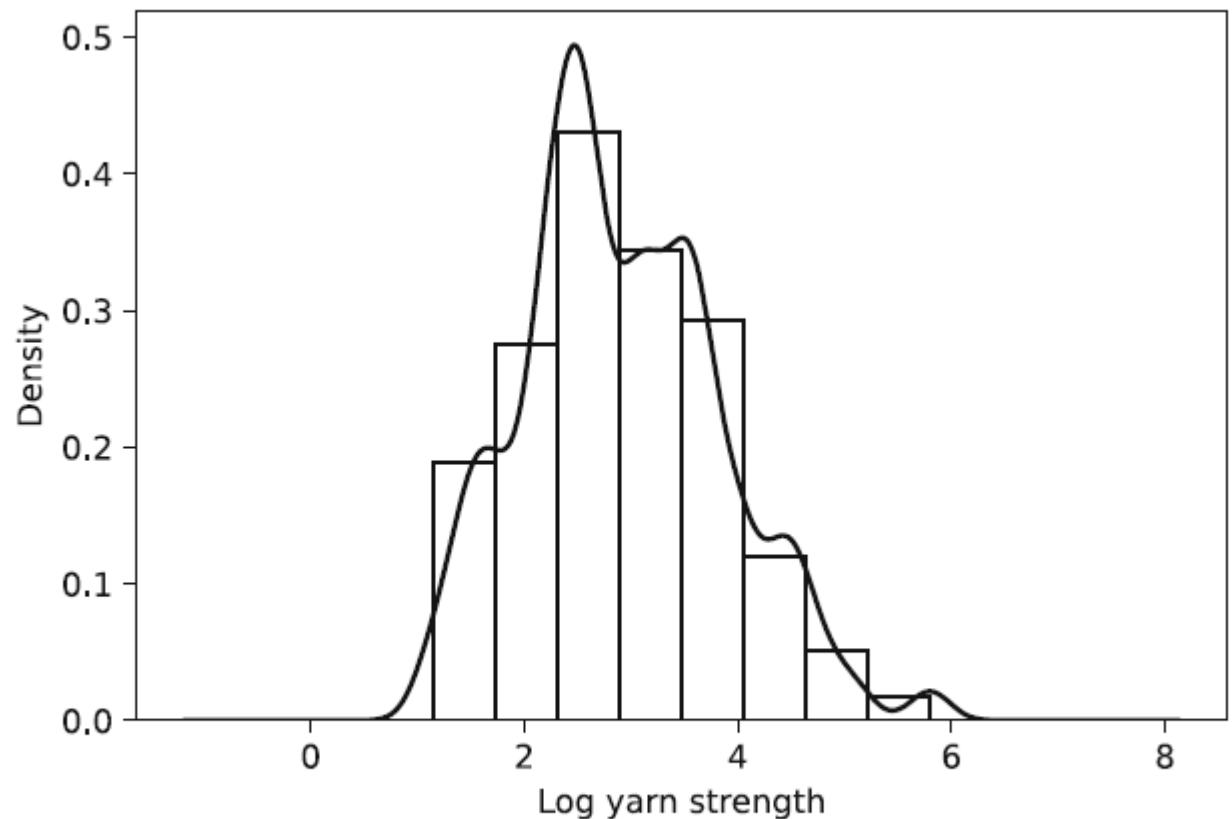


Fig. 1.13 Normal, steep, and flat distributions

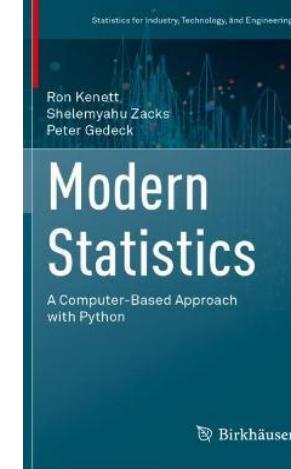
```
X = mistat.load_data('YARNSTRG')
ax = X.plot.hist(bins=8, color='white', edgecolor='black',
legend=False, density=True)
X.plot.density(bw_method=0.2, ax=ax, color='black')
ax.set_xlabel('Log yarn strength')
plt.show()
```



**Fig. 1.14** Comparison of histogram and density plot for the log yarn strength datasets

# Chapter 4

## Variability in Several Dimensions and Regression Models



**Preview** When surveys or experiments are performed, measurements are usually taken on several characteristics of the observation elements in the sample. In such cases we have multivariate observations, and the statistical methods which are used to analyze the relationships between the values observed on different variables are called multivariate methods. In this chapter we introduce some of these methods. In particular, we focus attention on graphical methods, linear regression methods, and the analysis of contingency tables. The linear regression methods explore the linear relationship between a variable of interest and a set of variables, by which we try to predict the values of the variable of interest. Contingency tables analysis studies the association between qualitative (categorical) variables, on which we cannot apply the usual regression methods.

```
# The following command would be sufficient to create the  
scatterplot matrix
```

```
# matplotlib has however a problem with scaling xDev  
sns.pairplot(place[['xDev', 'yDev', 'tDev']], markers=".",  
             plot_kws={'facecolors': 'none', 'edgecolor': 'black'},  
             diag_kws={'color': 'grey'})
```

```
#def panelPlot(x, y, **kwargs):  
#    plt.scatter(x, y, **kwargs,  
#                facecolors='none', edgecolor='black', s=20)  
#    dx = 0.05*(max(x) - min(x))  
#    plt.xlim(min(x)-dx, max(x) + dx)  
#    dy = 0.05*(max(y) - min(y))  
#    plt.ylim(min(y)-dy, max(y) + dy)  
#g = sns.PairGrid(place[['xDev', 'yDev', 'tDev']])  
#g.map_offdiag(panelPlot)  
plt.show()
```

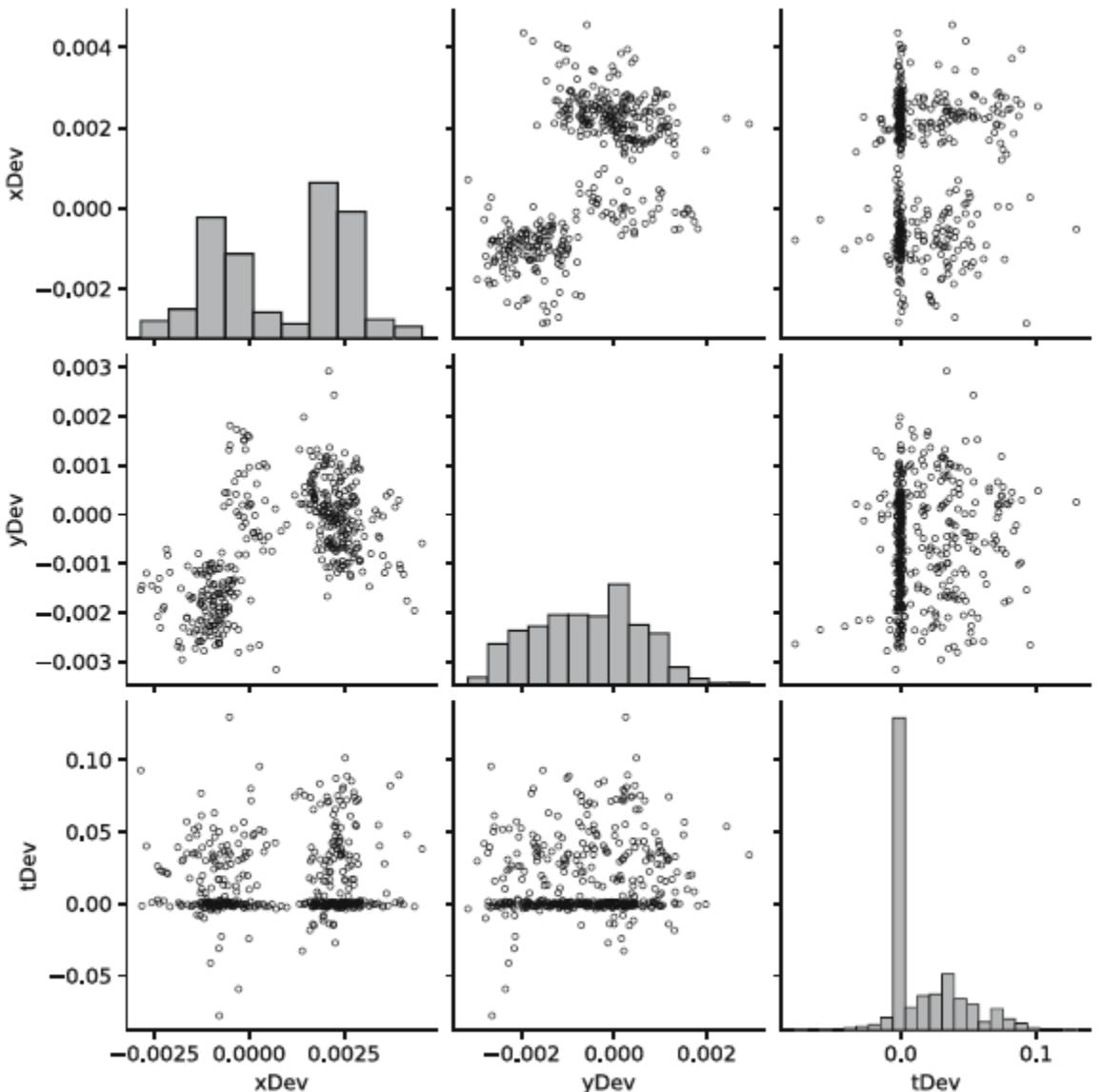


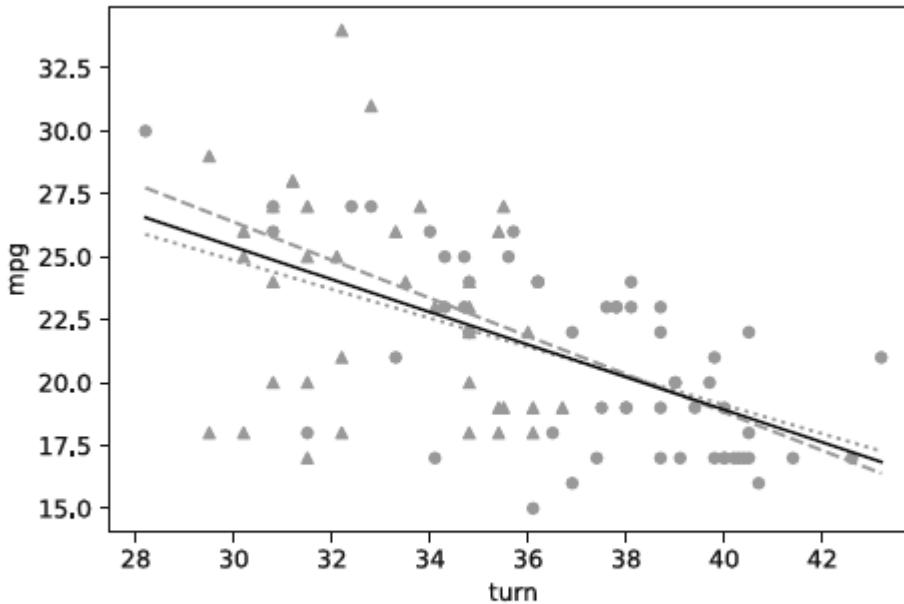
Fig. 4.2 Scatterplot matrix

```

# create visualization
ax = car_US.plot.scatter(x='turn', y='mpg', color='gray',
marker='o')
car_Asia.plot.scatter(x='turn', y='mpg', ax=ax, color='gray',
marker='^')

car_combined = car_combined.sort_values(['turn'])
ax.plot(car_combined['turn'],
model_US.predict(car_combined),
color='gray', linestyle='--')
ax.plot(car_combined['turn'],
model_Asia.predict(car_combined),
color='gray', linestyle=':')
ax.plot(car_combined['turn'],
model_simple.predict(car_combined),
color='black', linestyle='-' )
plt.show()

```



**Fig. 4.16** Linear regression analysis for US (filled circle, dashed line) and Japanese cars (filled triangle, dotted line). The solid line is the linear regression of the combined data set

# Information Quality

*The potential of a particular dataset to achieve a particular goal using a given empirical analysis method*



$$\text{Info}Q(f, X, g) = U(f(X|g))$$

*g* A specific analysis goal

*x* The available dataset

*f* An empirical analysis method

*U* A utility measure

1. Data resolution  
2. Data structure  
3. Data integration  
4. Temporal relevance  
5. Chronology of data and goal  
6. Generalizability  
7. Operationalization  
**8. Communication**

How

What



# Presenting uncertainty in data

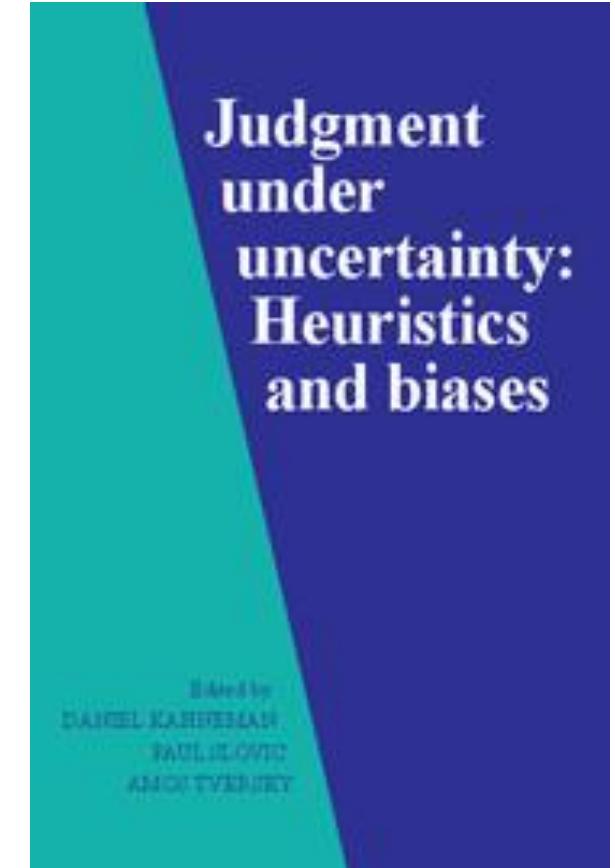
- Capability of human mind for solving complex problems is limited compared with the size of problems
- Lack of objectively rational behaviour in real world. Cognitive illusions.
- Use of simple “rules of thumb” to simplify decision making
- Heuristics can be helpful, but can also lead to biases, especially in uncertain situations where probabilities are encountered

# Presenting uncertainty in data

- “Nothing is certain”
- In many situations, decisions have to be based on probabilities
- Interpretation of probabilities is sometimes not straightforward
- Appropriate presentation can help to make the right decisions

# Presenting uncertainty in data

- formulating the problem:
  - probabilities vs. frequencies
  - the framing effect
  - the anchoring effect
- underweighting base rates
- hindsight and confirmation bias
- belief persistence: Primacy and inertia effect
- group conformity and decision regret



# Conditional probabilities

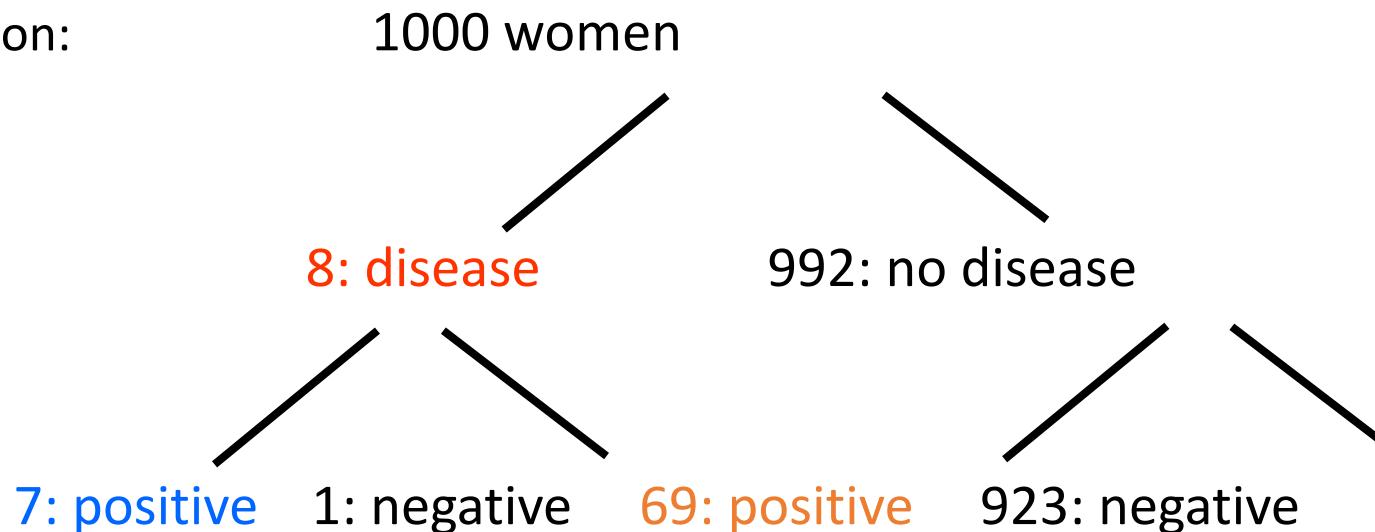
- Breast cancer screening: The facts:
  - Probability that a woman aged 40-50 has breast cancer = 0.8%
  - If a woman has breast cancer, probability of positive test = 90%
  - If a woman does not have breast cancer, prob. of positive test=7%
- Imagine a woman with a positive test.

What is the probability, that she actually has breast cancer?

- Solution:
  - $p(\text{disease}) = 0.008$
  - $p(\text{pos} | \text{disease}) = 0.90$
  - $p(\text{pos} | \text{no disease}) = 0.07$
  - $$p(\text{disease} | \text{pos}) = \frac{p(\text{disease}) * p(\text{pos} | \text{disease})}{p(\text{disease}) * p(\text{pos} | \text{disease}) + p(\text{no disease}) * p(\text{pos} | \text{no disease})}$$
$$= 0.09$$

# Frequency formulation

- Breast cancer screening: The facts:
  - Probability that a woman aged 40-50 has breast cancer = 0.8%
  - If a woman has breast cancer, probability of positive test = 90%
  - If a woman does not have breast cancer, prob. of positive test=7%
- Solution:

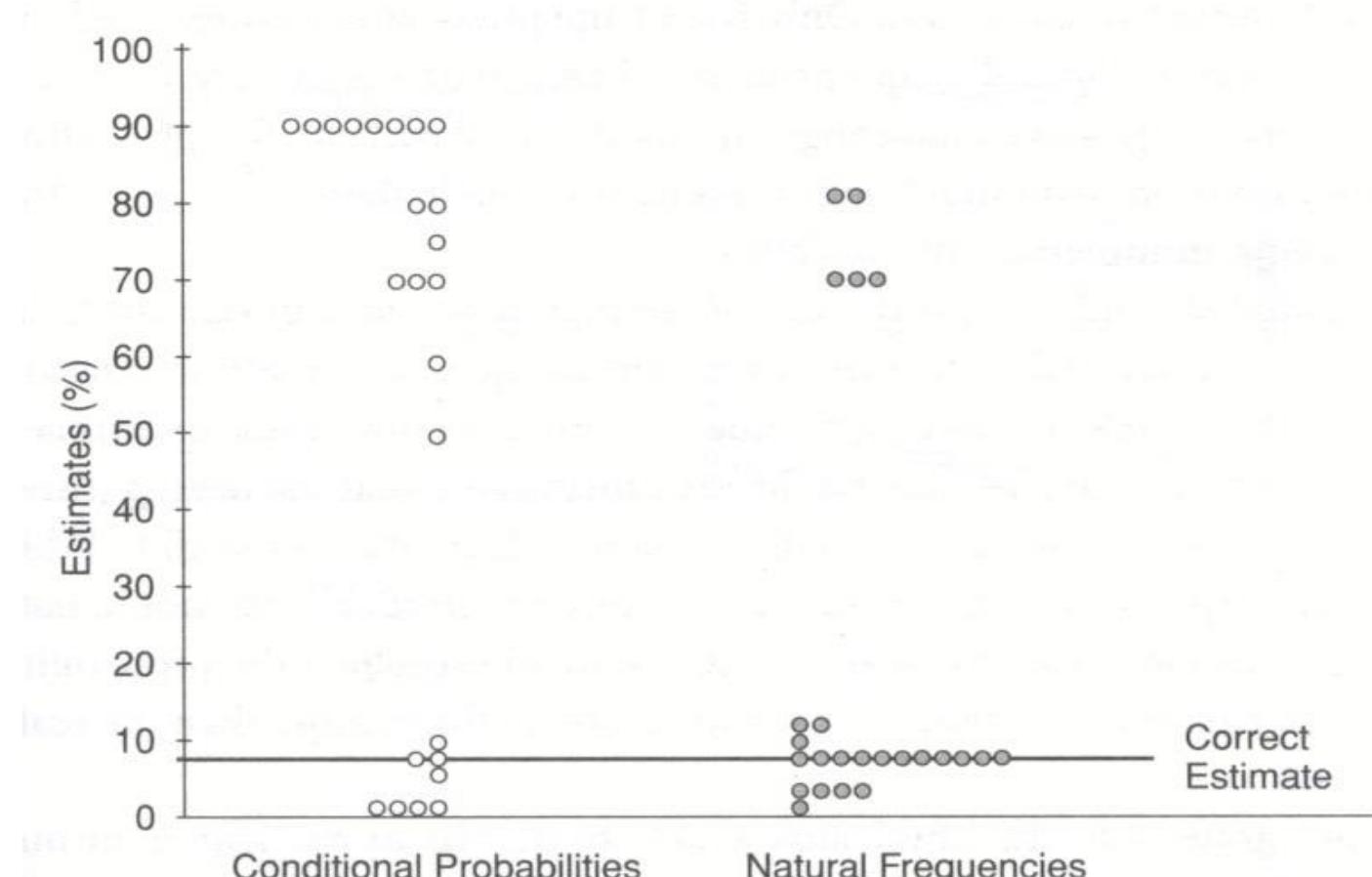


Handwritten calculations on a whiteboard:

$$\begin{aligned} & 5 \times 3 \times 8 = 120 \\ & 240 \times \frac{1}{2} = 120 \\ & 60 \text{ min} \times 8 \text{ hrs/day} = 480 \text{ min/day} \\ & 480 \text{ min/day} \times 240 \text{ days/yr} = 115,200 \\ & 115,200 \times 0.008 = 921.6 \\ & 115,200 \times 0.09 = 10368 \\ & 10368 / 921.6 = 11.3 \approx 13\% \end{aligned}$$

$$p(\text{disease} | \text{pos}) = 7 / (7+69) \approx 0.09$$

# Probabilities vs. frequencies



Estimated chances of breast cancer given a positive screening mammogram (Gigerenzer, 2002)

# The framing effect

- The way a problem (or forecast) is formulated can affect a decision
- Imagine that London faces an unusual disease that is expected to kill 600 people. Two alternative programs to combat disease:
  - Program A: 200 people will be saved
  - Program B: 1/3 probability 600 saved, 2/3 probability nobody saved

Tests indicate that 72% would select program A (risk-averse)

- Slightly changed wording:
  - Program C: 400 people will die
  - Program D: 1/3 prob. that nobody will die, 2/3 prob. that 600 will die

Tests indicate that 78% would select program D (risk-taking)

# The framing effect in real life

- Professionals, experienced in decision-making, are still affected
- E.g., information for doctors:
  - mortality rate of 7% within 5 years -> hesitant to recommend
  - survival rate after 5 years of 93% -> more inclined to recommend
- For weather predictions this suggests different response to forecasts expressed as likelihood of drought or non-liability of wet conditions
- E.g., different response to: 30% chance of drought and
  - 70% chance of normal or wet conditions
- Worded vs. numerical forecast:
  - 11% judge forecast “rain is likely” as poor if it did not rain
  - 37% judge forecast “70% chance of rain” as poor if it did not rain although they associate the word “likely” with probability of 70%

# Test your knowledge of history

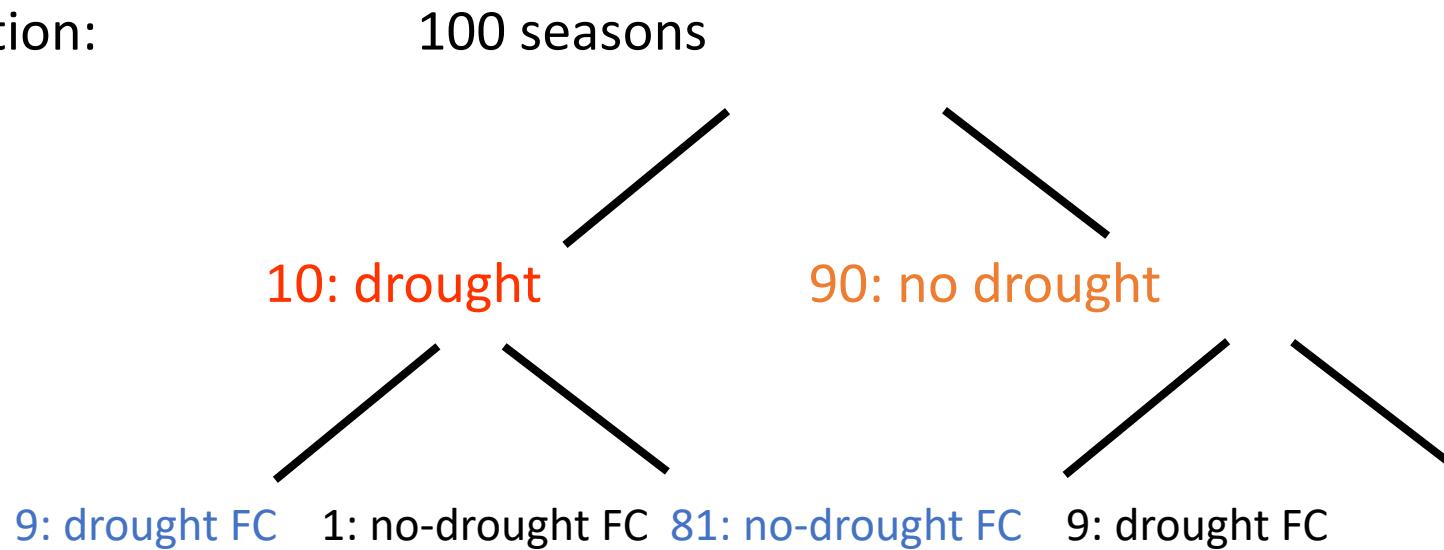
- What are the last three digits of your phone number?
- Add 400 to this number
- Do you think Attila the Hun was defeated in Europe before or after that year?
- In what year would you guess Attila the Hun was defeated?
- The correct answer is: A.D. 451



| Range of initial anchor | Average estimate |
|-------------------------|------------------|
| 400 – 599               | 629              |
| 600 – 799               | 680              |
| 800 – 999               | 789              |
| 1000 – 1199             | 885              |
| 1200 – 1399             | 988              |

# Underweighting base rates

- Imagine a climate model (with 90% accuracy) predicts drought
- Historically, there is 10% chance of drought
- What is the chance that drought will occur in next season?
- Solution:



$$p(\text{drought} \mid \text{drought FC}) = 9 / (9+9) = 0.50$$

# Underweighting base rates

- Imagine a climate model (with 90% accuracy) predicts drought
- Historically, there is 10% chance of drought
- What is the chance that drought will occur in next season?

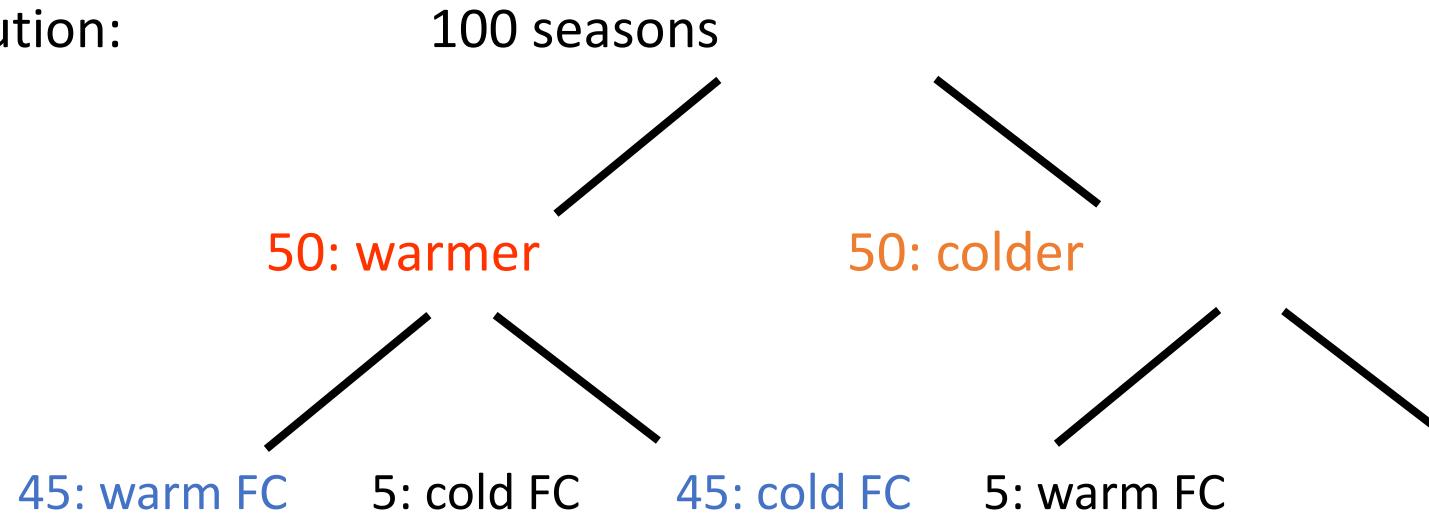
Challenge to convince user that

- Model was correct 90% of time
- the probability of a drought next season was only 50%

*Remember: only for equally likely events,  
accuracy translates into probabilities*

# Underweighting base rates

- Imagine a climate model (with 90% accuracy) predicts warmer than normal conditions
- There is a 50% chance of above normal
- What is the chance that warmer than normal conditions will occur?
- Solution:



$$p(\text{warmer} \mid \text{warm FC}) = 45 / (45+5) = 0.90$$

# Hindsight and confirmation bias

Men mark where they hit, and not where they miss. (Jevons, 1958)

- After finding out whether or not an event occurred, individuals tend to overestimate the degree to which they would have predicted the correct outcome
- Reported outcomes seem less surprising in hindsight than in foresight
- Example: El Nino 1997 regarded as “stunning success”, although only one model was reported in the March 1997 NOAA Long-Lead Forecast Bulletin predicting more than slight warming. Some of the very poor forecasts simply ignored in hindsight.
- Considerable evidence that people tend to ignore (and not search for) disconfirming information of any hypothesis
- Introduce “double-blind test” for model assessment?

# Belief persistence

- Primacy and inertia also tend to weight evidence inaccurately.
- People tend to weight more heavily evidence presented first, e.g. persons described as:
  - intelligent, industrious, impulsive, critical, stubborn, envious
  - are more favourable perceived than persons described as
  - envious, stubborn, critical, impulsive, industrious, intelligent
- Inertia may lead people to ignore evidence that contradicts their prior belief (e.g. that a particular forecast system produces useful forecasts)
- Forecast producers may not recognise the disparity of model predictions, and instead rely too heavily on a forecast that supports their intuitive understanding of the current state of climate

# Probabilities in gambling

## **Monty Hall: Let's Make a Deal**

- in one of the boxes is a bottle of wine
- choose 1, 2, or 3

1

2

3

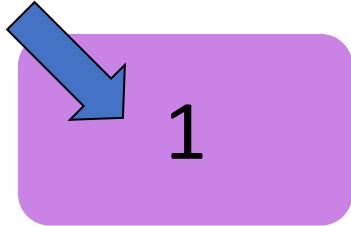
- after choosing, one of the empty boxes will be opened, so that only one empty and one full box are left
- you can choose again (stay with first choice or switch)
- what is the best strategy?

# Probabilities in gambling

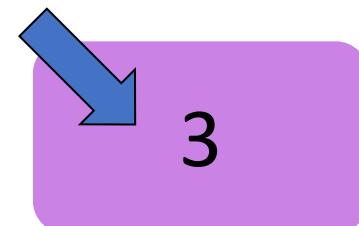
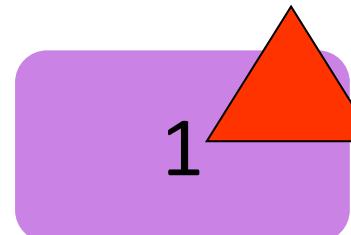
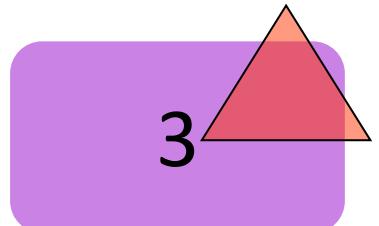
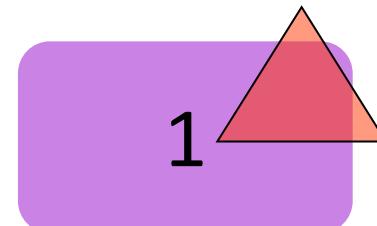
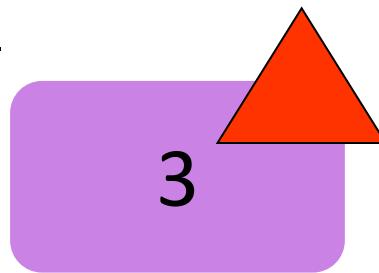
## Monty Hall: Let's Make a Deal

- in one of the boxes is a bottle of wine
- choose 1, 2, or 3

stay



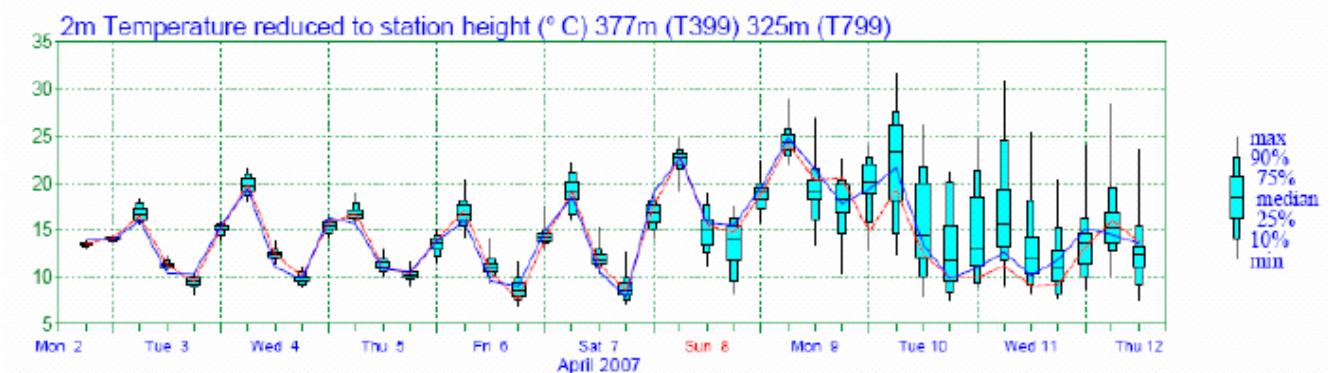
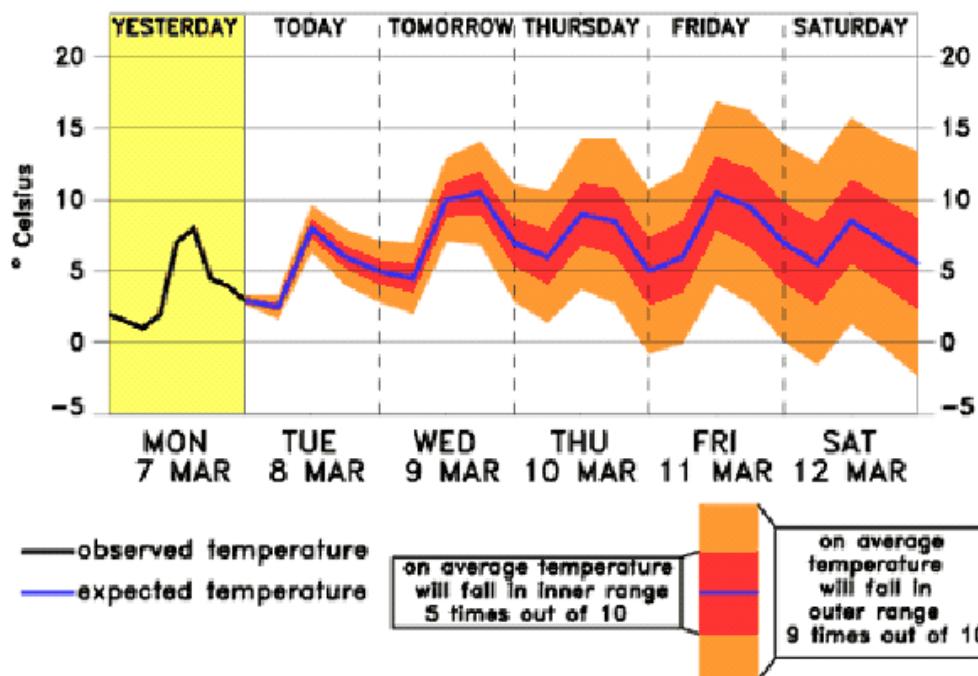
switch



# Strategies to reduce cognitive illusions

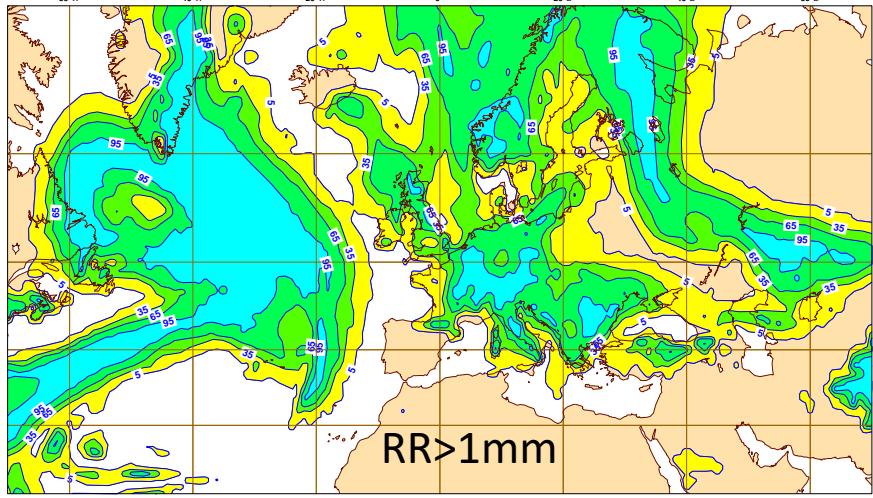
- Recognition that decision-making is inherently biased
- Understanding how written forecasts, and numerical probability forecasts are interpreted by potential users
- Try to reduce impact of cognitive illusions by
  - encouraging forecaster groups to de-bias forecasts by e.g. reducing overconfidence or hindsight bias
  - taking care that media reports and forecasts do not cause anchoring to extreme events (e.g. El Nino 82/83)
  - taking care in wording forecasts to avoid framing
  - avoid “intuitive” approach when combining forecasts, objective approaches exist and are more successful
  - ensuring that base-rates are not ignored
  - using additional visual aids to convey real levels of skill

# Visualization of time series



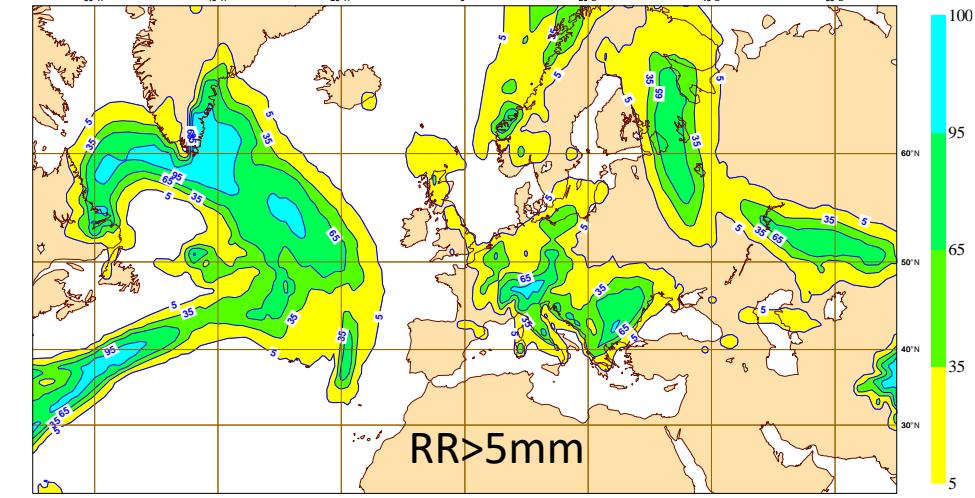
# Probability Maps

Sunday 13 April 2008 00UTC ©ECMWF Forecast probability t+036-060 VT: Monday 14 April 2008 12UTC - Tuesday 15 April 2008 12UTC  
Surface: Total precipitation of at least 1 mm

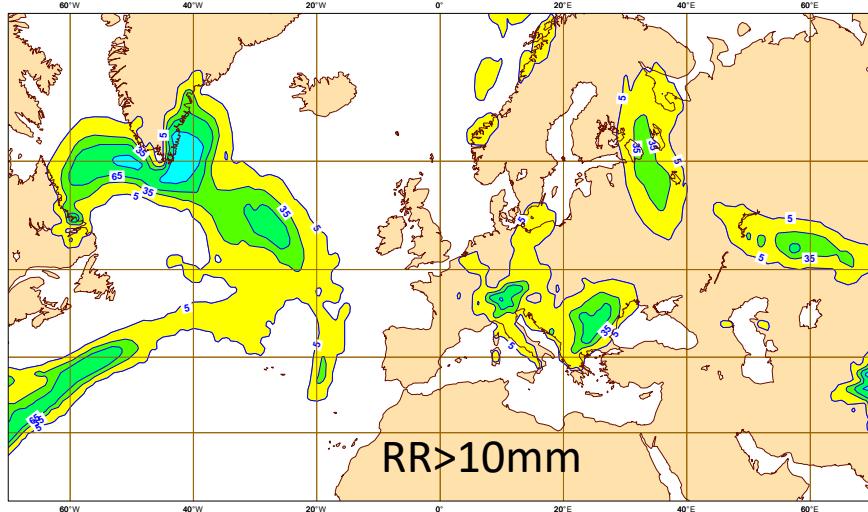


Sunday 13 April 2008 00UTC ©ECMWF Forecast probability t+036-060 VT: Monday 14 April 2008 12UTC - Tuesday 15 April 2008 12UTC  
Surface: Total precipitation of at least 5 mm

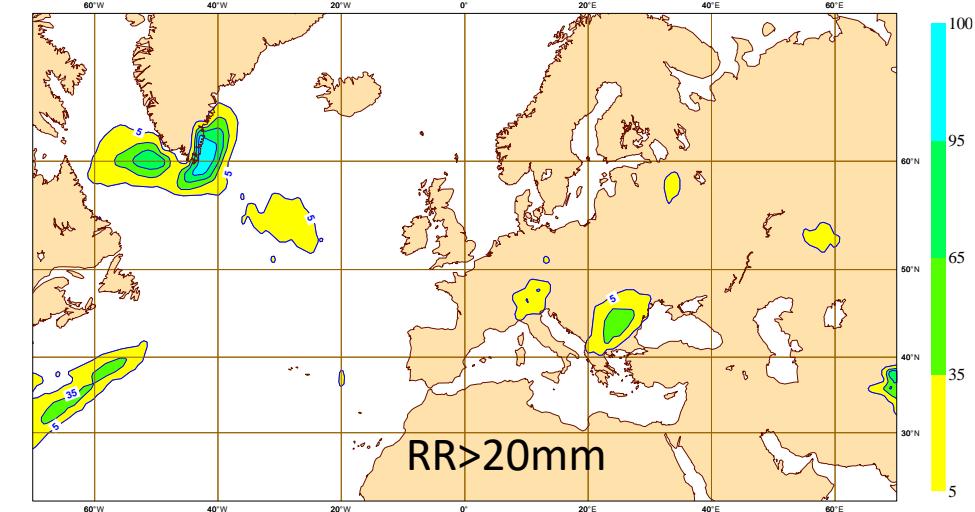
Sunday 13 April 2008 00UTC ©ECMWF Forecast probability t+036-060 VT: Monday 14 April 2008 12UTC - Tuesday 15 April 2008 12UTC  
Surface: Total precipitation of at least 5 mm



Sunday 13 April 2008 00UTC ©ECMWF Forecast probability t+036-060 VT: Monday 14 April 2008 12UTC - Tuesday 15 April 2008 12UTC  
Surface: Total precipitation of at least 10 mm



Sunday 13 April 2008 00UTC ©ECMWF Forecast probability t+036-060 VT: Monday 14 April 2008 12UTC - Tuesday 15 April 2008 12UTC  
Surface: Total precipitation of at least 20 mm



# Communication Checklist

1. Why was this work done?
2. For whom was it done?
3. To whom do you want to communicate information about the work?
4. Why would they be interested?
5. What information for what audiences?
6. Who may benefit from the work?
7. Who originated it?

# Communication Checklist

8. What exchange, style and content of memoranda were needed to clarify the purpose of the project?
9. What communication measures were needed to establish high quality and timely data collection?
10. What support was needed from colleagues or specialists?
11. What progress memoranda and reports were written and for whom?

# Tables

- Right justify numbers in tables;
- Line up decimal points in columns;
- Round numbers so that the two most effective digits are visible;
- Avoid distortion of the information in the data;
- Add rows and column averages or total where these are appropriate and may help;
- Consider re-ordering rows and/or columns to make the table clearer;
- Consider transposing the table;
- Give attention to the spacing and layout of the table.

# Graphics

- Use graphs when the shape of the data, such as trends or groups, are more important than exact values;
- Be sure that the graphic shows the data, so that you persuade the reader to think about the substance rather than the methodology or graphic design;
- Design the graphic so that it encourages the reader's eye to compare different pieces of data;
- Reveal the data at several levels of detail, from a broad overview to the fine structure
- Give every graph a clear, self-explanatory title

# Graphics

- State all measurement units;
- Choose scales on graphs carefully;
- Label axes clearly;
- Avoid chart junk;
- Improve by trial-and-error since you rarely get the graphic right first time;
- Beware of the graphic artist who aims to beautify the image but fails to elucidate the data. So insist on checking the figures after the artist has done the work.
- Beware of misleading scales.

HOWARD WAINER\*

Methods for displaying data badly have been developing for many years, and a wide variety of interesting and inventive schemes have emerged. Presented here is a synthesis yielding the 12 most powerful techniques that seem to underlie many of the realizations found in practice. These 12 (the dirty dozen) are identified and illustrated.

KEY WORDS: Graphics; Data display; Data density; Data-ink ratio.

#### 1. INTRODUCTION

The display of data is a topic of substantial contemporary interest and one that has occupied the thoughts of many scholars for almost 200 years. During this time there have been a number of attempts to codify standards of good practice (e.g., ASME Standards 1915; Cox 1978; Ehrenberg 1977) as well as a number of books that have illustrated them (i.e., Bertin 1973, 1977, 1981; Schmid 1954; Schmid and Schmid 1979; Tufte 1983). The last decade or so has seen a tremendous increase in the development of new display techniques and tools that have been reviewed recently (Macdonald-Ross 1977; Fienberg 1979; Cox 1978; Wainer and Thissen 1981). We wish to concentrate on methods of data display that leave the viewers as uninformed as they were before seeing the display or, worse, those that induce confusion. Although such techniques are broadly practiced, to my knowledge they have not as yet been gathered into a single source or carefully

#### How to Display Data Badly

categorized. This article is the beginning of such a compendium.

The aim of good data graphics is to display data accurately and clearly. Let us use this definition as a starting point for categorizing methods of bad data display. The definition has three parts. These are (a) showing data, (b) showing data accurately, and (c) showing data clearly. Thus, if we wish to display data badly, we have three avenues to follow. Let us examine them in sequence, parse them into some of their component parts, and see if we can identify means for measuring the success of each strategy.

#### 2. SHOWING DATA

Obviously, if the aim of a good display is to convey information, the less information carried in the display,

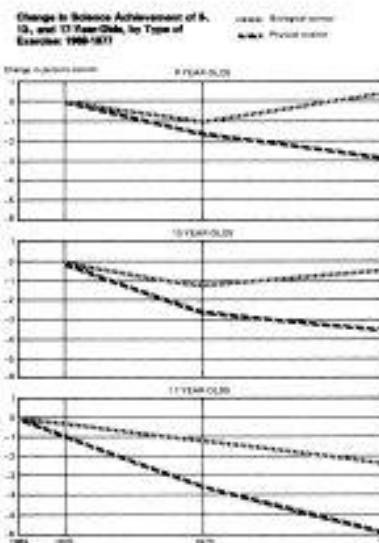


Figure 1. An example of a low density graph (from S13 ( $dd = .3$ )).

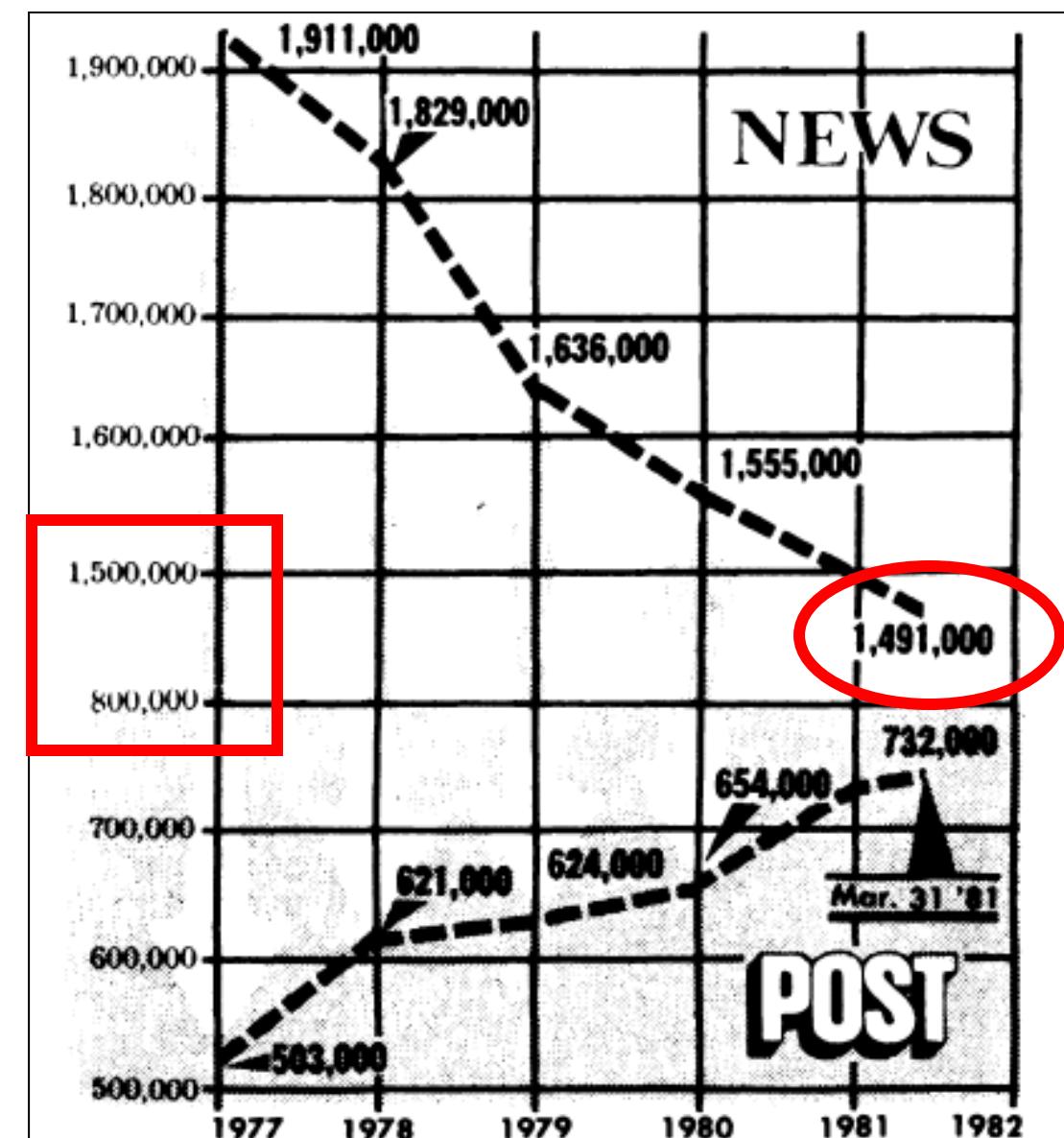
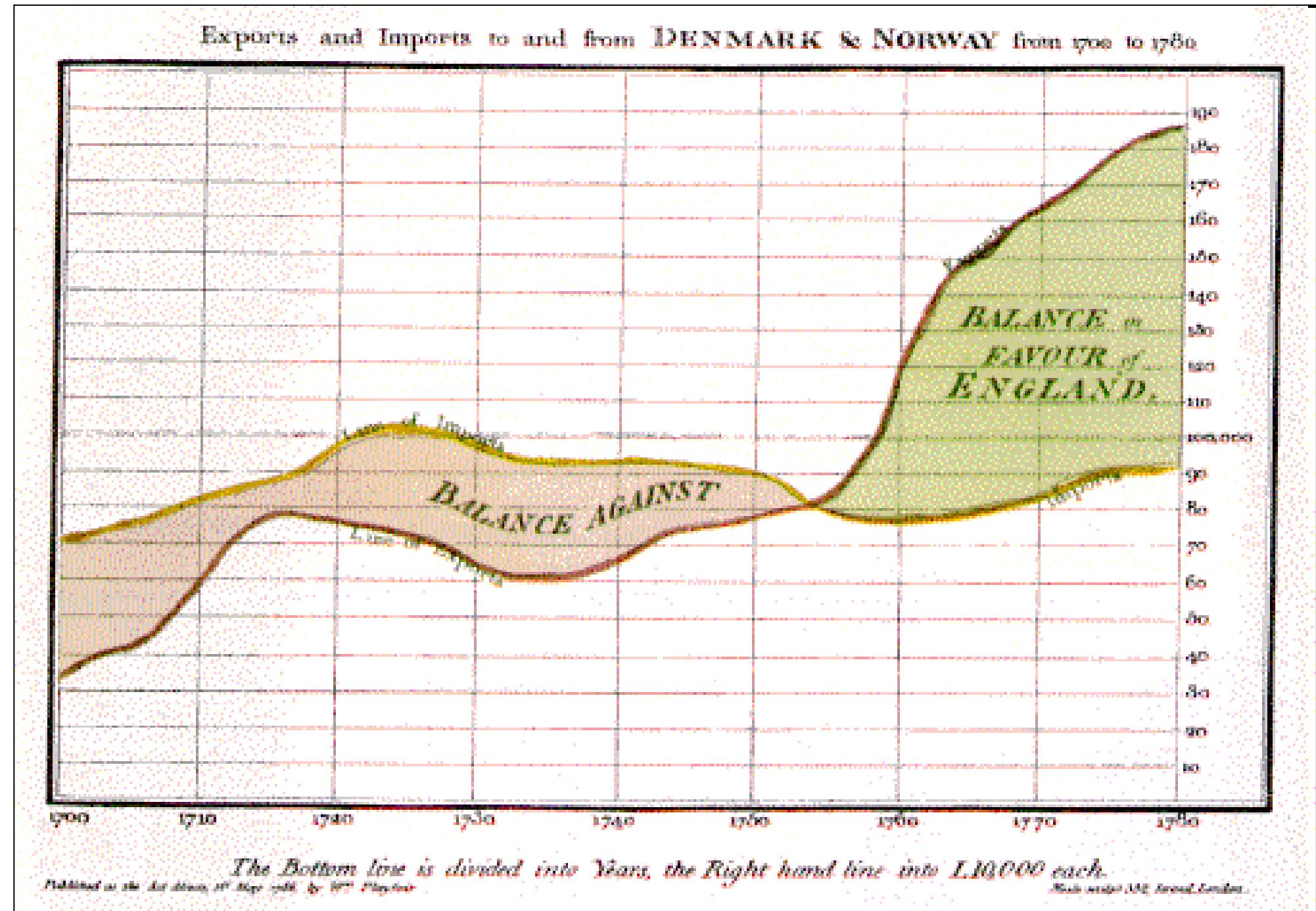


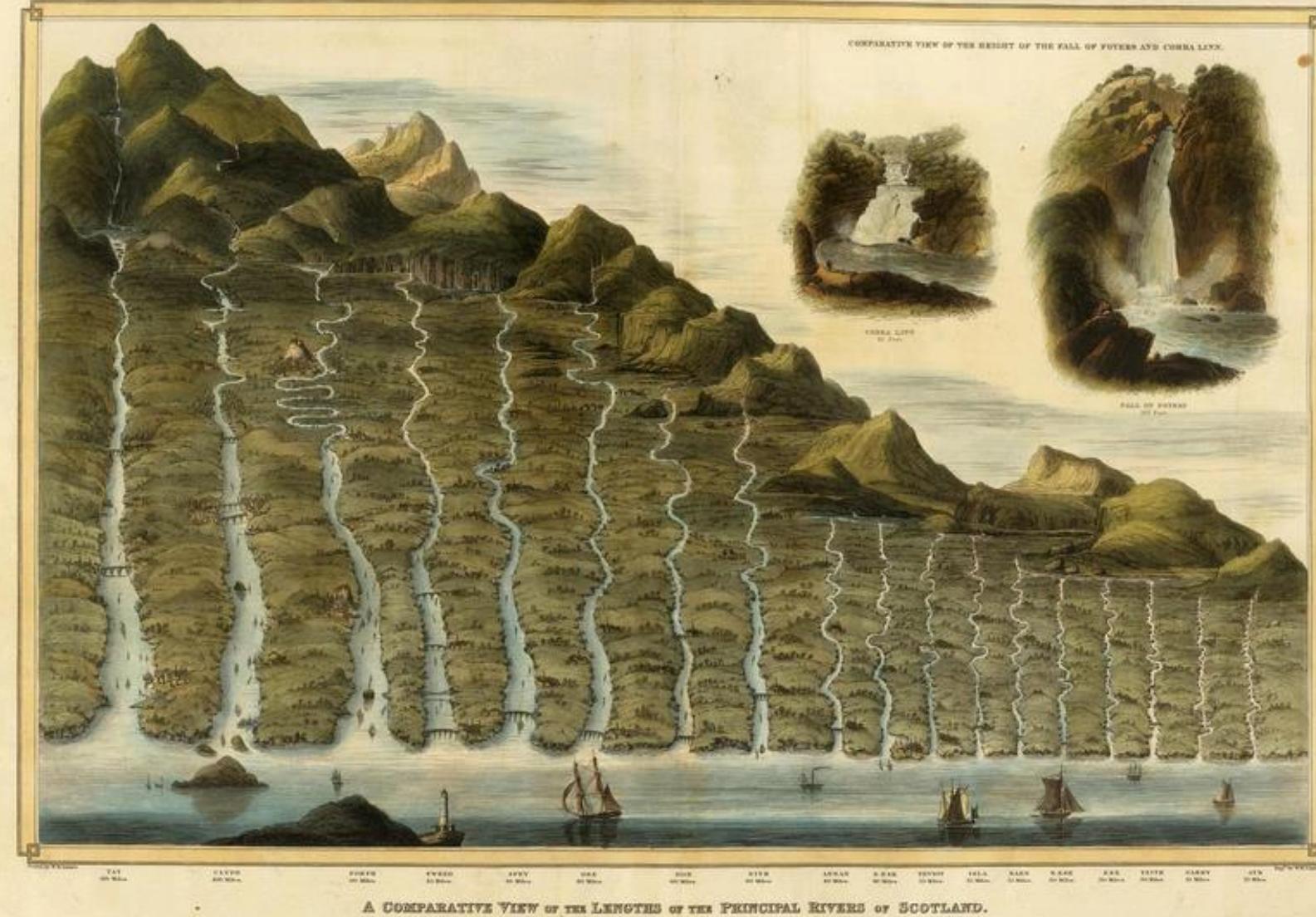
Figure 12. Changing scale in mid-axis to make large differences small (© 1981, New York Post).

\*Howard Wainer is Senior Research Scientist, Educational Testing Service, Princeton, NJ 08541. This is the text of an invited address to the American Statistical Association. It was supported in part by the Program Statistics Research Project of the Educational Testing Service. The author would like to express his gratitude to the numerous friends and colleagues who read or heard this article and offered valuable suggestions for its improvement. Especially helpful were David Andrews, Paul Holland, Bruce Kaplan, James O. Ramsay, Edward Tufte, the participants in the Stanford Workshop on Advanced Graphical Presentation, two anonymous referees, the long-suffering associate editor, and Gary Koch.

William Playfair's trade-balance time-series chart, published in his Commercial and Political Atlas, 1786



A comparative view of the lengths of the principal rivers of Scotland with comparative view of the height of the falls of Foyers and Corra Linn (1831).



After witnessing deplorable sanitary conditions in the Crimea, Florence Nightingale wrote *Notes on Matters Affecting the Health, Efficiency and Hospital Administration of the British Army* (1858), including several graphs of her own design, which she called "Coxcombs". This figure makes it clear that far more deaths were attributable to non-battle causes ("preventable causes") than to battle-related causes

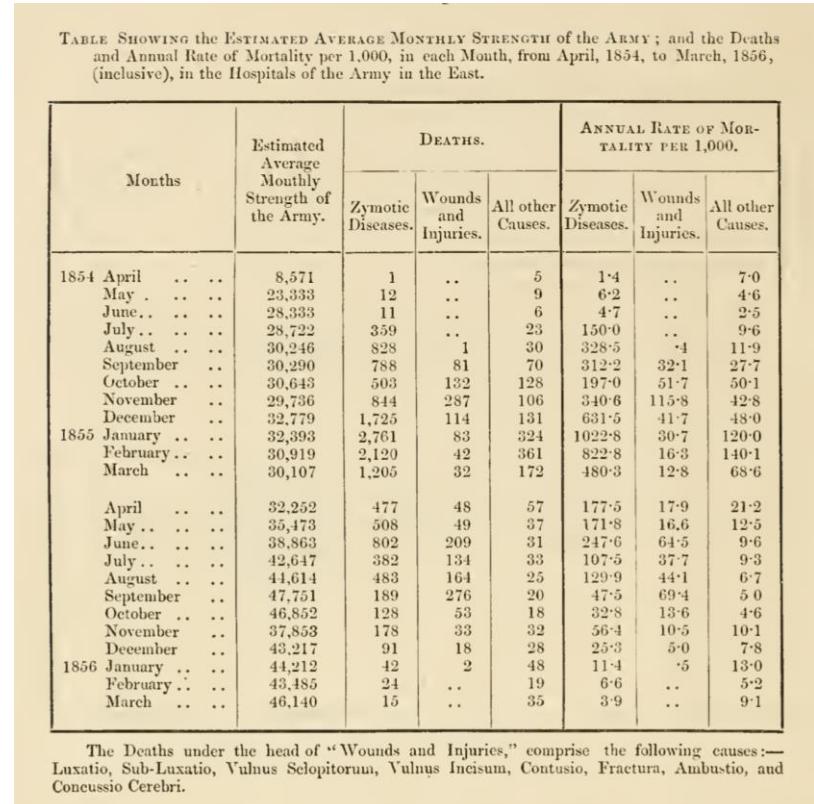
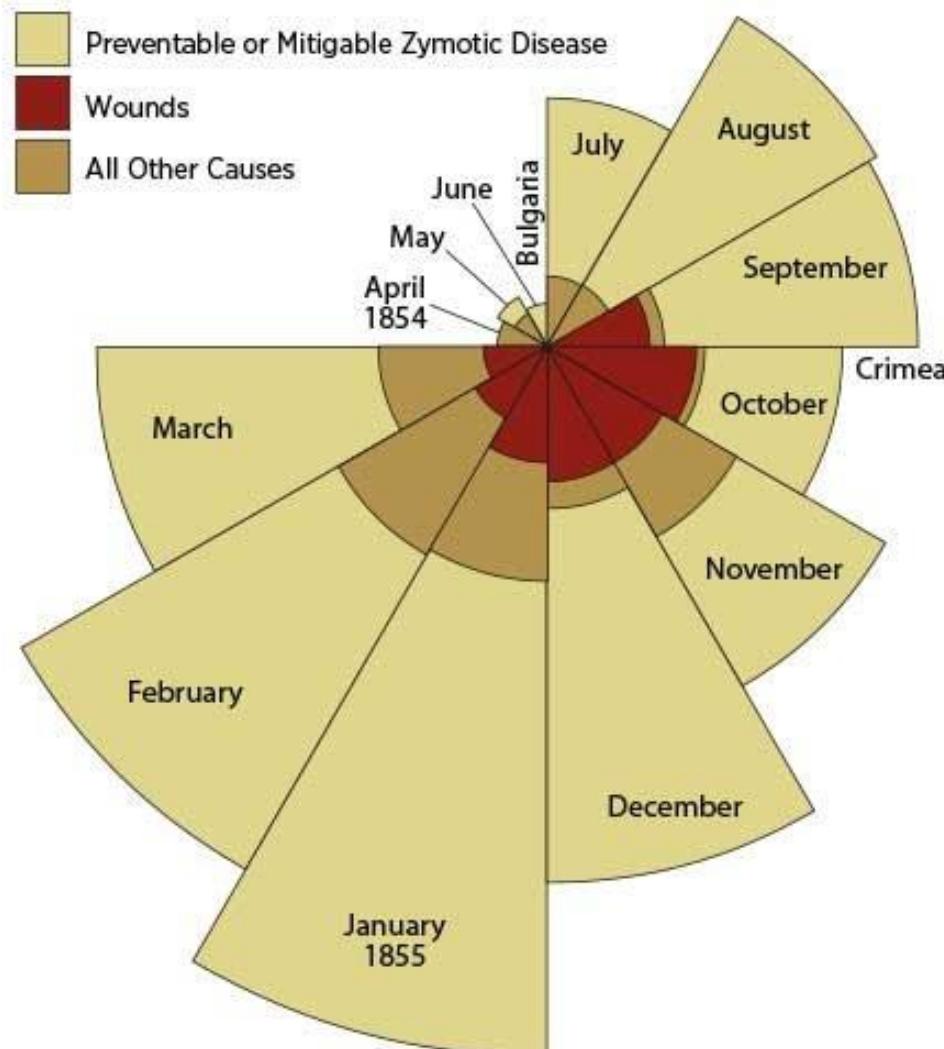


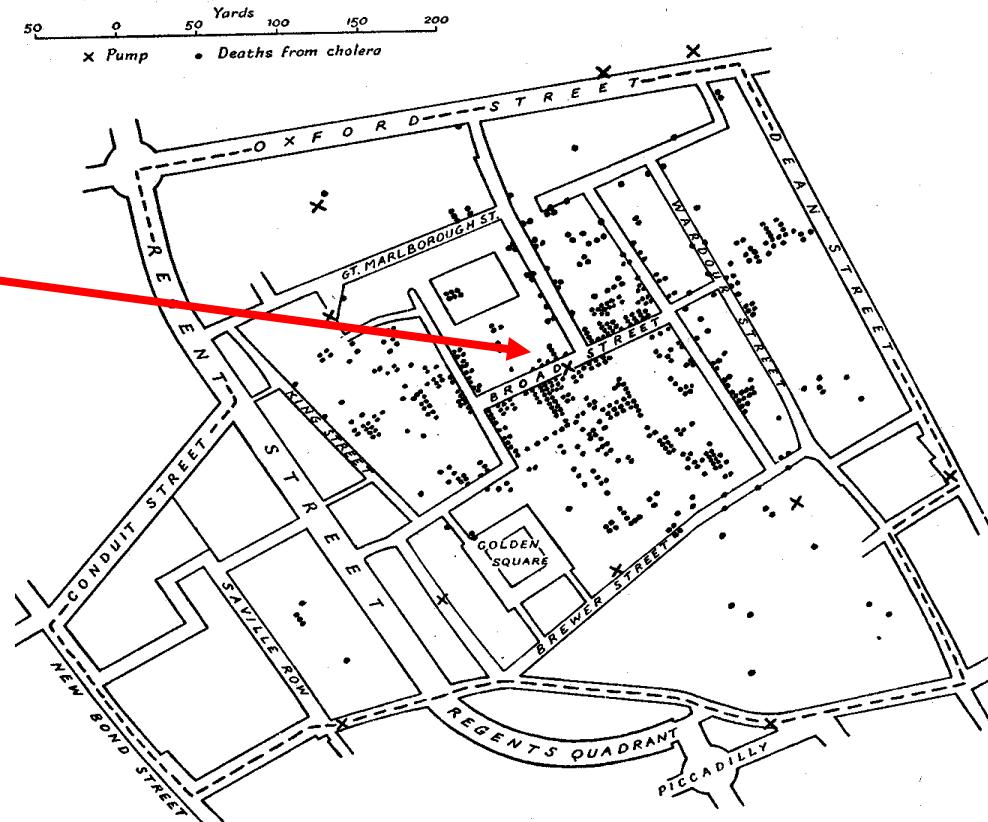
Diagram of the Causes of Mortality in the Army in the East



The black line across November 1854 marks the boundary of the deaths from all other causes during that month. In October 1854, the black coincides with the red.

Florence Nightingale  
1856

The cholera outbreak in Soho, England, in 1854.  
John Snow (1813 –1858)



# Hans Rosling

<http://www.gapminder.org/videos/ted-talks/hans-rosling-ted-2006-debunking-myths-about-the-third-world/>

Debunking myths about the "third world" - Gapminder.org - Mozilla Firefox

File Edit View History Bookmarks Tools Help

Most Visited Links

City Hotel Ljubljana Debunking myths about the ...

Unveiling the beauty of statistics for a fact based world view.

HOME  
GAPMINDER WORLD  
BLOG  
VIDEOS  
DOWNLOADS  
FAQ  
ABOUT

Search

DEBUNKING MYTHS ABOUT THE "THIRD WORLD"

Posted November 14, 2008 Comments(19)

TED Ideas worth spreading

About this talk

You've never seen data presented like this. Hans Rosling's presentation at the TED-conference in 2006 has been seen by millions over the internet, at [TED's web-page](#), at [Google Video](#) or [Youtube](#).

With the urgency of a sportscaster, Hans Rosling debunks myths about the so-called "developing world" using the animation software that powers [Gapminder World](#).

Related Content

[Flash-presentation used.](#)

Transferring data from video.ted.com...

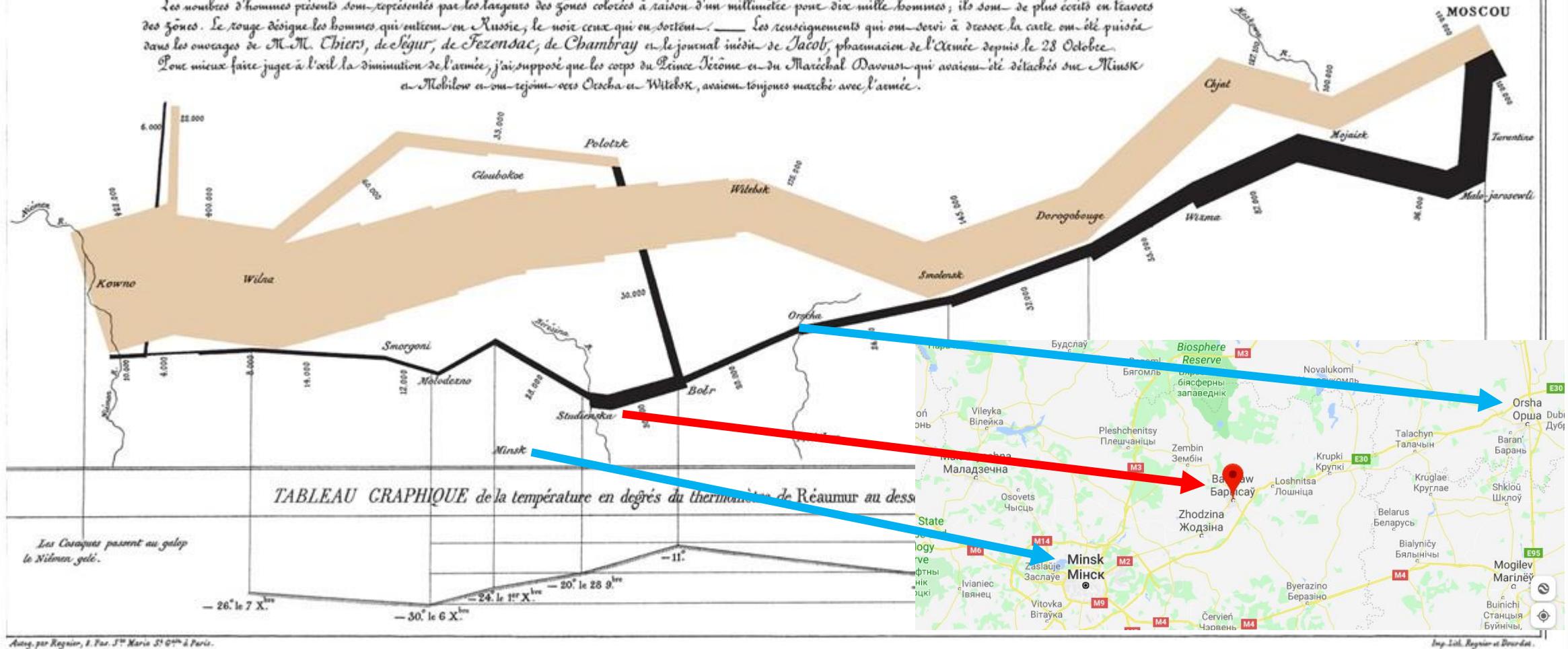
Start Deb... 2009 Dror... Deb... 10 1... Mic... 78% 11:18 AM

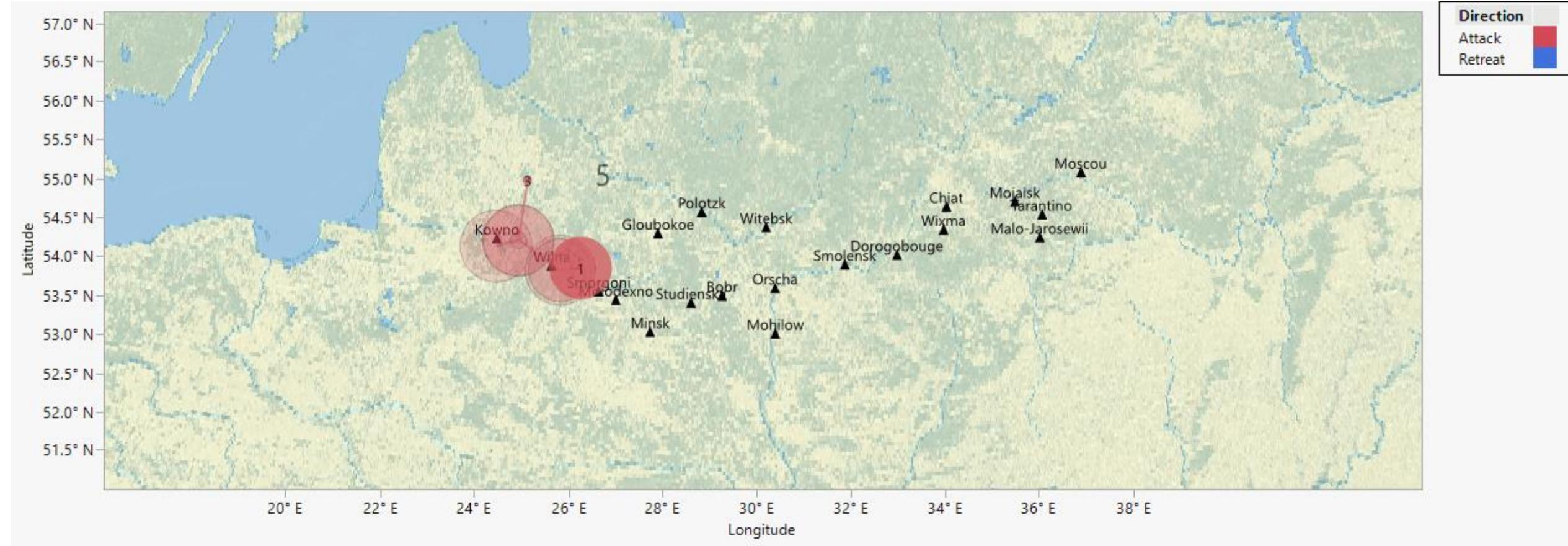
# Carte Figurative des pertes successives en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.

Dessinée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite Paris, le 20 Novembre 1869.

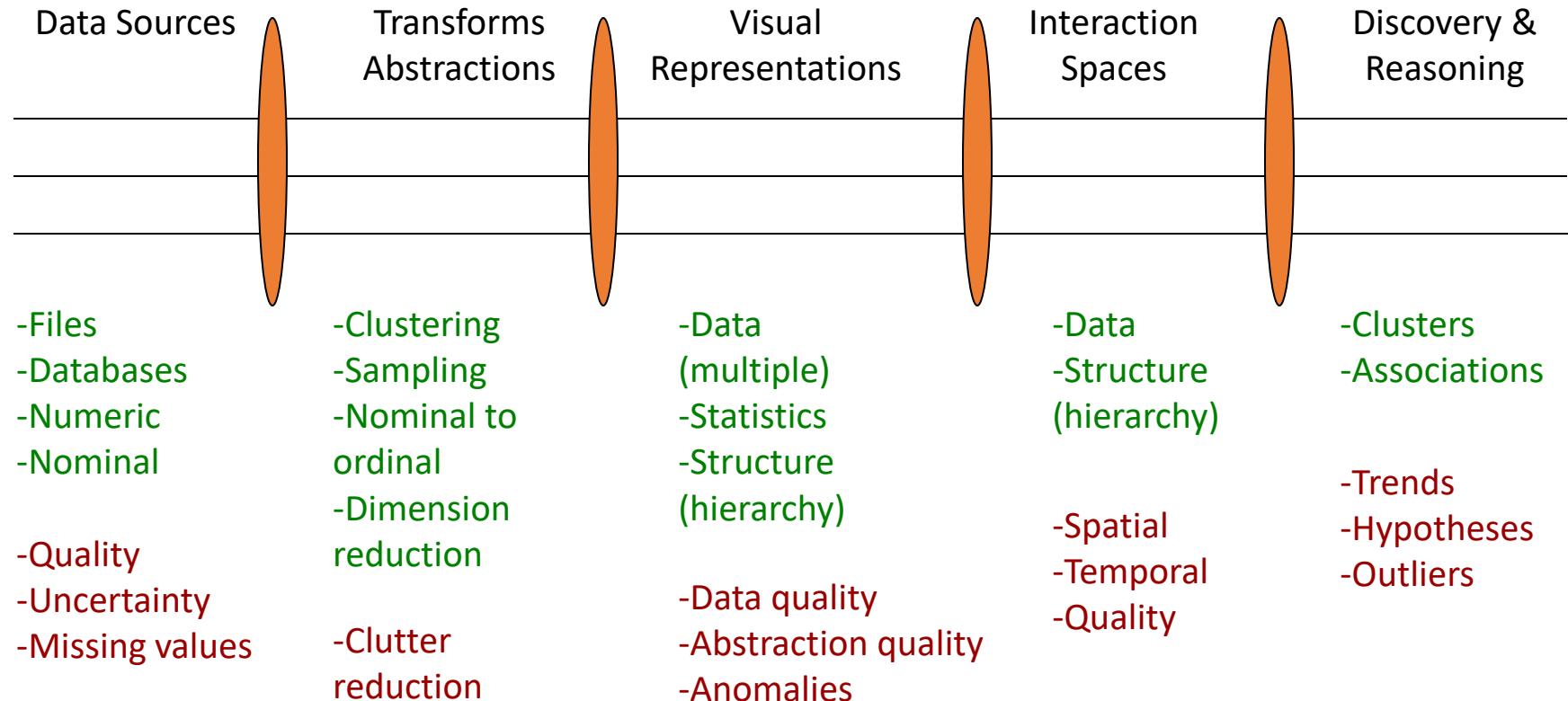
Les nombres d'hommes perdus sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en lettres des zones. Le rouge désigne les hommes qui ont péri en Russie; le noir ceux qui en sont sortis. — Les renseignements qui ont servi à dessiner la carte ont été puisés dans les ouvrages de M. M. Chiers, de Séguir, de Ferondac, de Chambray et le journal intime de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk au Mobilier et avaient rejoint vers Orsha et Wilno, avaient toujours marché avec l'armée.



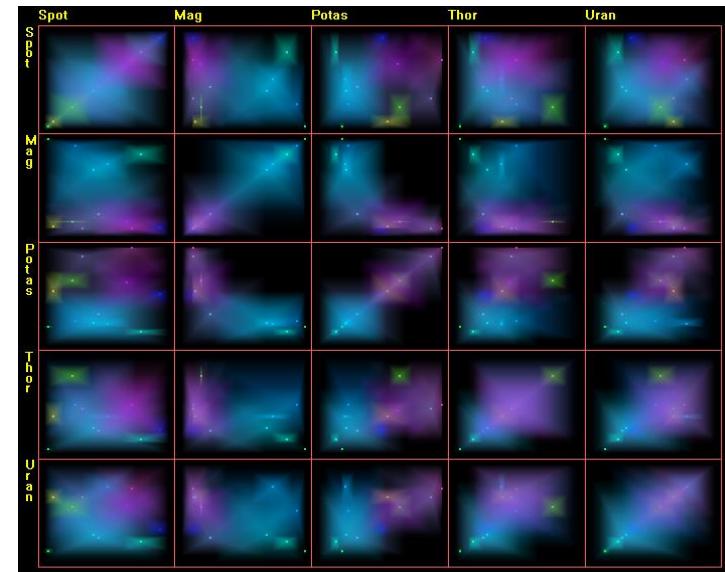
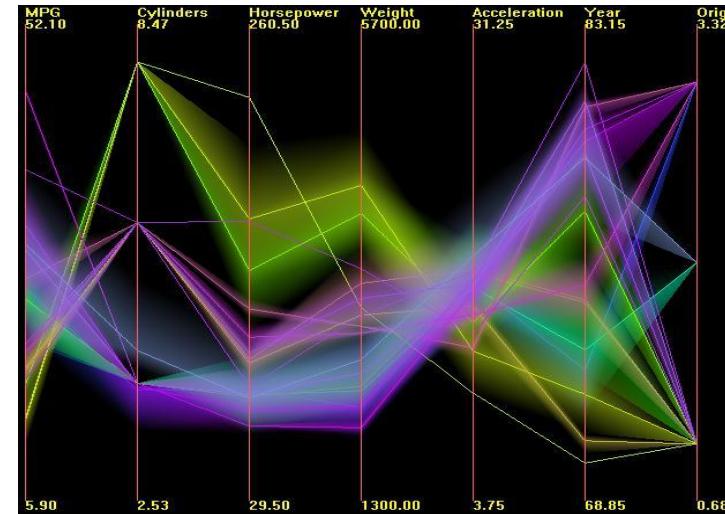


# Visual Analytics



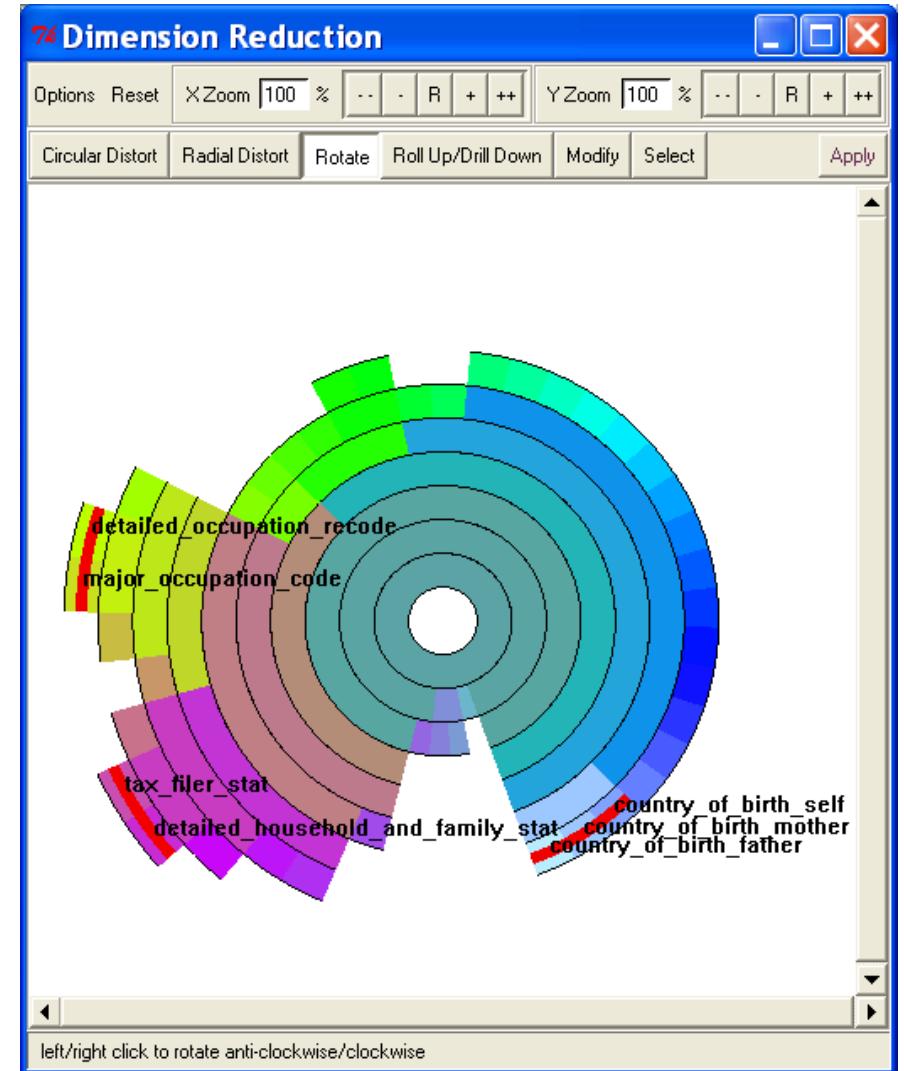
# Multiresolution Visualization

- For large datasets, visualizations quickly get cluttered
- Hierarchical clustering generates many levels of detail
- User can select areas of interest to view at full resolution while the rest of the data is shown via cluster centers and extents (shown as bands of variable opacity)



# Dimension Reduction

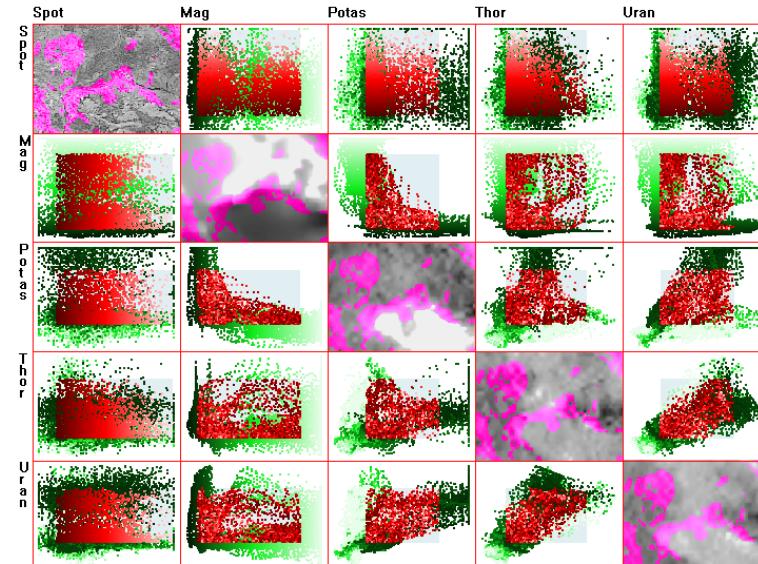
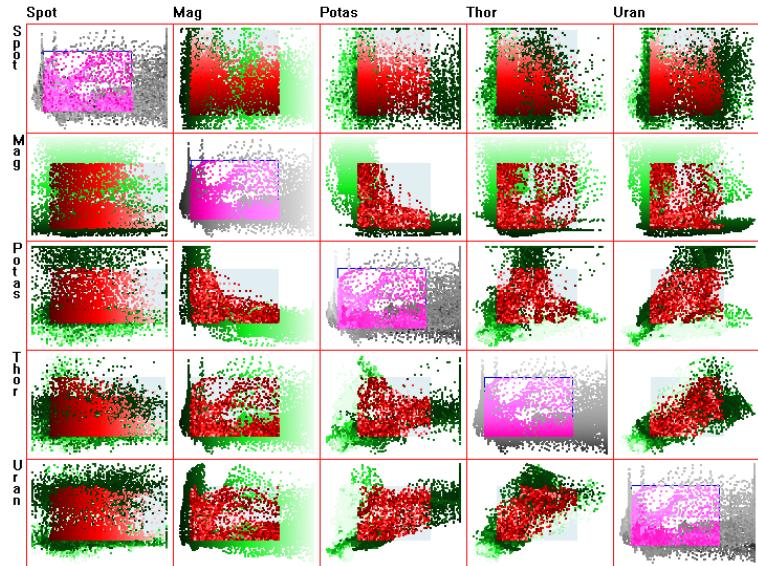
- Dimensions are hierarchically clustered based on similarity measures
- Hierarchy displayed using Inter Ring
- Users select clusters of dimensions or representative dimensions for detailed analysis



42 dimension census dataset.

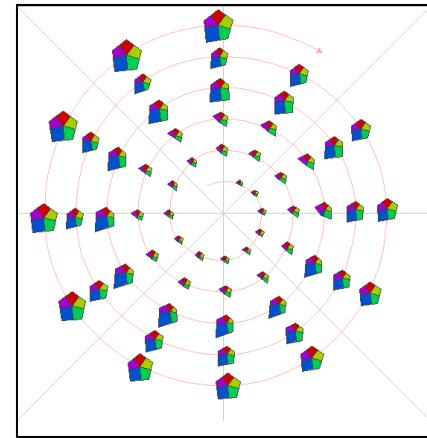
# Linking Spatial and Non-Spatial

- Diagonal plots of scatterplot matrix can have numerous uses
- Example shows multispectral remote sensing data, 1 layer per diagonal plot
- User can select in either 2-D or parameter space and see corresponding elements in other views.

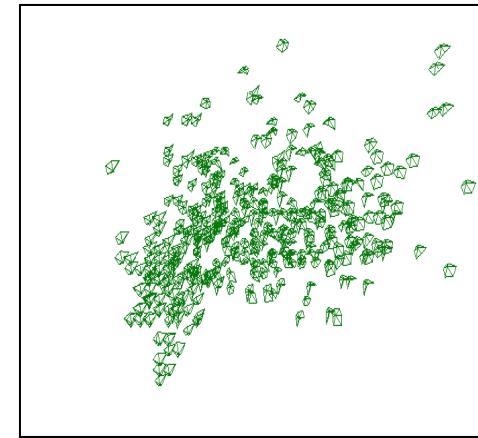


# Layout Strategies

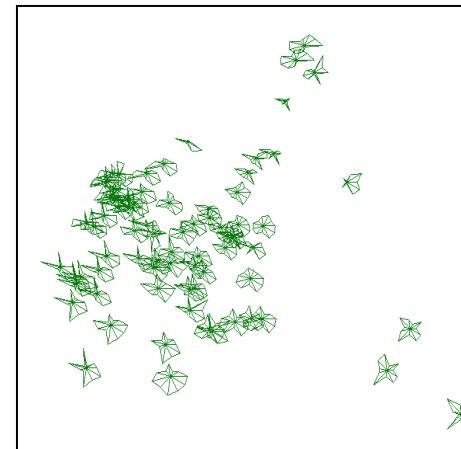
- Different layout strategies can reveal different patterns in the data
- Detecting, classifying, and measuring trends, outliers, repeated patterns, clusters, and correlations can be facilitated via specific layouts



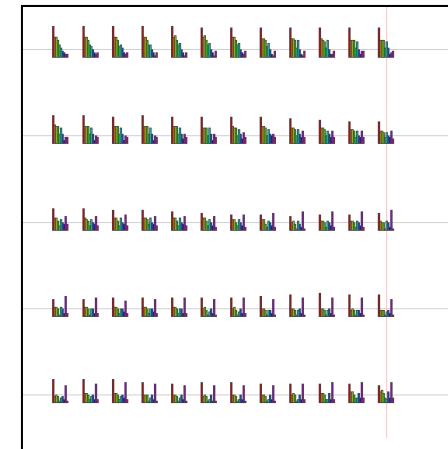
Cyclic



Data Driven



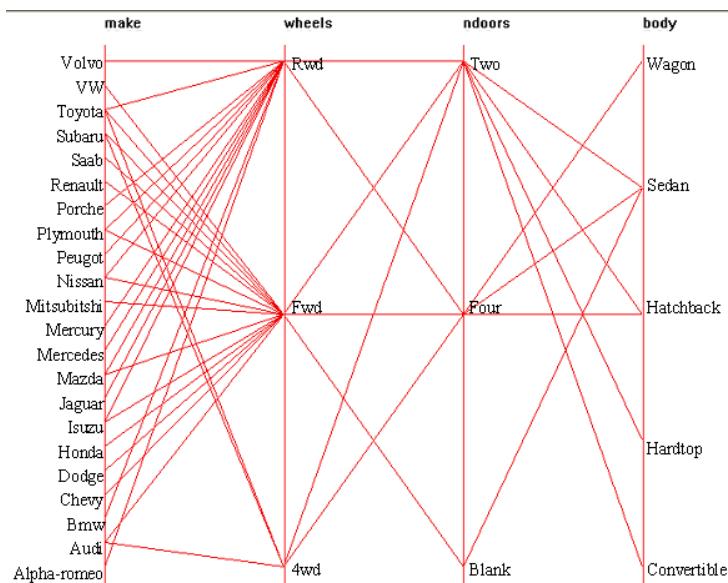
Principal Components



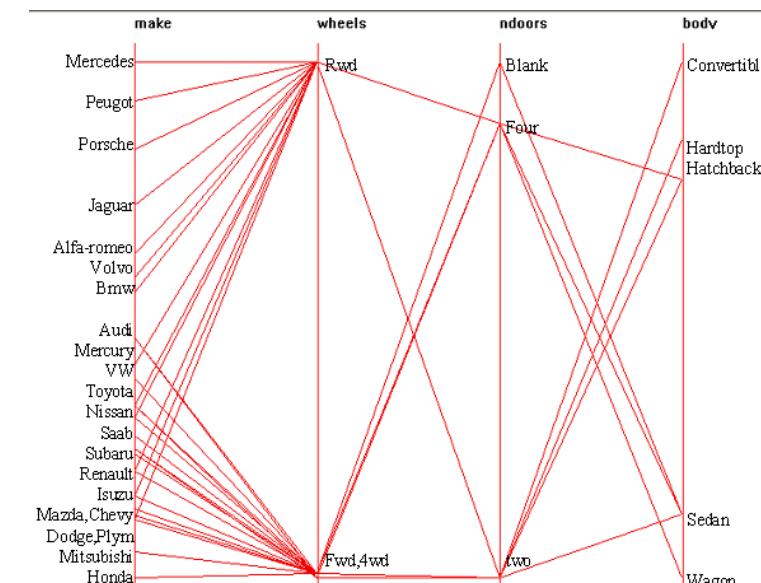
Order Driven

# Visualizing Data with Parallel Plots

- Arbitrary assignment of non-numeric fields to numbers can lead to misinterpretation, lost patterns
- By looking at similarities in distributions across all dimensions, we can group values of a nominal variable with similar global characteristics
- Assignments used to convey order and relative distance



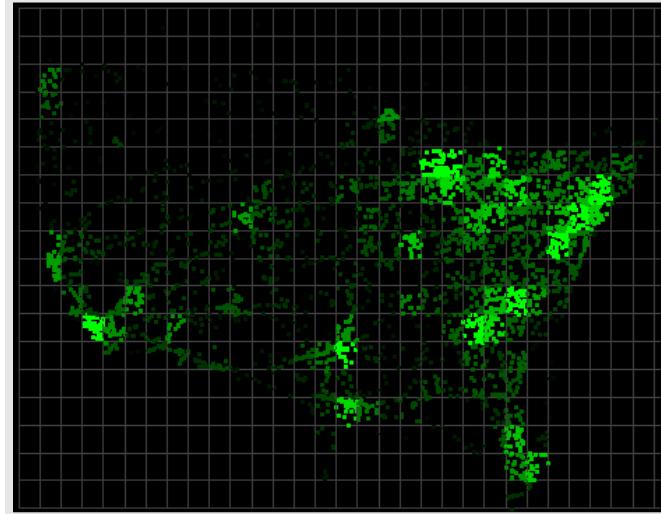
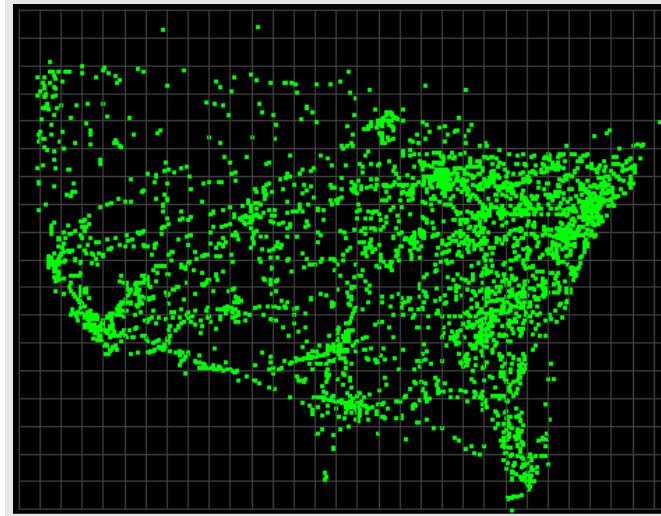
Original Assignment



Assignment after Correspondence Analysis

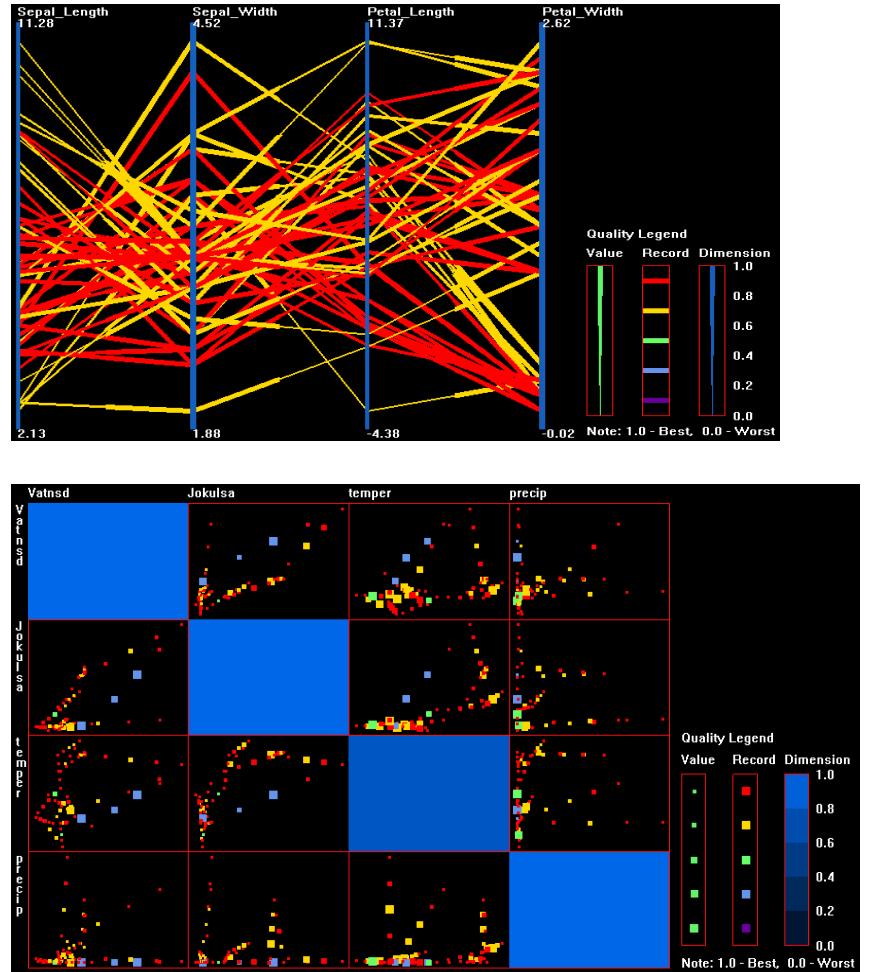
# Visual Clutter Reduction

- In scenes with thousands of moving objects, there is need to reduce clutter
- Many strategies, including:
  - Information-preserving
  - Information-reducing
  - Visual remapping



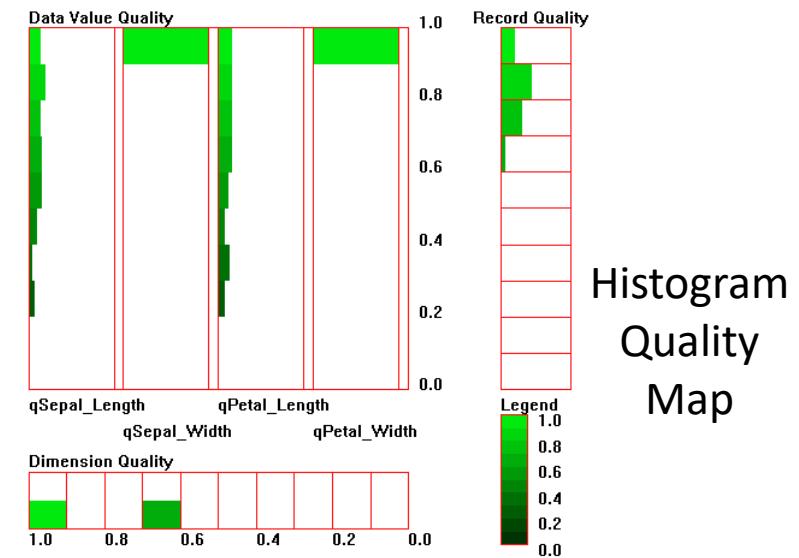
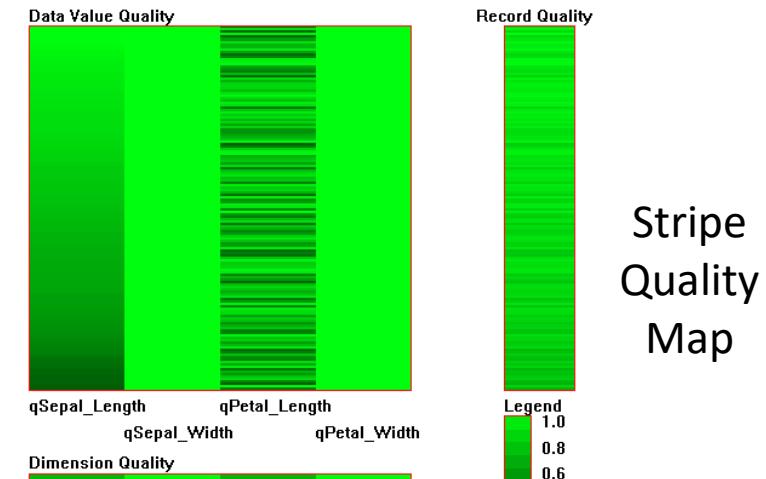
# Data Quality Visual Encoding

- Data quality refers to the degree of uncertainty of data
- Quality measures are visually encoded into existing visualizations
- This helps users focus on high quality data to draw reliable conclusions



# Quality Space Visualization

- Quality space is visualized separately to convey patterns in the data quality measures
- Records or dimensions can be ordered by quality to reveal structure and relations
- Stripe view shows individual data value quality; Histogram view shows summarization and distribution

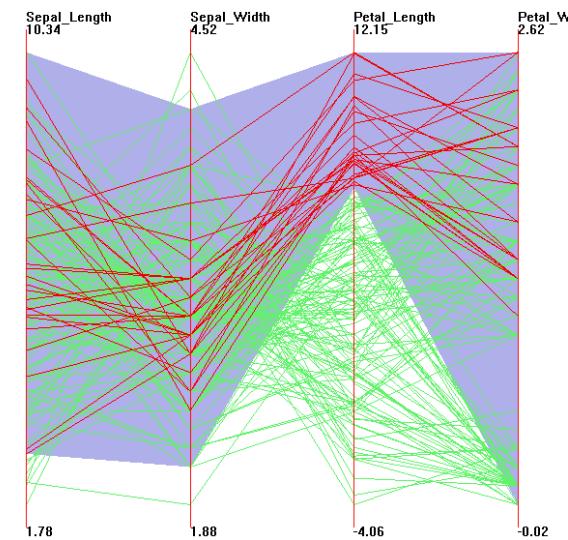


Stripe  
Quality  
Map

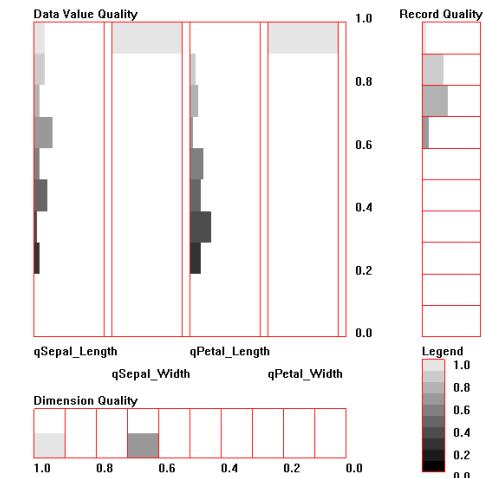
Histogram  
Quality  
Map

# Interactions between Data Space and Quality Space

- Linking brush: When users select a subset in one space, the corresponding subset in the other space will be highlighted accordingly.
- Sample figures: The data points in the data space with high values in the third dimension are highlighted, then the distribution of quality measures for this subset is rendered in the quality map.

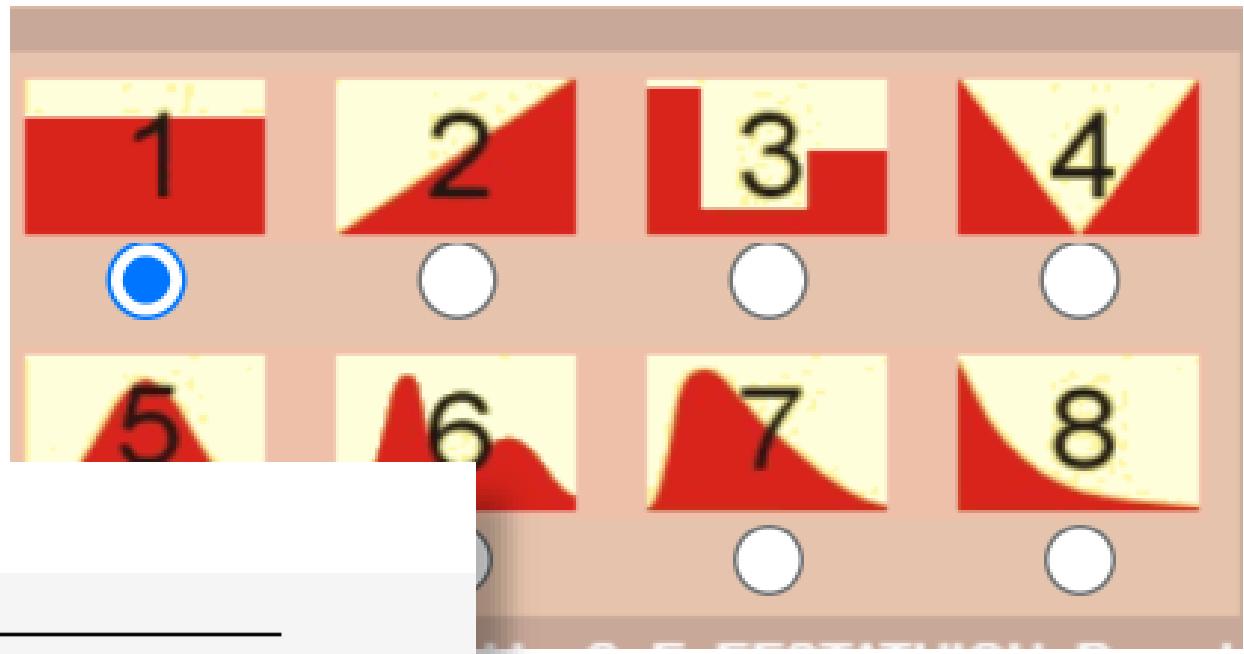


Data space  
with highlighting



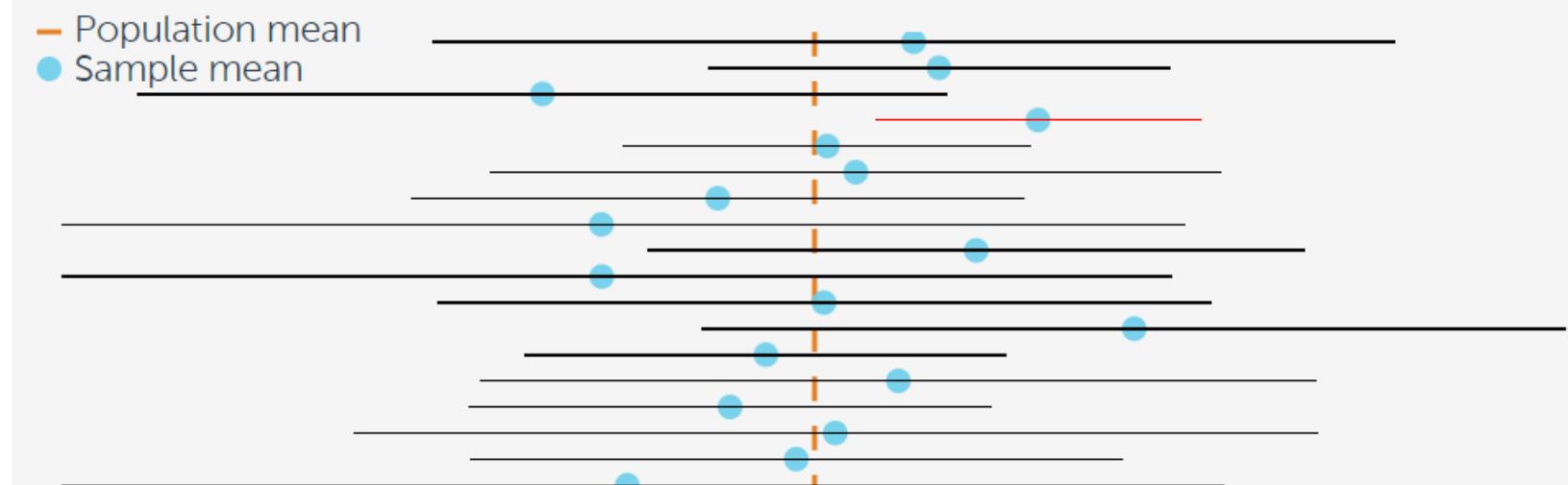
Linked  
Quality space

[http://195.134.76.37/applets/AppletCentralLimit/App\\_CentralLimit2.html](http://195.134.76.37/applets/AppletCentralLimit/App_CentralLimit2.html)



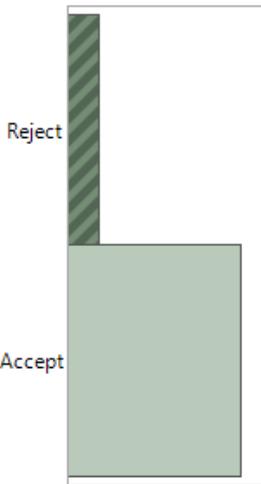
<https://rpsychologist.com/d3/CI/>

## 95% confidence intervals

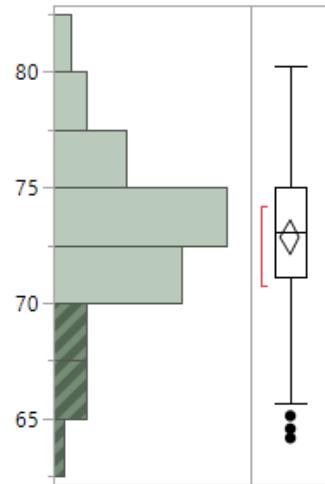


## Distributions

### Lot Acceptance

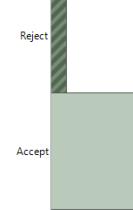


### Disso

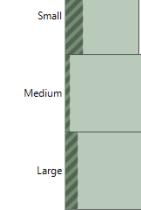


## Distributions

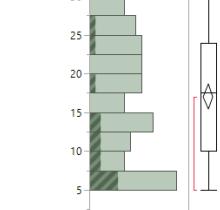
### Lot Acceptance



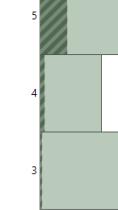
### API Particle Size



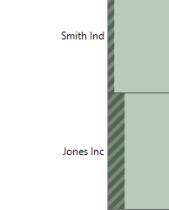
### Mill Time



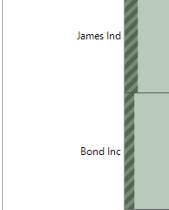
### Screen Size



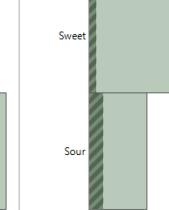
### Mag. Stearate Supplier



### Lactose Supplier



### Sugar Supplier



Tablet Production.jmp

### Frequencies

| Level     | Count | Prob    |
|-----------|-------|---------|
| Accept    | 76    | 0.84444 |
| Reject    | 14    | 0.15556 |
| Total     | 90    | 1.00000 |
| N Missing | 0     |         |
| 2 Levels  |       |         |

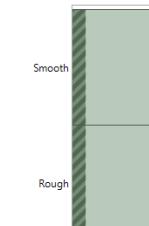
### Quantiles

|                |         |
|----------------|---------|
| 100.0% maximum | 80.23   |
| 99.5%          | 80.23   |
| 97.5%          | 80.1025 |
| 90.0%          | 77.532  |
| 75.0% quartile | 74.9925 |
| 50.0% median   | 73.05   |
| 25.0% quartile | 71.11   |
| 10.0%          | 67.844  |
| 2.5%           | 64.704  |
| 0.5%           | 64.15   |
| 0.0% minimum   | 64.15   |

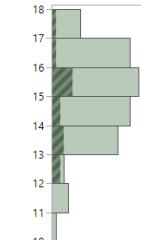
### Summary Statistics

|                |           |
|----------------|-----------|
| Mean           | 72.860556 |
| Std Dev        | 3.5121345 |
| Std Err Mean   | 0.3702115 |
| Upper 95% Mean | 73.596158 |
| Lower 95% Mean | 72.124953 |
| N              | 90        |

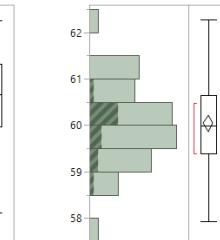
### Talc Supplier



### Blend Time



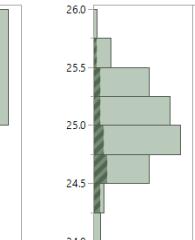
### Blend Speed



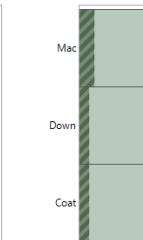
### Compressor



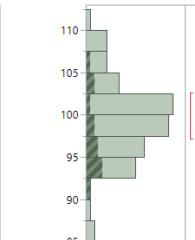
### Force



### Coating Supplier



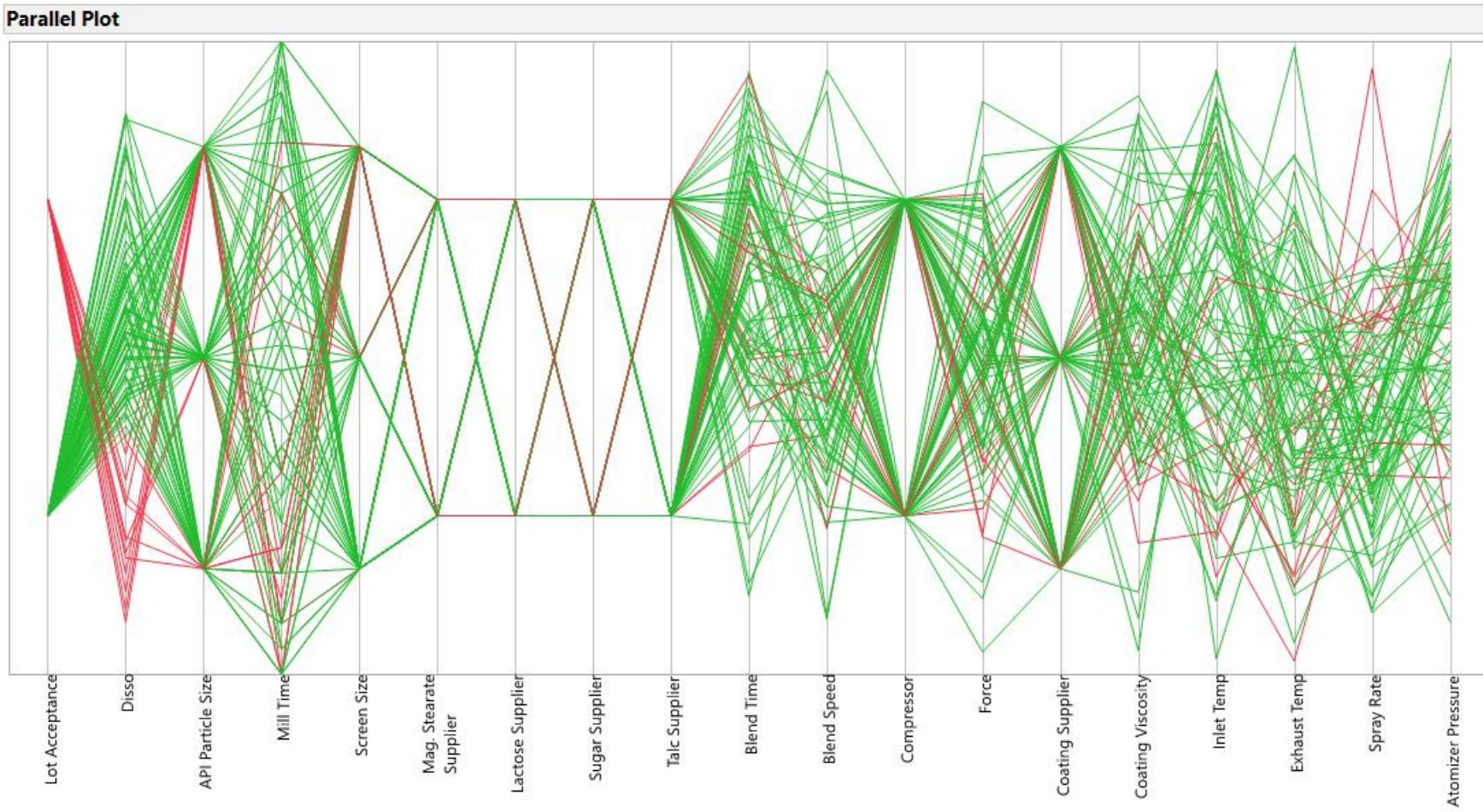
### Coating Viscosity

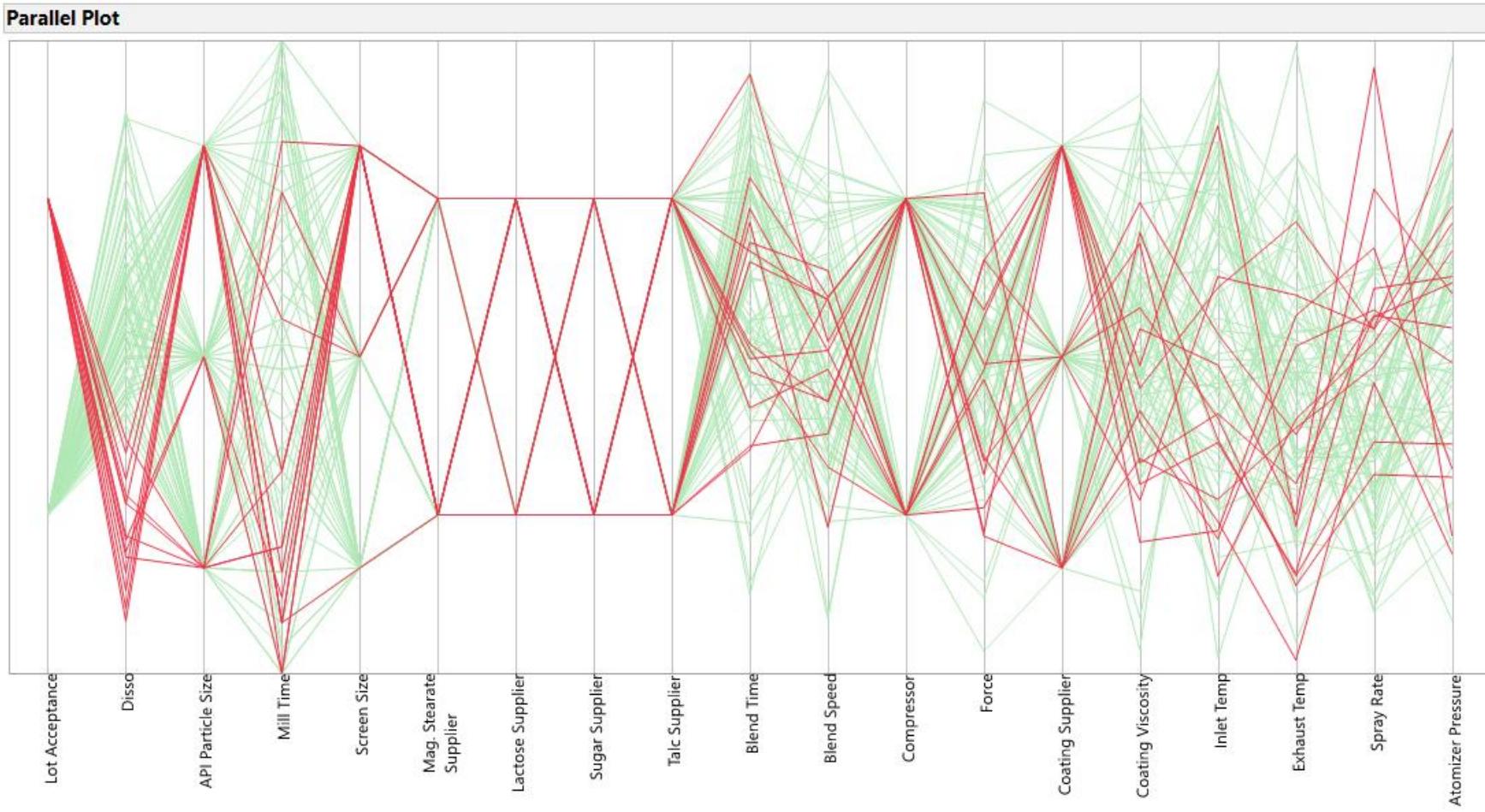


### Summary Statistics

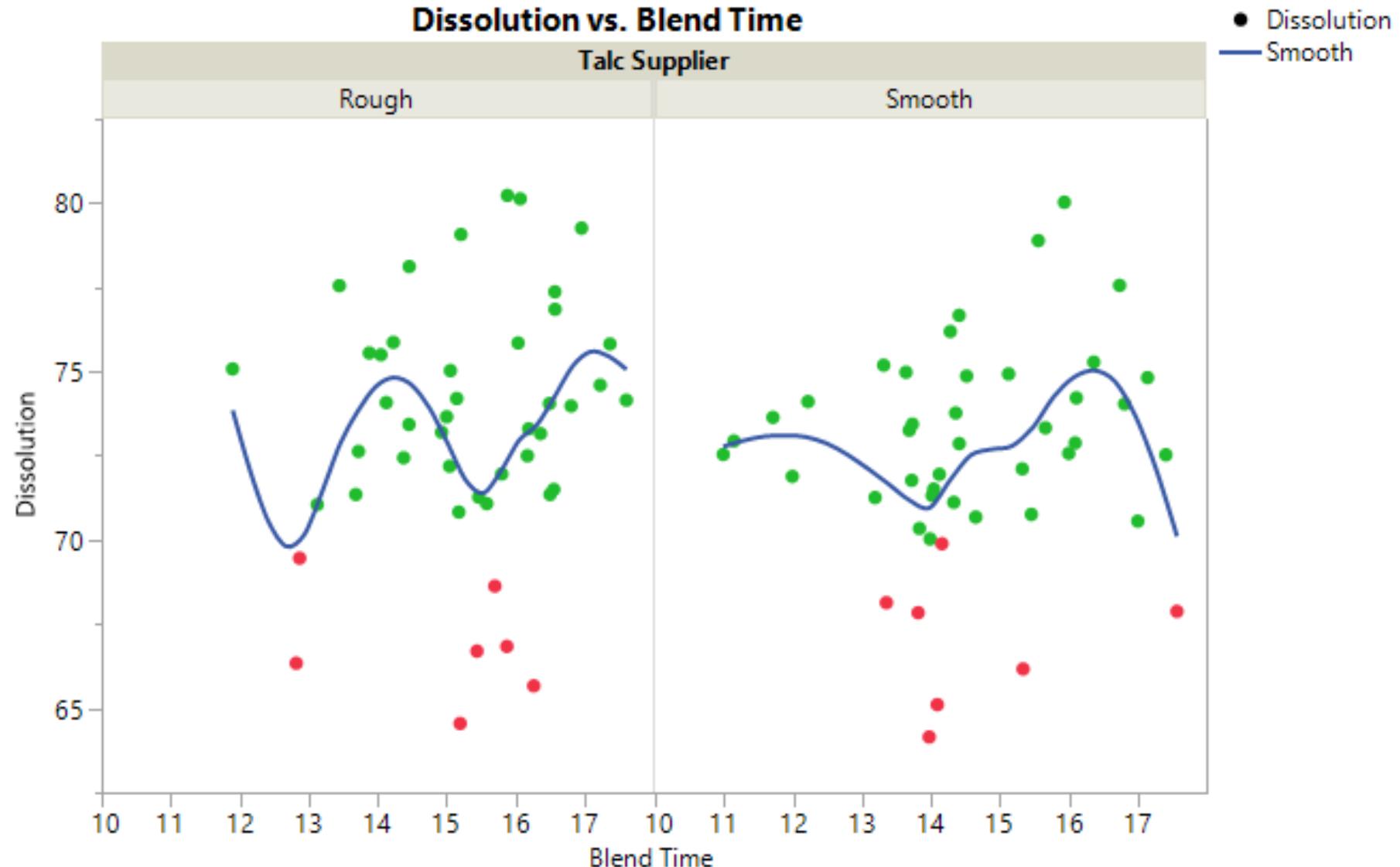
|                |           |
|----------------|-----------|
| Mean           | 99.315625 |
| Std Dev        | 4.8466816 |
| Std Err Mean   | 0.0352212 |
| Upper 95% Mean | 100.33074 |
| Lower 95% Mean | 98.300507 |
| N              | 90        |



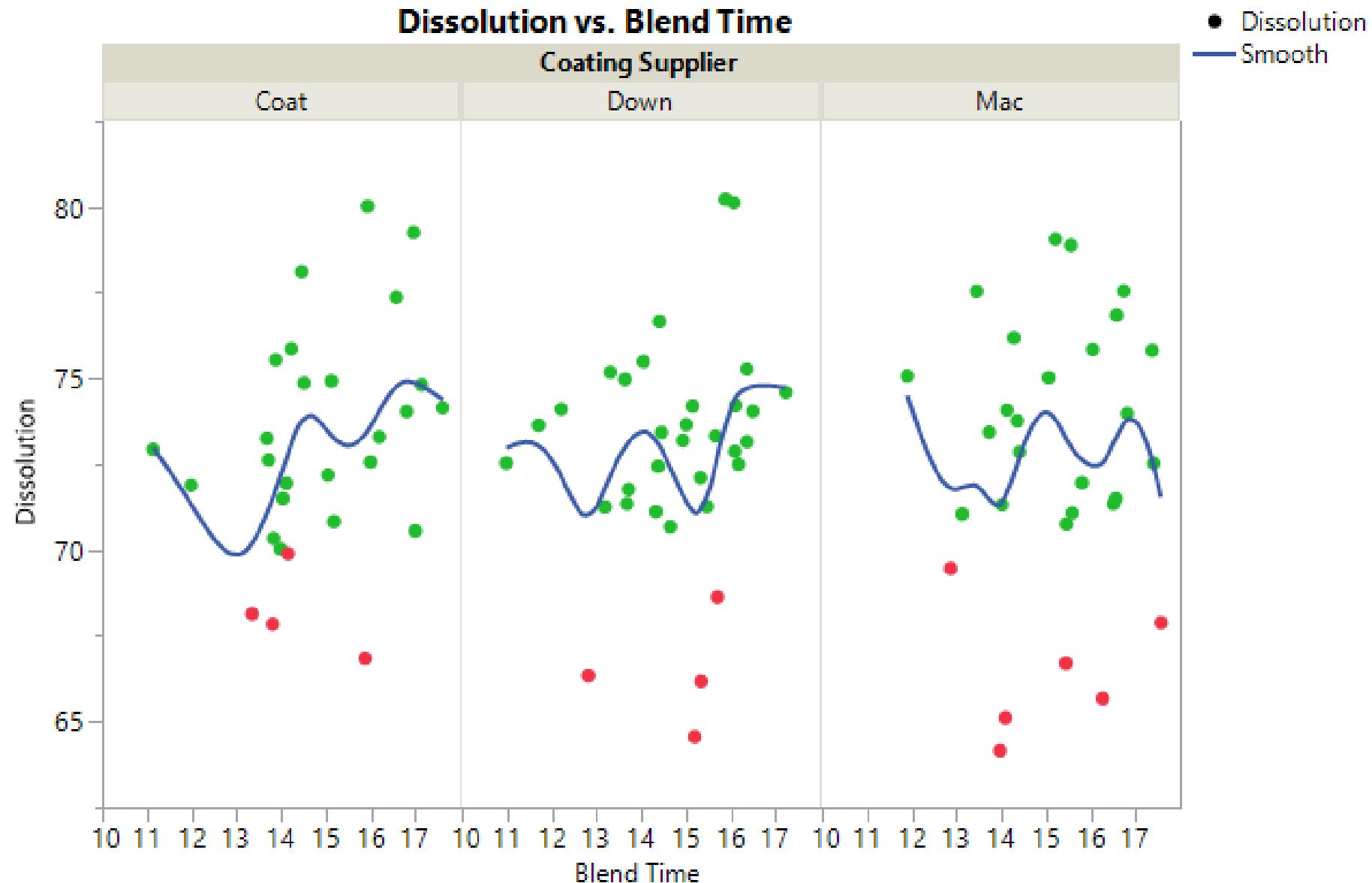




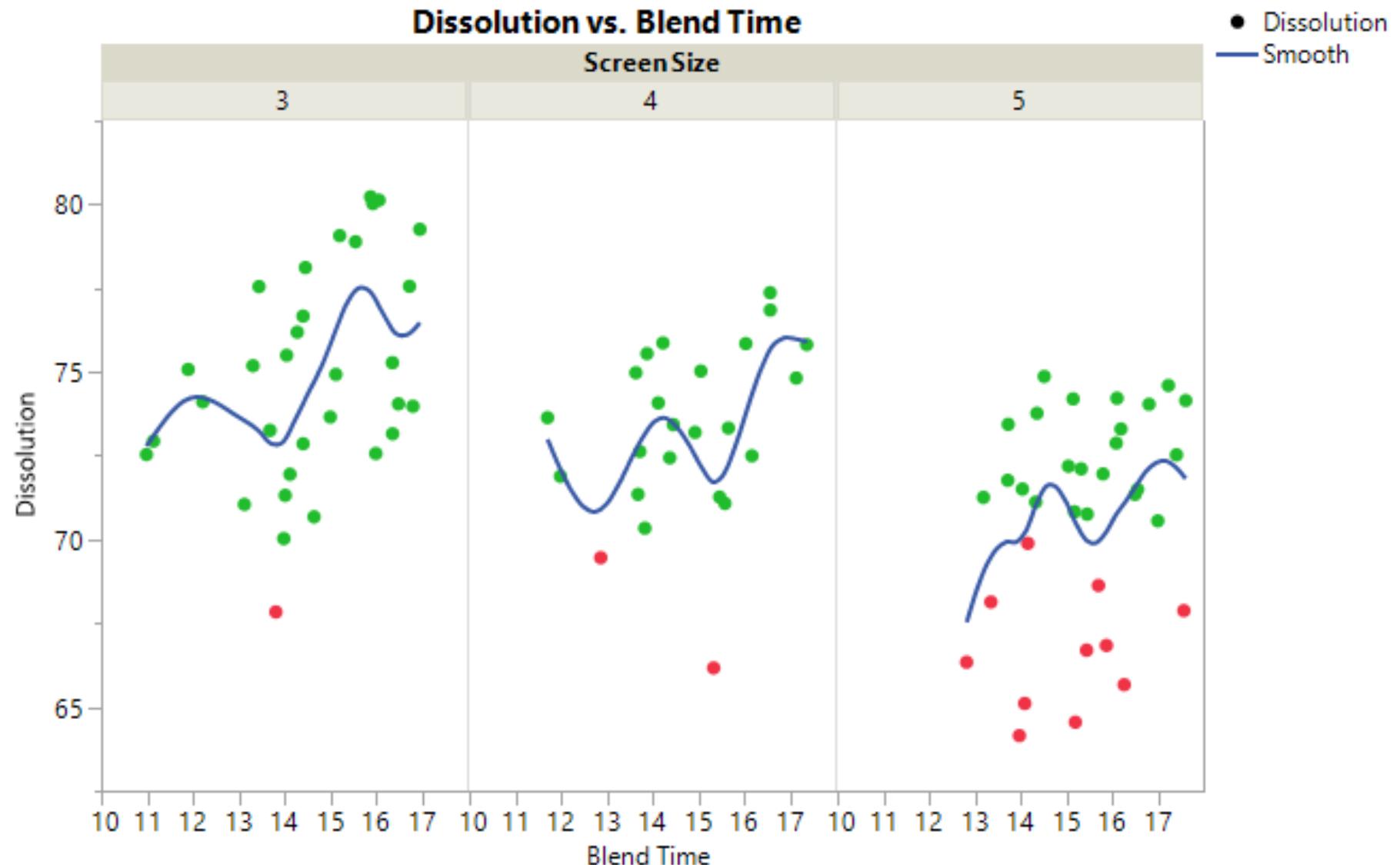
## Graph Builder

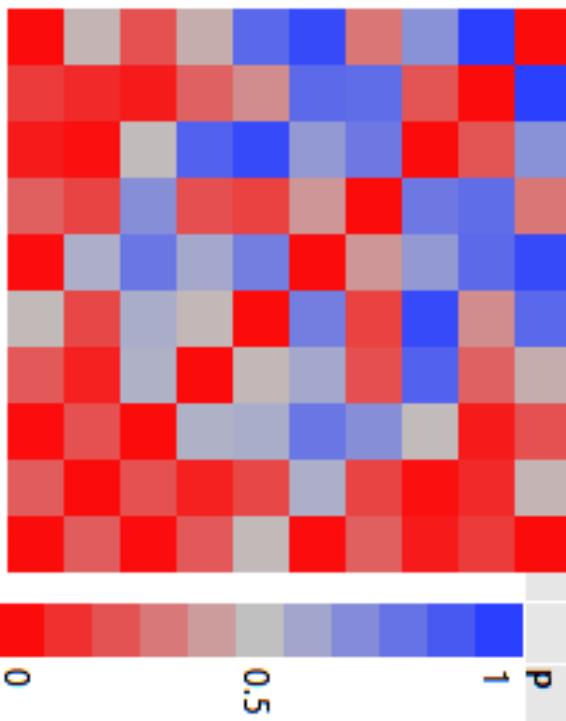


## Graph Builder



## Graph Builder





Mill Time  
Blend Time  
Blend Speed  
Force  
Coating Viscosity  
Inlet Temp  
Exhaust Temp  
Spray Rate  
Atomizer Pressure  
Dissolution

#### Explore Outliers

##### Quantile Range Outliers

Outliers are values Q times the interquartile range past the lower and upper quantiles.

Tail Quantile

Select columns and choose an action.

Q

Restrict search to integers

Show only columns with outliers

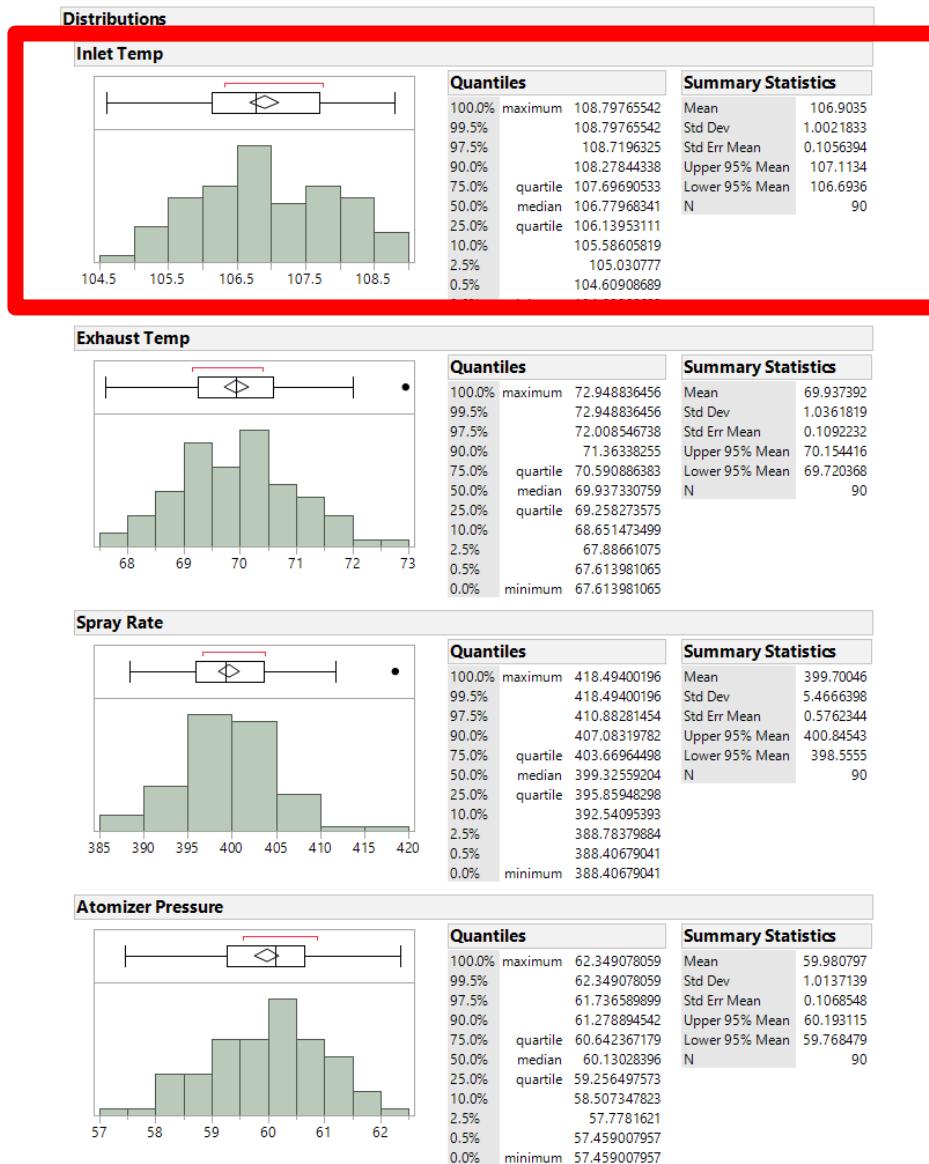
| Column            | 10% Quantile | 90% Quantile | Low Threshold | High Threshold | Number of Outliers | Outliers (Count) |
|-------------------|--------------|--------------|---------------|----------------|--------------------|------------------|
| Mill Time         | 6.1          | 28           | -59.6         | 93.7           | 0                  | 0                |
| Blend Time        | 13.1348      | 16.8123      | 2.10227       | 27.8448        | 0                  | 0                |
| Blend Speed       | 59.0087      | 61.2002      | 52.434        | 67.7748        | 0                  | 0                |
| Force             | 24.5967      | 25.4662      | 21.9885       | 28.0744        | 0                  | 0                |
| Coating Viscosity | 93.9009      | 106.366      | 56.5066       | 143.76         | 0                  | 0                |
| Inlet Temp        | 105.586      | 108.278      | 97.5089       | 116.356        | 0                  | 0                |
| Exhaust Temp      | 68.6515      | 71.3634      | 60.5157       | 79.4991        | 0                  | 0                |
| Spray Rate        | 392.541      | 407.083      | 348.914       | 450.71         | 0                  | 0                |
| Atomizer Pressure | 58.5073      | 61.2789      | 50.1927       | 69.5935        | 0                  | 0                |
| Dissolution       | 67.844       | 77.532       | 38.78         | 106.596        | 0                  | 0                |

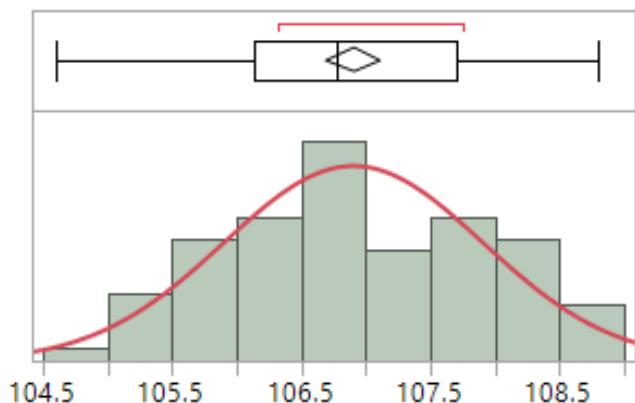
#### Correlations

|                   | Mill Time | Blend Time | Blend Speed | Force   | Coating Viscosity | Inlet Temp | Exhaust Temp | Spray Rate | Atomizer Pressure | Dissolution |
|-------------------|-----------|------------|-------------|---------|-------------------|------------|--------------|------------|-------------------|-------------|
| Mill Time         | 1.0000    | 0.0004     | -0.0436     | 0.1116  | 0.0057            | 0.0217     | 0.0810       | -0.1381    | -0.0775           | 0.3638      |
| Blend Time        | 0.0004    | 1.0000     | 0.1348      | 0.0241  | 0.0223            | 0.0977     | -0.1257      | 0.2145     | 0.1841            | 0.1598      |
| Blend Speed       | -0.0436   | 0.1348     | 1.0000      | -0.0301 | 0.0482            | 0.0059     | -0.0184      | -0.0745    | 0.2632            | 0.2143      |
| Force             | 0.1116    | 0.0241     | -0.0301     | 1.0000  | 0.0928            | 0.1535     | 0.1402       | 0.0421     | 0.1506            | 0.1271      |
| Coating Viscosity | 0.0057    | 0.0223     | 0.0482      | 0.0928  | 1.0000            | -0.0331    | 0.0570       | 0.0287     | -0.0615           | 0.3194      |
| Inlet Temp        | 0.0217    | 0.0977     | 0.0059      | 0.1535  | -0.0331           | 1.0000     | 0.0761       | 0.0603     | 0.1476            | 0.0755      |
| Exhaust Temp      | 0.0810    | -0.1257    | -0.0184     | 0.1402  | 0.0570            | 0.0761     | 1.0000       | -0.0628    | 0.1977            | 0.1327      |
| Spray Rate        | -0.1381   | 0.2145     | -0.0745     | 0.0421  | 0.0287            | 0.0603     | -0.0628      | 1.0000     | 0.1380            | -0.3292     |
| Atomizer Pressure | -0.0775   | 0.1841     | 0.2632      | 0.1506  | -0.0615           | 0.1476     | 0.1977       | 0.1380     | 1.0000            | 0.1288      |
| Dissolution       | 0.3638    | 0.1598     | 0.2143      | 0.1271  | 0.3194            | 0.0755     | 0.1327       | -0.3292    | 0.1288            | 1.0000      |

The correlations are estimated by Row-wise method.

| Inlet Temp | Exhaust Temp | Spray Rate | Atomizer Pressure |
|------------|--------------|------------|-------------------|
| 107.9      | 70.5         | 404.6      | 61.0              |
| 107.5      | 70.8         | 407.4      | 60.6              |
| 106.6      | 69.2         | 399.3      | 59.1              |
| 106.1      | 68.8         | 403.7      | 58.8              |
| 108.3      | 69.4         | 396.7      | 59.6              |
| 106.3      | 69.1         | 404.7      | 60.4              |
| 106.1      | 69.7         | 399.3      | 58.4              |
| 107.6      | 70.0         | 398.5      | 61.6              |
| 107.2      | 71.4         | 404.0      | 61.1              |
| 106.8      | 70.4         | 394.9      | 59.5              |
| 105.2      | 69.7         | 403.3      | 60.9              |
| 105.5      | 71.9         | 395.4      | 59.7              |
| 106.6      | 69.3         | 397.7      | 57.5              |
| 106.6      | 70.1         | 388.4      | 58.2              |
| 105.3      | 68.7         | 391.6      | 58.5              |
| 106.7      | 68.9         | 418.5      | 58.2              |
| 108.1      | 69.4         | 402.3      | 60.5              |
| 105.9      | 69.1         | 396.8      | 60.2              |
| 106.5      | 69.1         | 397.2      | 62.3              |
| 105.5      | 70.6         | 408.5      | 58.8              |
| 107.6      | 70.6         | 401.1      | 60.4              |
| 106.9      | 69.6         | 404.8      | 60.0              |
| 106.7      | 69.6         | 407.4      | 61.4              |
| 107.4      | 72.0         | 403.6      | 60.1              |
| 105.0      | 70.3         | 390.9      | 58.5              |
| 107.7      | 71.2         | 400.9      | 61.4              |





Normal(106.903, 1.00218)

#### Compare Distributions

| Show                                | Distribution     | Number of Parameters | -2*LogLikelihood | AICc       |
|-------------------------------------|------------------|----------------------|------------------|------------|
| <input type="checkbox"/>            | SHASH            | 4                    | 248.053662       | 256.52425  |
| <input type="checkbox"/>            | Normal 2 Mixture | 5                    | 247.568499       | 258.282785 |
| <input type="checkbox"/>            | Gamma            | 2                    | 254.794139       | 258.93207  |
| <input type="checkbox"/>            | LogNormal        | 2                    | 254.79609        | 258.934021 |
| <input checked="" type="checkbox"/> | Normal           | 2                    | 254.801502       | 258.939433 |
| <input type="checkbox"/>            | Johnson Sl       | 3                    | 254.794544       | 261.073614 |
| <input type="checkbox"/>            | GLog             | 3                    | 254.795904       | 261.074974 |
| <input type="checkbox"/>            | Johnson Su       | 4                    | 254.795964       | 263.266553 |
| <input type="checkbox"/>            | Normal 3 Mixture | 8                    | 247.982016       | 265.759794 |
| <input type="checkbox"/>            | Weibull          | 2                    | 262.445187       | 266.583118 |
| <input type="checkbox"/>            | Extreme Value    | 2                    | 262.445187       | 266.583118 |
| <input type="checkbox"/>            | Exponential      | 1                    | 1020.94678       | 1022.99223 |

#### Fitted Normal

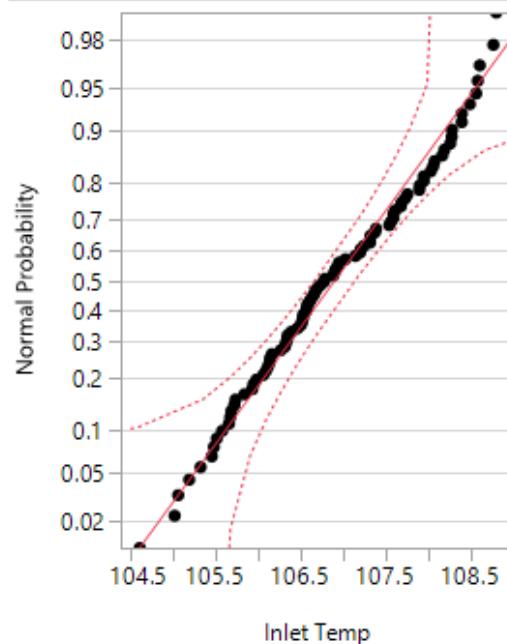
##### Parameter Estimates

| Type       | Parameter | Estimate  | Lower 95% | Upper 95% |
|------------|-----------|-----------|-----------|-----------|
| Location   | $\mu$     | 106.9035  | 106.6936  | 107.1134  |
| Dispersion | $\sigma$  | 1.0021833 | 0.8741169 | 1.1745646 |

##### Measure

|                  |           |
|------------------|-----------|
| -2*LogLikelihood | 254.8015  |
| AICc             | 258.93943 |
| BIC              | 263.80112 |

#### Diagnostic Plot

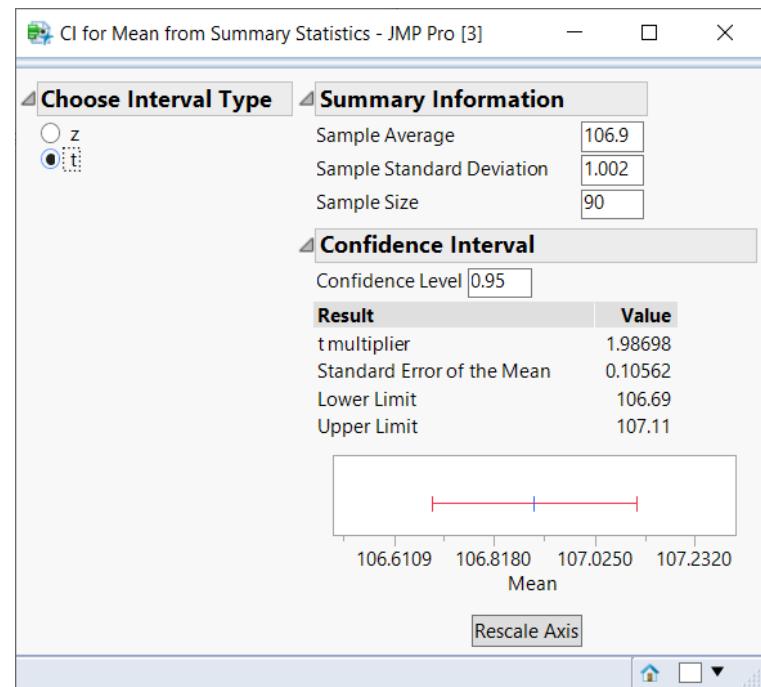


#### Goodness-of-Fit Test

Shapiro-Wilk W Test

| W        | Prob<W |
|----------|--------|
| 0.980334 | 0.1906 |

Note:  $H_0$  = The data is from the Normal distribution. Small p-values reject  $H_0$ .



## Confidence Intervals

| Parameter | Estimate | Lower CI | Upper CI | 1-Alpha |
|-----------|----------|----------|----------|---------|
| Mean      | 106.9035 | 106.6254 | 107.1816 | 0.990   |
| Std Dev   | 1.002183 | 0.838607 | 1.237307 | 0.990   |

Confidence Intervals

Enter (1-alpha) for confidence interval

Two-sided  
 One-sided lower limit  
 One-sided upper limit  
 Use known Sigma

OK Cancel Help

## Confidence Intervals

| Parameter | Estimate | Lower CI | Upper CI | 1-Alpha |
|-----------|----------|----------|----------|---------|
| Mean      | 106.9035 | 106.6936 | 107.1134 | 0.950   |
| Std Dev   | 1.002183 | 0.874117 | 1.174565 | 0.950   |

## Confidence Intervals

| Parameter | Estimate | Lower CI | Upper CI | 1-Alpha |
|-----------|----------|----------|----------|---------|
| Mean      | 106.9035 | 106.7279 | 107.0791 | 0.900   |
| Std Dev   | 1.002183 | 0.893286 | 1.14444  | 0.900   |

**Prediction Intervals**

Enter (1-alpha) for prediction interval

Enter number of future samples

Two-sided  
 One-sided lower limit  
 One-sided upper limit

**OK** **Cancel** **Help**

### Prediction Interval

| Parameter  | Future N | Lower PI | Upper PI | 1-Alpha |
|------------|----------|----------|----------|---------|
| Individual | 1        | 104.9012 | 108.9058 | 0.950   |
| Mean       | 1        | 104.9012 | 108.9058 | 0.950   |
| Std Dev    | 1        | .        | .        | 0.950   |

**Tolerance Intervals**

Computes an interval that contains at least the specified proportion of the population with (1-Alpha) confidence.

Specify confidence (1-Alpha):

Specify Proportion to cover:

Two-sided  
 One-sided lower limit  
 One-sided upper limit

Method

Assume Normal Distribution  
 Nonparametric

**OK** **Cancel** **Help**

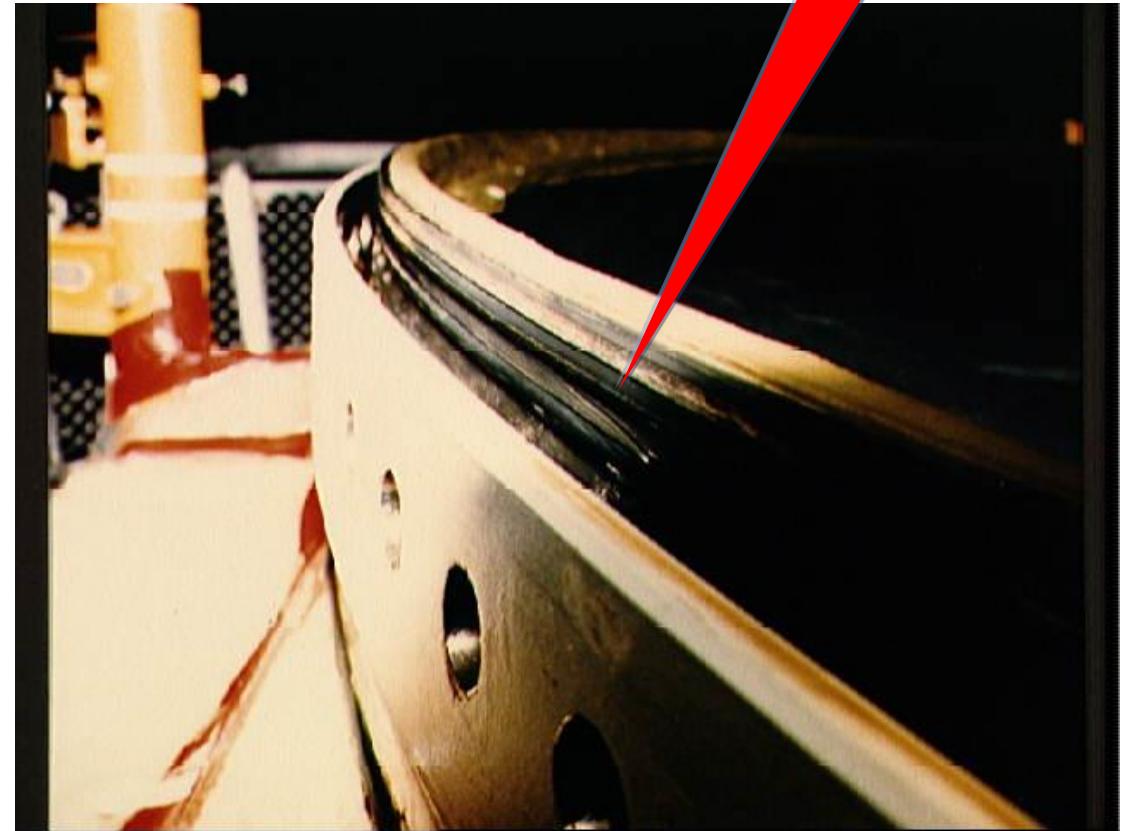
### Prediction Interval

| Parameter  | Future N | Lower PI | Upper PI | 1-Alpha |
|------------|----------|----------|----------|---------|
| Individual | 10       | 104.0025 | 109.8045 | 0.950   |
| Mean       | 10       | 106.2397 | 107.5673 | 0.950   |
| Std Dev    | 10       | 0.542555 | 1.506759 | 0.950   |

### Tolerance Intervals

| Proportion | Lower TI | Upper TI | 1-Alpha |
|------------|----------|----------|---------|
| 0.900      | 105.0095 | 108.7975 | 0.950   |

# The Challenger



O-ring

Kenett, R. and Thyregod, P. (2006) Aspects of statistical consulting not taught by academia, *Statistica Neerlandica*, special issue on Industrial Statistics, 30, 3, pp. 396-412.

# The Challenger

The US space shuttle Challenger was scheduled to take-off on January 28th, 1986, with seven crew members. Engineers from Morton Thiokol, manufacturers of the rocket motors, had been worried about problems with the O-ring seals. They feared that low temperatures greatly and adversely affected the ability of O-rings to create a seal on solid rocket booster joints.

On the night before the flight, the temperature predicted at launch time was 3° C, and the engineers expressed their concerns over the effect of the unseasonable cold weather on the O-rings and suggested to abort the flight.

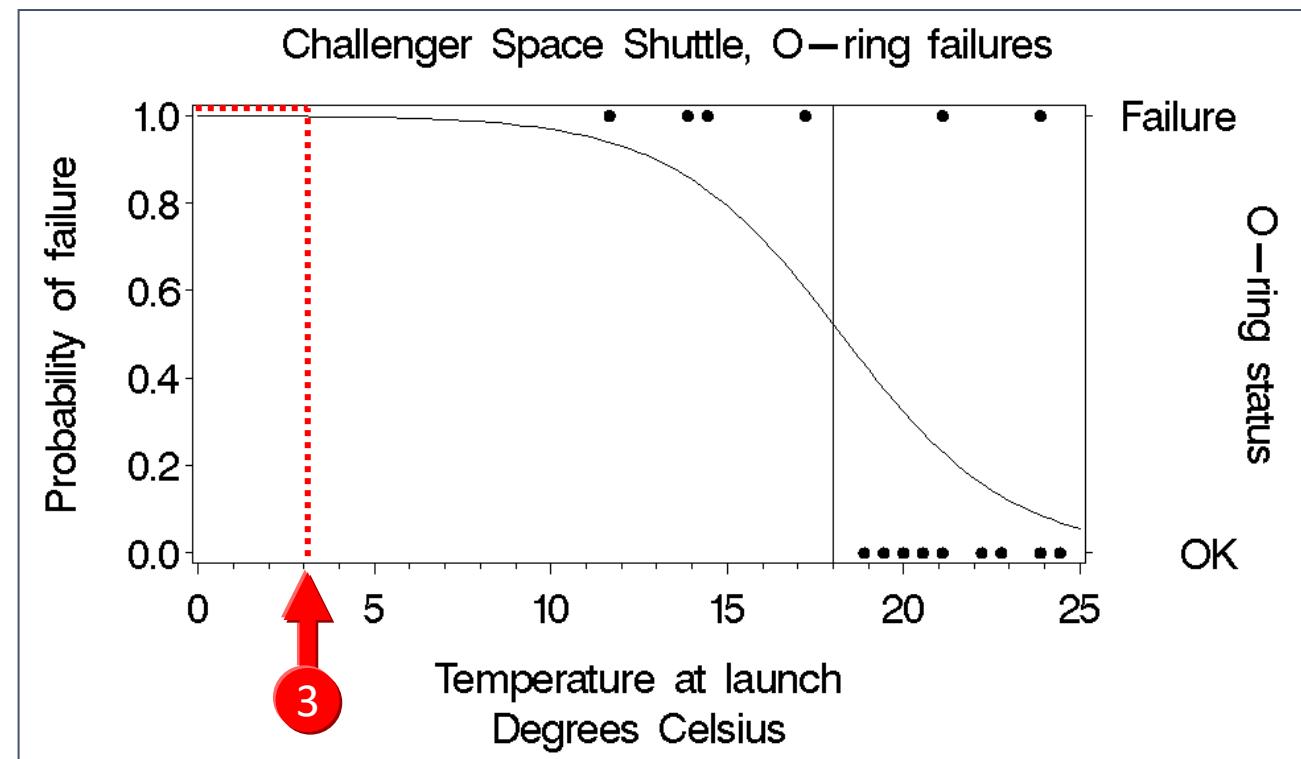
# The Challenger

A telephone conference was held between NASA engineers and managers and Thiokol engineers and managers.

With short notice, the Thiokol engineers presented their case via 13 telefaxed charts and their commentary and argument.

However, they failed to convince the managers that temperature was a factor in O-ring performance or damage, and it was decided to **go ahead** with the launching.

# Probability of O-ring Failure





- Tablet Production
- Reference Based on tablet produc
  - Control Chart and Distribution
  - Distribution
  - Multivariate
  - Oneway
  - Parallel Plot
  - Partition
  - Fit Model
  - Fit Model with Interactions
  - Fit Model with...action Profiles
  - Generalized Regression
  - Generalized R...duced Model
  - Fit Y by X of L... by Inlet Temp

## Columns (21/1)

- API Lot No
- API Particle Size
- Mill Time
- Screen Size
- Mag. Stearate Supplier
- Lactose Supplier
- Sugar Supplier
- Talc Supplier
- Blend Time
- Blend Speed
- Compressor
- Force
- Coating Supplier
- Coating Viscosity

- Rows
- |          |    |
|----------|----|
| All rows | 90 |
| Selected | 0  |
| Excluded | 0  |
| Hidden   | 0  |
| Labeled  | 0  |

20/1 Cols

|    | API Particle Size | Mill Time | Screen Size | Mag. Stearate Supplier | Lactose Supplier | Sugar Supplier | Talc Supplier | Blend Time | Blend Speed | Compressor | Force | Co    |
|----|-------------------|-----------|-------------|------------------------|------------------|----------------|---------------|------------|-------------|------------|-------|-------|
| 1  | Small             | 27        | 4           | Smith Ind              | James Ind        | Sour           | Rough         | 16.0       | 59.9        | Compress2  | 25.5  | Ma    |
| 2  | Small             | 11        | 5           | Jones Inc              | James Ind        | Sour           | Smooth        | 14.4       | 59.8        | Compress2  | 24.9  | Ma    |
| 3  | Small             | 20        | 4           | Jones Inc              | Bond Inc         | Sour           | Rough         | 14.5       | 60.8        | Compress2  | 25.5  | Dov   |
| 4  | Small             | 13        | 3           | Smith Ind              | Bond Inc         | Sweet          | Smooth        | 14.4       | 59.4        | Compress1  | 24.8  | Ma    |
| 5  | Small             | 13        | 5           | Smith Ind              | James Ind        | Sweet          | Smooth        | 16.1       | 59.9        | Compress2  | 25.3  | Dov   |
| 6  | Small             | 19        | 4           | Smith Ind              | Bond Inc         | Sweet          | Rough         | 12.9       | 59.4        | Compress2  | 24.6  | Ma    |
| 7  | Small             | 10        | 4           | Jones Inc              | Bond Inc         | Sweet          | Smooth        | 13.6       | 59.8        | Compress2  | 25.0  | Dov   |
| 8  | Small             | 24        | 4           | Jones Inc              | James Ind        | Sour           | Rough         | 15.1       | 61.1        | Compress2  | 24.9  | Ma    |
| 9  | Small             |           |             |                        | James Ind        | Sour           |               |            |             |            |       | ress1 |
| 10 | Small             |           |             |                        | Bond Inc         | Sweet          |               |            |             |            |       | ress2 |
| 11 | Small             |           |             |                        | Bond Inc         | Sweet          |               |            |             |            |       | ress1 |
| 12 | Small             |           |             |                        | James Ind        | Sour           |               |            |             |            |       | ress1 |
| 13 | Small             |           |             |                        | James Ind        | Sour           |               |            |             |            |       | ress2 |
| 14 | Small             |           |             |                        | James Ind        | Sour           |               |            |             |            |       | ress1 |
| 15 | Small             | 22        | 5           | Jones Inc              | James Ind        | Sweet          |               |            |             |            |       | ress2 |
| 16 | Small             | 7         | 3           | Jones Inc              | James Ind        | Sour           |               |            |             |            |       | ress2 |
| 17 | Small             | 6         | 3           | Jones Inc              | James Ind        | Sweet          |               |            |             |            |       | ress1 |
| 18 | Small             | 30        | 3           | Jones Inc              | Bond Inc         | Sweet          | Smooth        | 16.4       | 61.2        | Compress1  | 24.7  | Dov   |
| 19 | Small             | 29        | 3           | Smith Ind              | Bond Inc         | Sour           | Smooth        | 12.2       | 59.8        | Compress1  | 25.2  | Dov   |
| 20 | Small             | 7         | 5           | Jones Inc              | Bond Inc         | Sour           | Smooth        | 14.0       | 60.0        | Compress1  | 25.1  | Ma    |
| 21 | Small             | 25        | 5           | Jones Inc              | James Ind        | Sour           | Smooth        | 17.4       | 59.8        | Compress1  | 25.8  | Ma    |
| 22 | Small             | 13        | 5           | Jones Inc              | Bond Inc         | Sour           | Rough         | 15.7       | 58.7        | Compress2  | 25.0  | Dov   |
| 23 | Small             | 18        | 4           | Jones Inc              | James Ind        | Sweet          | Rough         | 17.4       | 61.2        | Compress2  | 24.9  | Ma    |
| 24 | Small             | 24        | 3           | Smith Ind              | Bond Inc         | Sweet          | Smooth        | 15.1       | 57.9        | Compress2  | 25.0  | Coa   |
| 25 | Small             | 13        | 5           | Smith Ind              | Bond Inc         | Sour           | Rough         | 15.2       | 61.5        | Compress2  | 24.7  | Coa   |
| 26 | Small             | 28        | 3           | Jones Inc              | Bond Inc         | Sour           | Smooth        | 15.9       | 61.1        | Compress2  | 25.1  | Coa   |
| 27 | Small             | 19        | 5           | Smith Ind              | James Ind        | Sweet          | Rough         | 15.8       | 60.2        | Compress2  | 24.7  | Ma    |
| 28 | Small             | 9         | 5           | Jones Inc              | James Ind        | Sour           | Rough         | 16.3       | 60.5        | Compress2  | 24.4  | Ma    |



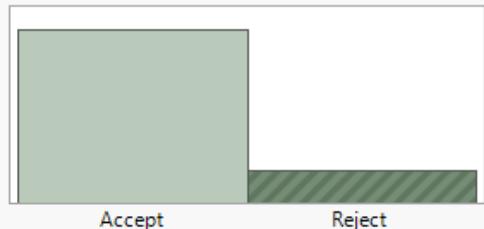
Tablet Production.jmp



## Logistic Fit of Lot Acceptance By Inlet Temp

### Distributions

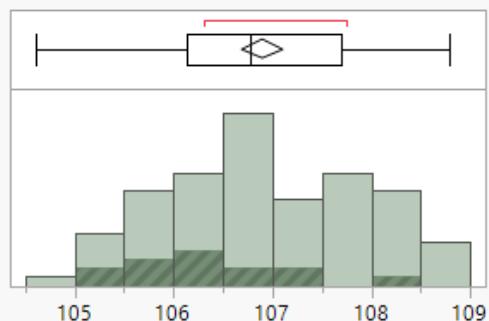
#### Lot Acceptance



#### Frequencies

| Level     | Count | Prob    |
|-----------|-------|---------|
| Accept    | 76    | 0.84444 |
| Reject    | 14    | 0.15556 |
| Total     | 90    | 1.00000 |
| N Missing | 0     |         |
| 2 Levels  |       |         |

#### Inlet Temp

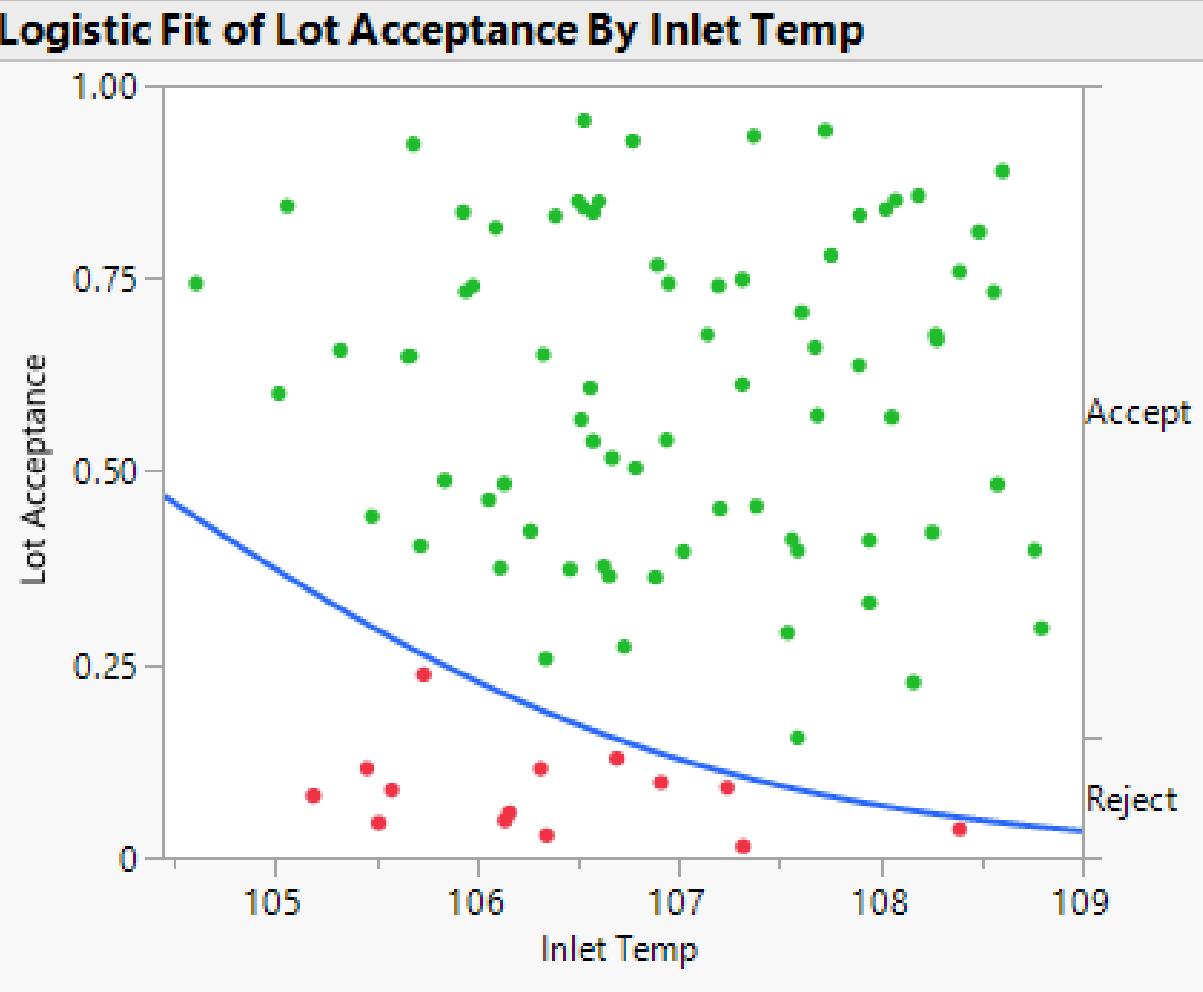


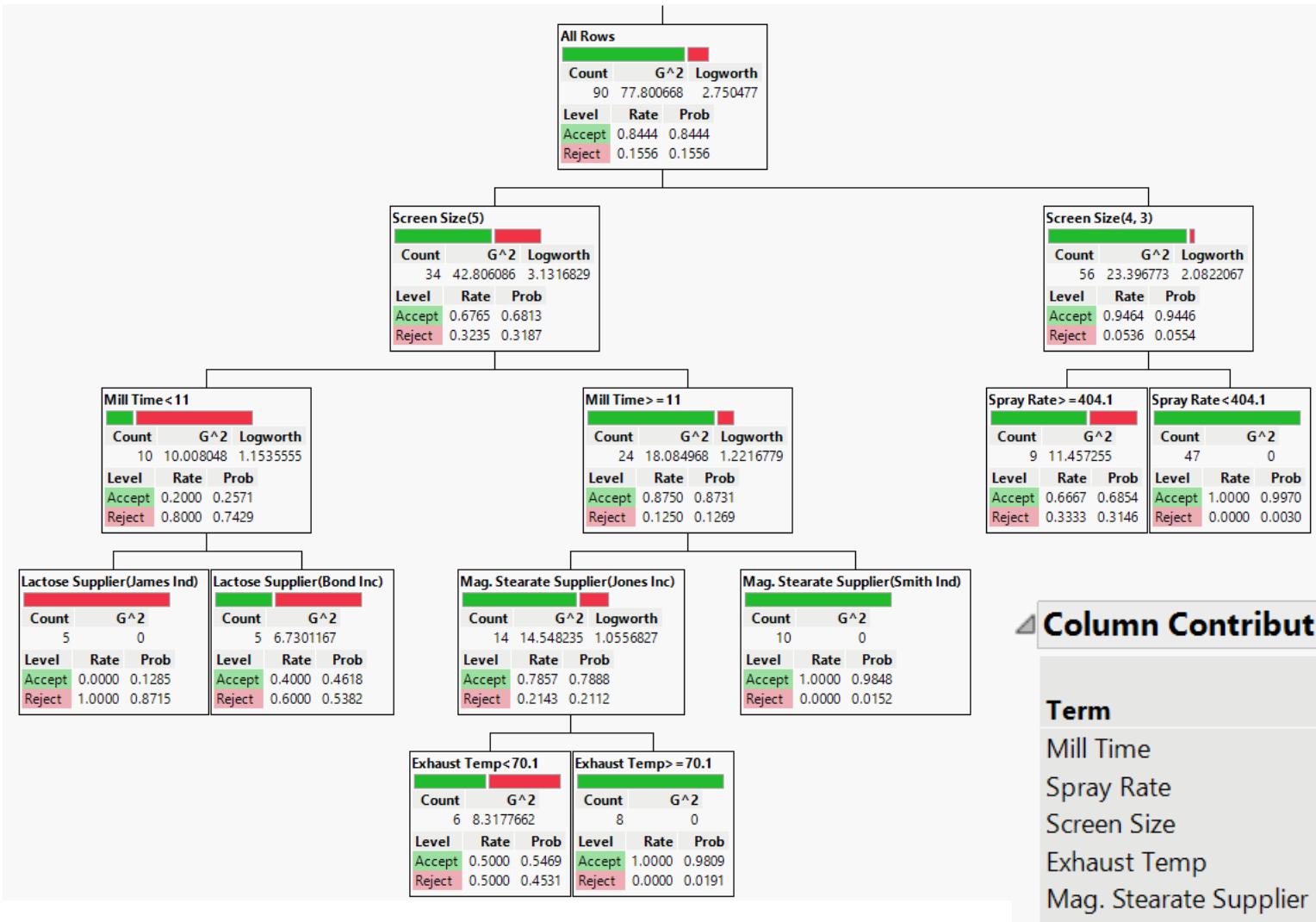
#### Quantiles

|        |          |              |
|--------|----------|--------------|
| 100.0% | maximum  | 108.79765542 |
| 99.5%  |          | 108.79765542 |
| 97.5%  |          | 108.7196325  |
| 90.0%  |          | 108.27844338 |
| 75.0%  | quartile | 107.69690533 |
| 50.0%  | median   | 106.77968341 |
| 25.0%  | quartile | 106.13953111 |
| 10.0%  |          | 105.58605819 |
| 2.5%   |          | 105.030777   |
| 0.5%   |          | 104.60908689 |
| 0.0%   | minimum  | 104.60908689 |

#### Summary Statistics

|                |           |
|----------------|-----------|
| Mean           | 106.9035  |
| Std Dev        | 1.0021833 |
| Std Err Mean   | 0.1056394 |
| Upper 95% Mean | 107.1134  |
| Lower 95% Mean | 106.6936  |
| N              | 90        |
| N Missing      | 0         |





**Fit Details**

| Measure                | Training Definition  |
|------------------------|--|
| Entropy RSquare        | 0.6282 $1 - \text{Loglike}(\text{model})/\text{Loglike}(0)$    |
| Generalized RSquare    | 0.7241 $(1 - L(0)/L(\text{model}))^{(2/n)}/(1 - L(0)^{(2/n)})$ |
| Mean - Log p           | 0.1607 $\sum -\text{Log}(p_{ij})/n$                            |
| RASE                   | 0.2315 $\sqrt{\sum (y_{ij} - \hat{p}_{ij})^2/n}$               |
| Mean Abs Dev           | 0.1166 $\sum  y_{ij} - \hat{p}_{ij} /n$                        |
| Misclassification Rate | 0.0889 $\sum (p_{ij} \neq p_{\text{Max}})/n$                   |
| N                      | 90 n   |

**Confusion Matrix**

|        |  | Training |                 |
|--------|--|----------|-----------------|
|        |  | Actual   | Predicted Count |
|        |  | Lot      | Acceptance      |
| Accept |  | 74       | 2               |
| Reject |  | 6        | 8               |

|        |  | Training |                |
|--------|--|----------|----------------|
|        |  | Actual   | Predicted Rate |
|        |  | Lot      | Acceptance     |
| Accept |  | 0.974    | 0.026          |
| Reject |  | 0.429    | 0.571          |

## Column Contributions

