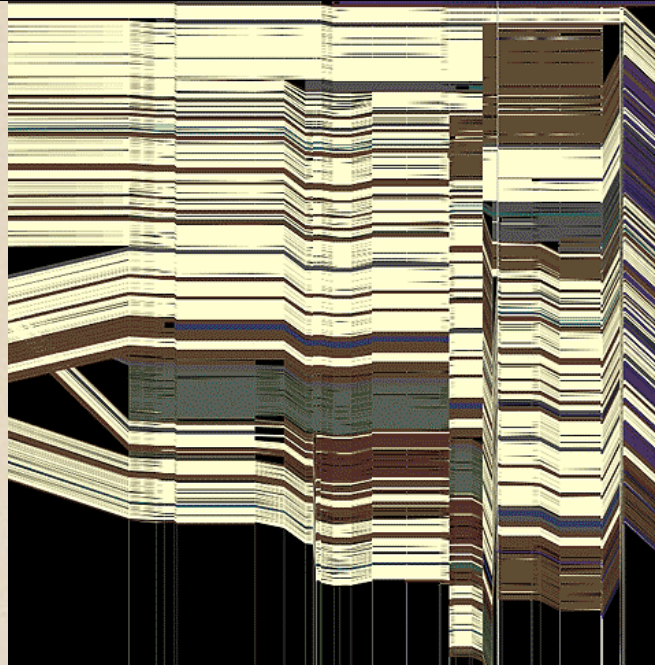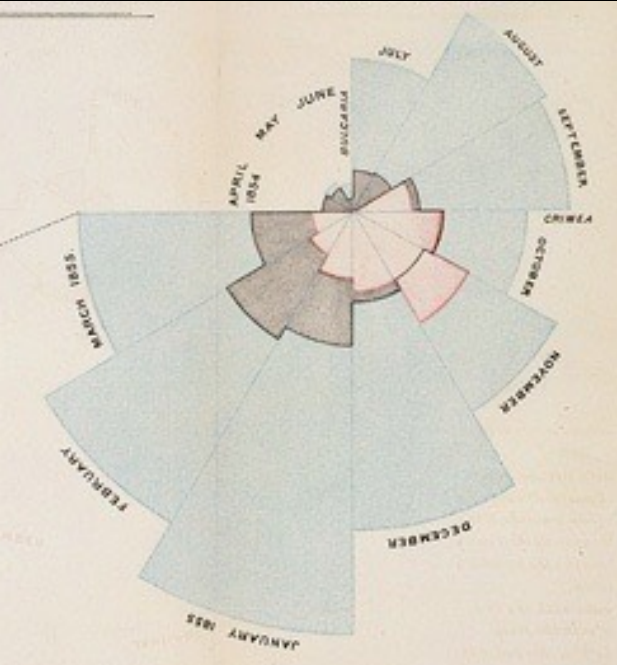**CSE 442** - Data Visualization

# Exploratory Data Analysis



Jeffrey Heer  University of Washington

# What was the **first** data visualization?
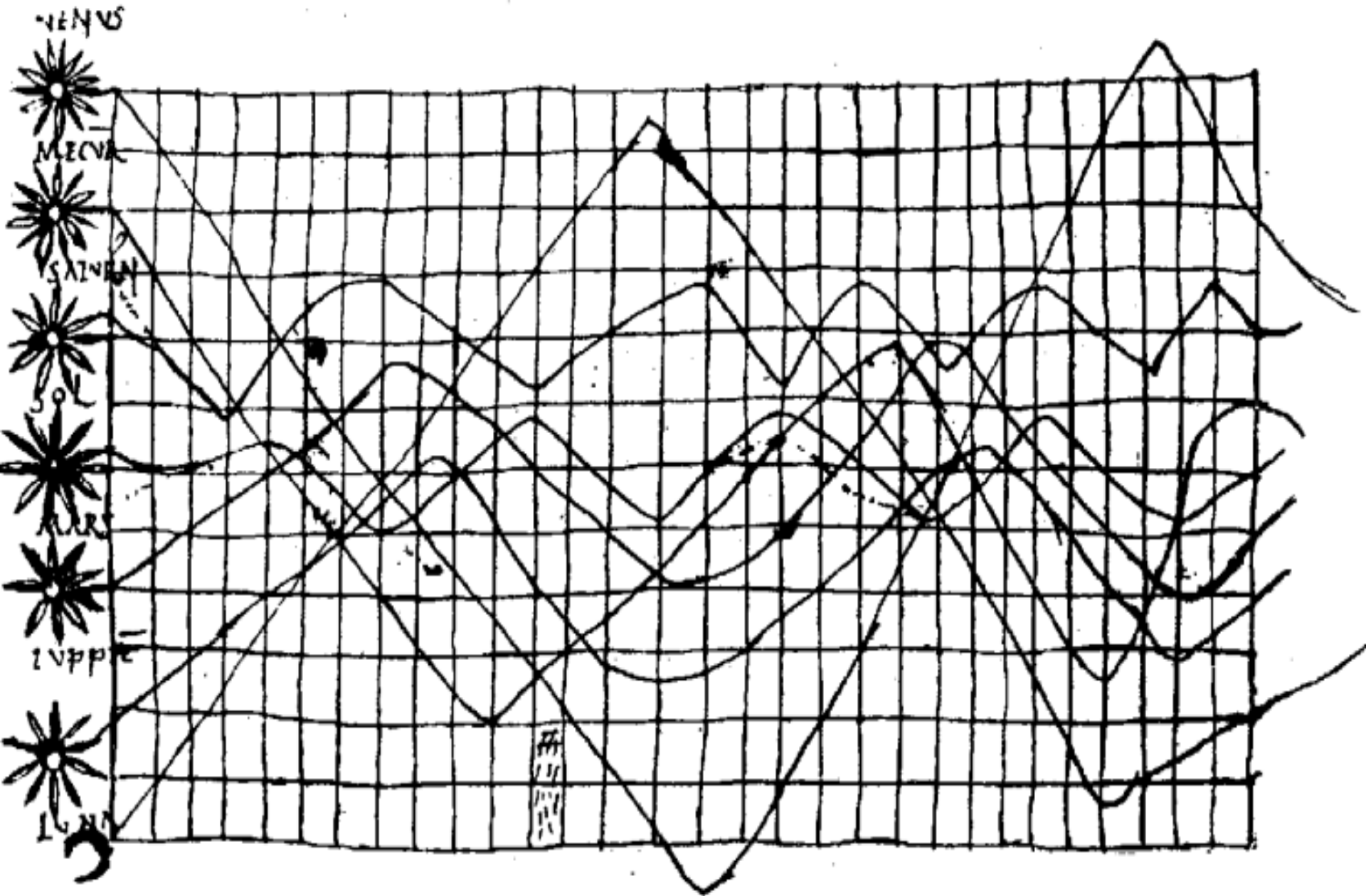
0 BC

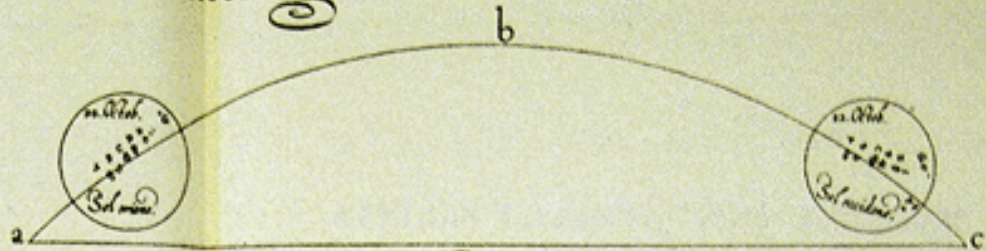~6200 BC Town Map of Catal Hyük, Konya Plain, Turkey            0 BC

~950 AD Position of Sun, Moon and Planets

Sunspots over time, Scheiner 1626

Longitudinal distance between Toledo and Rome, van Langren 1644

The Rate of Water Evaporation, Lambert 1765

The Rate of Water Evaporation, Lambert 1765

# The **Golden Age** of Data Visualization

1786   1900

Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.

The Commercial and Political Atlas, William Playfair 1786

Statistical Breviary, William Playfair 1801

1786    1826(?) Illiteracy in France, Pierre Charles Dupin

DIAGRAM of the CAUSES of MORTALITY IN THE ARMY in the EAST.

2. APRIL 1855 to MARCH 1856.

1. APRIL 1854 to MARCH 1855.

"to affect thro' the Eyes what we fail to convey to the public through their word-proof ears"

1786

1856 "Coxcomb" of Crimean War Deaths, Florence Nightingale

1864 British Coal Exports, Charles Minard

# Consommations approximatives de la Houille dans la Grande Bretagne de 1850 à 1864.

*Les abscisses représentent les années et les ordonnées les quantités annuelles de houille consommée.*

*Les couleurs indiquent les espèces de consommations. Les longueurs d'ordonnées comprises dans une couleur sont les quantités de houille consommées à raison de deux millimètres pour un million de tonnes.*

## Données admises pour former le Tableau ci-contre.

**Consommations.** —— **Sources des Renseignements.**

**Exportations.** — *Mineral statistics 1865 page 214 et Renseignements Parlementaires.*
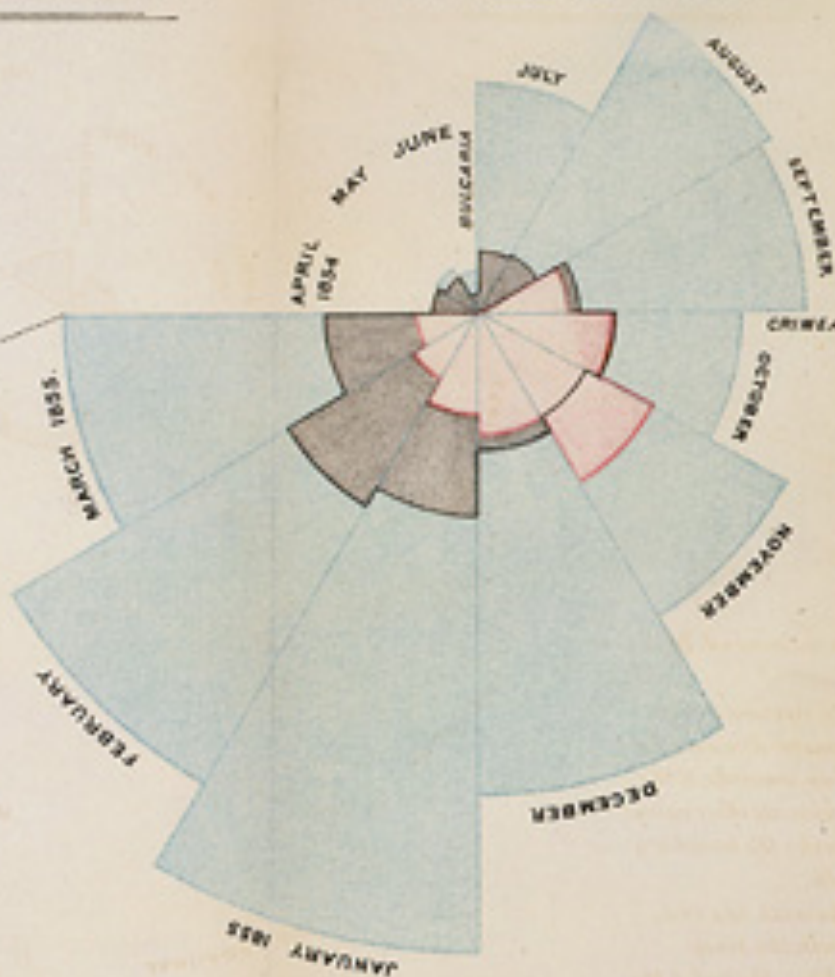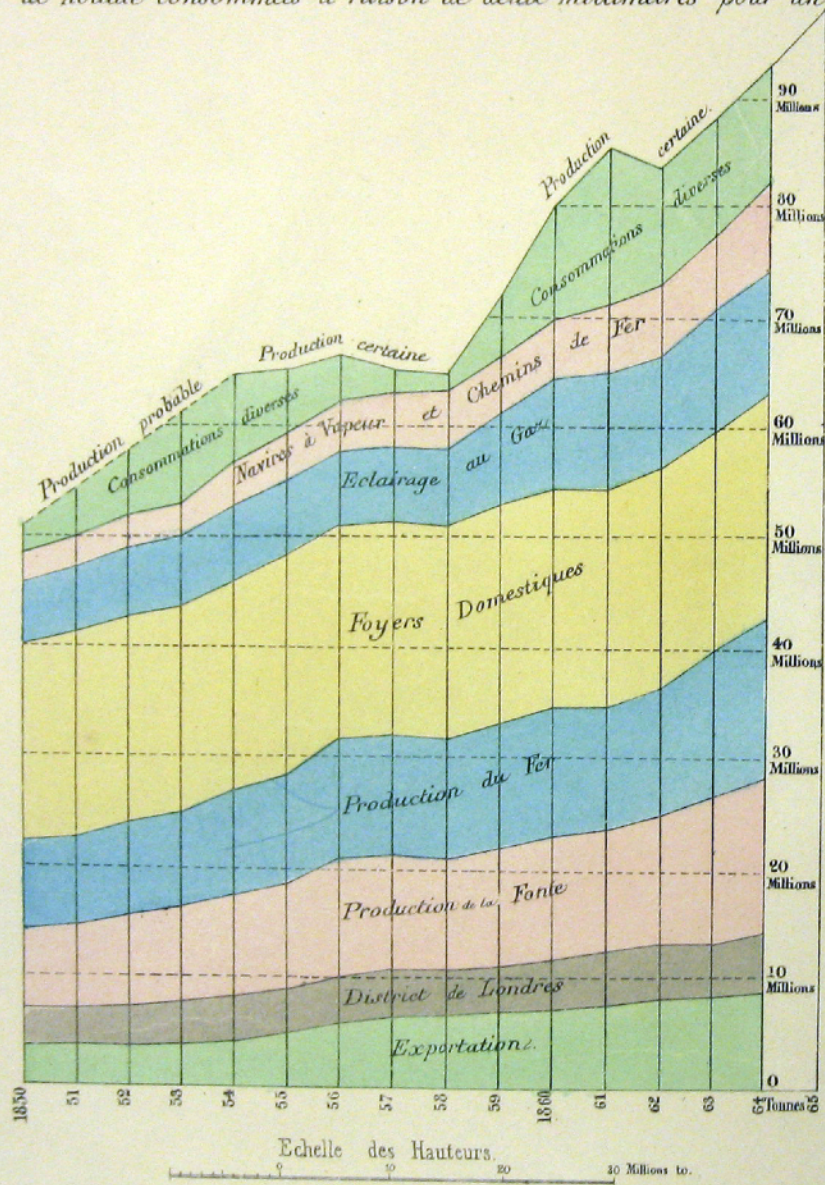
**District de Londres.** —— *id.* —————— *page 213*

**Produits de la Fonte.** —— *id* —————— *page 215 et pour les années avant 1855 calculée à raison de 3 t. de houille pour 1 t. de fonte, en admettant les quantités annuelles de fonte du Coal question page 192.*

**Production du fer** — *Mineral statistics* —— *page 215 et pour les années avant 1855* —— *calculée à raison de 3 t. 35 de houille pour 1 tonne de fonte convertie en fer; et admettant 2/10 es de la fonte produite convertis en fer.*

**Foyers domestiques :** —— *En y comprenant les petites manufactures. On l'estimait en 1848 à 19 millions de tonnes, (A) qu'on peut réduire à 18 millions to. pour les foyers seuls, mais qu'on peut porter à 20 millions pour la population de 1864.*
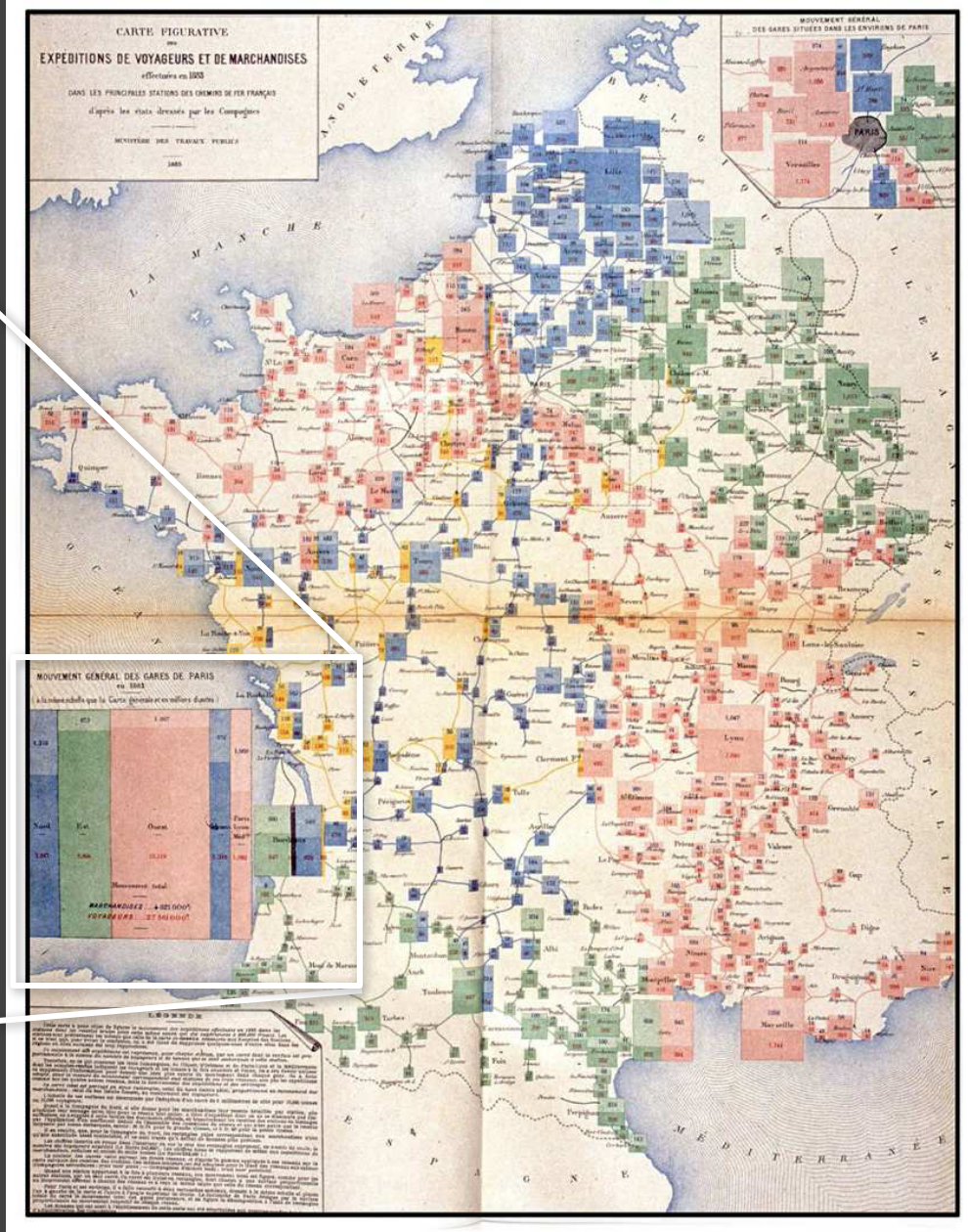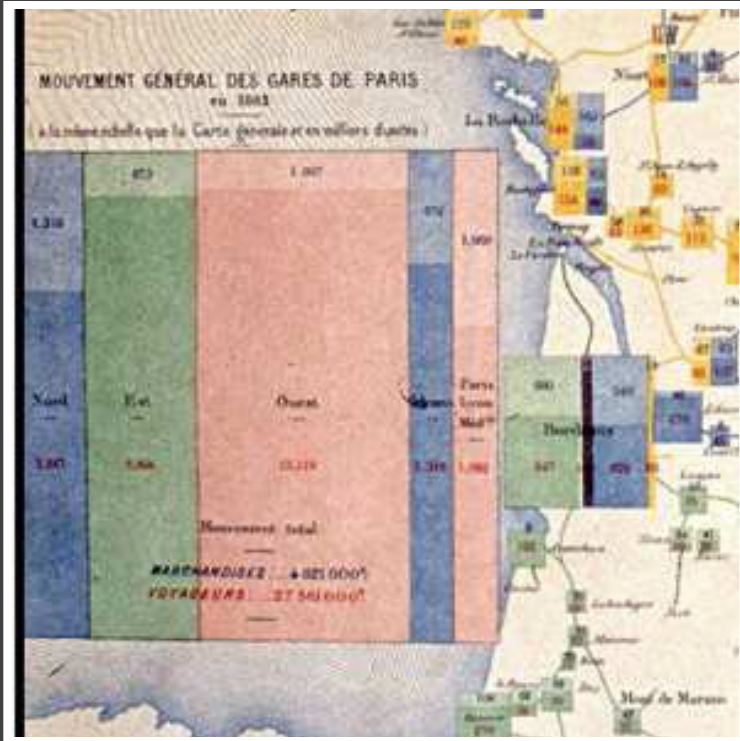
**Eclairage au Gaz.** — *Consommation estimée généralement du 1/3 au 1/8 de la production totale.*

**Exploitation des Chemins de Fer.** — *En supposant pour consommation totale 10 t. par Kilomètre parcouru par les trains d'après les renseignements parlementaires.*

**Navigation à vapeur.** — *Calculée à raison de 5 t. houille par cheval vapeur et par heure, le nombre de chevaux étant celui du Steam Vessels pour 1864, et les steamers étant supposés marcher la moitié de l'année;*

*Avant 1864 j'ai supposé les consommations proportionnelles aux tonnages annuels des steamers du statistical abstract et du Board of trade.*

**(A)** *Voir l'excellent article houille de M.r Lamé Fleury, Dictionnaire du Commerce Page III.*

### Chart labels

Production probable · Production certaine · Production certaine diverse

Consommations diverses · Navires à Vapeur et Chemins de Fer · Consommations diverses

Eclairage au Gaz

Foyers Domestiques

Production du Fer

Production de la Fonte

District de Londres

Exportation.

Axis (right): 90 Millions · 80 Millions · 70 Millions · 60 Millions · 50 Millions · 40 Millions · 30 Millions · 20 Millions · 10 Millions · 0 Tonnes

Axis (bottom): 1850 · 51 · 52 · 53 · 54 · 55 · 56 · 57 · 58 · 59 · 1860 · 61 · 62 · 63 · 64 Tonnes 65

Echelle des Hauteurs. 0 · 10 · 20 · 30 Millions to.

1786

1884 Rail Passengers and Freight from Paris

66. INTERSTATE MIGRATION—NUMBER OF NATIVE IMMIGRANTS AND NATIVE EMIGRANTS, BY STATES AND TERRITORIES: 1890.

1786

1890 Statistical Atlas of the Eleventh U.S. Census

# The Rise of Statistics

Rise of **formal methods** in statistics and social science – Fisher, Pearson, …

**Little innovation** in graphical methods

A period of **application and popularization**

Graphical methods enter textbooks, curricula, and **mainstream use**

1786                    1900                    1950

Data Analysis & Statistics, Tukey 1962

Four major influences act on data analysis today:

1. The formal theories of statistics.

2. Accelerating developments in computers and display devices.

3. The challenge, in many fields, of more and larger bodies of data.

4. The emphasis on quantification in a wider variety of disciplines.

The last few decades have seen the rise of formal theories of statistics, "legitimizing" variation by confining it by assumption to random sampling, often assumed to involve tightly specified distributions, and restoring the appearance of security by emphasizing narrowly optimized techniques and claiming to make statements with "known" probabilities of error.

While some of the influences of statistical theory on data analysis have been helpful, others have not.

**Exposure**, the effective laying open of the data to display the unanticipated, is to us a major portion of data analysis. Formal statistics has given almost no guidance to exposure; indeed, it is not clear how the **informality** and **flexibility** appropriate to the **exploratory character of exposure** can be fitted into any of the structures of formal statistics so far proposed.

Nothing - not the careful logic of mathematics, not statistical models and theories, not the awesome arithmetic power of modern computers - nothing can substitute here for the **flexibility of the informed human mind**.

Accordingly, both approaches and techniques need to be structured so as to **facilitate human involvement and intervention**.

| Set A | | Set B | | Set C | | Set D | |
|---|---|---|---|---|---|---|---|
| X | Y | X | Y | X | Y | X | Y |
| 10 | 8.04 | 10 | 9.14 | 10 | 7.46 | 8 | 6.58 |
| 8 | 6.95 | 8 | 8.14 | 8 | 6.77 | 8 | 5.76 |
| 13 | 7.58 | 13 | 8.74 | 13 | 12.74 | 8 | 7.71 |
| 9 | 8.81 | 9 | 8.77 | 9 | 7.11 | 8 | 8.84 |
| 11 | 8.33 | 11 | 9.26 | 11 | 7.81 | 8 | 8.47 |
| 14 | 9.96 | 14 | 8.1 | 14 | 8.84 | 8 | 7.04 |
| 6 | 7.24 | 6 | 6.13 | 6 | 6.08 | 8 | 5.25 |
| 4 | 4.26 | 4 | 3.1 | 4 | 5.39 | 19 | 12.5 |
| 12 | 10.84 | 12 | 9.11 | 12 | 8.15 | 8 | 5.56 |
| 7 | 4.82 | 7 | 7.26 | 7 | 6.42 | 8 | 7.91 |
| 5 | 5.68 | 5 | 4.74 | 5 | 5.73 | 8 | 6.89 |

**Summary Statistics**     **Linear Regression**

$u_X = 9.0$  $\sigma_X = 3.317$      $Y = 3 + 0.5\,X$

$u_Y = 7.5$  $\sigma_Y = 2.03$      $R^2 = 0.67$

[Anscombe 1973]

# Topics

**Exploratory Data Analysis**
Data Wrangling
Exploratory Analysis Examples
Polaris / Tableau

# Data Wrangling

I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any "analysis" at all.

Anonymous Data Scientist

[Kandel et al. '12]

**Big Data Borat**
@BigDataBorat

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

Reported crime in Alabama

| Year | Population | Property crime rate | | | Burglary rate | Larceny-theft rate | Motor vehicle theft rate |
|------|------------|------|------|------|---------------|--------------------|--------------------------|
| 2004 | 4525375 | 4029.3 | 987 | 2732.4 | 309.9 | | |
| 2005 | 4548327 | 3900 | 955.8 | 2656 | 289 | | |
| 2006 | 4599030 | 3937 | 968.9 | 2645.1 | 322.9 | | |
| 2007 | 4627851 | 3974.9 | 980.2 | 2687 | 307.7 | | |
| 2008 | 4661900 | 4081.9 | 1080.7 | 2712.6 | 288.6 | | |

Reported crime in Alaska

| Year | Population | Property crime rate | | | Burglary rate | Larceny-theft rate | Motor vehicle theft rate |
|------|------------|------|------|------|---------------|--------------------|--------------------------|
| 2004 | 657755 | 3370.9 | 573.6 | 2456.7 | 340.6 | | |
| 2005 | 663253 | 3615 | 622.8 | 2601 | 391 | | |
| 2006 | 670053 | 3582 | 615.2 | 2588.5 | 378.3 | | |
| 2007 | 683478 | 3373.9 | 538.9 | 2480 | 355.1 | | |
| 2008 | 686293 | 2928.3 | 470.9 | 2219.9 | 237.5 | | |

Reported crime in Arizona

| Year | Population | Property crime rate | | | Burglary rate | Larceny-theft rate | Motor vehicle theft rate |
|------|------------|------|------|------|---------------|--------------------|--------------------------|
| 2004 | 5739879 | 5073.3 | 991 | 3118.7 | 963.5 | | |
| 2005 | 5953007 | 4827 | 946.2 | 2958 | 922 | | |
| 2006 | 6166318 | 4741.6 | 953 | 2874.1 | 914.4 | | |
| 2007 | 6338755 | 4502.6 | 935.4 | 2780.5 | 786.7 | | |
| 2008 | 6500180 | 4087.3 | 894.2 | 2605.3 | 587.8 | | |

Reported crime in Arkansas

| Year | Population | Property crime rate | | | Burglary rate | Larceny-theft rate | Motor vehicle theft rate |
|------|------------|------|------|------|---------------|--------------------|--------------------------|
| 2004 | 2750000 | 4033.1 | 1096.4 | 2699.7 | 237 | | |
| 2005 | 2775708 | 4068 | 1085.1 | 2720 | 262 | | |
| 2006 | 2810872 | 4021.6 | 1154.4 | 2596.7 | 270.4 | | |
| 2007 | 2834797 | 3945.5 | 1124.4 | 2574.6 | 246.5 | | |
| 2008 | 2855390 | 3843.7 | 1182.7 | 2433.4 | 227.6 | | |

Reported crime in California

| Year | Population | Property crime rate | | | Burglary rate | Larceny-theft rate | Motor vehicle theft rate |
|------|------------|------|------|------|---------------|--------------------|--------------------------|
| 2004 | 35842038 | 3423.9 | 686.1 | 2033.1 | 704.8 | | |
| 2005 | 36154147 | 3321 | 692.9 | 1915 | 712 | | |
| 2006 | 36457549 | 3175.2 | 676.9 | 1831.5 | 666.8 | | |
| 2007 | 36553215 | 3032.6 | 648.4 | 1784.1 | 600.2 | | |
| 2008 | 36756666 | 2940.3 | 646.8 | 1769.8 | 523.8 | | |

Reported crime in Colorado

| Year | Population | Property crime rate | | | Burglary rate | Larceny-theft rate | Motor vehicle theft rate |
|------|------------|------|------|------|---------------|--------------------|--------------------------|
| 2004 | 4601821 | 3918.5 | 717.3 | 2679.5 | 521.6 | | |

# **Data**Wrangler



**Wrangler: Interactive Visual Specification
of Data Transformation Scripts**
Sean Kandel et al. *CHI'11*

cn16 - Transformer - Trifacta

⊕ Campaign Finance 2016 > ▦ cn16 ⌄                                    1 → ▦ → 0    **Generate Results**   👤 ❓

Grid    Columns    Full Dataset - 461.78kB ⌄    15 Columns    4,864 Rows    3 Data Types              🔍 Filter in grid    ↺ ↻ ⫶

| | ABC CAND_ID ⌄ | ABC CAND_NAME ⌄ | ABC CAND_PARTY_AFFILIATION ⌄ | 🕓 CAND_ELECTION_YEAR ⌄ | ABC CAND_OFFICE_STATE ⌄ | ABC CAND_OFFICE |
|---|---|---|---|---|---|---|
| | 4,864 Categories | 4,760 Categories | 76 Categories | 1986 - 2052 | 57 Categories | 3 Categories |
| · | H0AK00097 | COX, · JOHN · R. | REP | 2014 | AK | H |
| · | H0AL02087 | ROBY, · MARTHA | REP | 2016 | AL | H |
| · | H0AL02095 | JOHN, · ROBERT · E · JR | IND | 2016 | AL | H |
| · | H0AL05049 | CRAMER, · ROBERT · E · "BUD" · JR | DEM | 2008 | AL | H |
| · | H0AL05163 | BROOKS, · MO | REP | 2016 | AL | H |
| · | H0AL06088 | COOKE, · STANLEY · KYLE | REP | 2010 | AL | H |
| · | H0AL07086 | SEWELL, · TERRI · A. | DEM | 2016 | AL | H |
| · | H0AL07094 | HILLIARD, · EARL · FREDERICK · JR | DEM | 2010 | AL | H |
| · | H0AL07177 | CHAMBERLAIN, · DON | REP | 2012 | AL | H |
| · | H0AR01083 | CRAWFORD, · ERIC · ALAN · RICK | REP | 2016 | AR | H |
| · | H0AR01091 | GREGORY, · JAMES · CHRISTOPHER | DEM | 2010 | AR | H |
| · | H0AR01109 | CAUSEY, · CHAD | DEM | 2010 | AR | H |
| · | H0AR01125 | SMITH, · PRINCELLA · D | REP | 2010 | AR | H |
| · | H0AR02107 | GRIFFIN, · JOHN · TIMOTHY | REP | 2014 | AR | H |
| · | H0AR02131 | ELLIOTT, · JOYCE · ANN | DEM | 2010 | AR | H |
| · | H0AR03022 | SKOCH, · BERNARD · KURT · 'BERNIE' | REP | 2010 | AR | H |
| · | H0AR03030 | WHITAKER, · DAVID · JEFFREY | DEM | 2010 | AR | H |
| · | H0AR03055 | WOMACK, · STEVE | REP | 2016 | AR | H |
| · | H0AS00018 | FALEOMAVAEGA, · ENI | DEM | 2014 | AS | H |
| · | H0AZ01184 | FLAKE, · JEFF · MR. | REP | 2012 | AZ | H |
| · | H0AZ01259 | GOSAR, · PAUL · ANTHONY | REP | 2016 | AZ | H |
| · | H0AZ01283 | MEHTA, · STEVE | REP | 2010 | AZ | H |
| · | H0AZ01325 | TOBIN, · ANDY · HON. | REP | 2014 | AZ | H |
| · | H0AZ01333 | GRESSLEY, · FORREST · DAYL | REP | 2010 | AZ | H |
| · | H0AZ03321 | PARKER, · VERNON | REP | 2014 | AZ | H |

**New Step**  Switch to editor                                          Cancel    **Add to Recipe**

**Choose a transformation**

[ Choose transformation ]

**TRIFACTA**

cn16 - Transformer - Trifacta

Campaign Finance 2016 > cn16

1 → □ → 0    Generate Results

| Grid | Columns | Full Dataset - 461.78kB | 15 Columns | 4,864 Rows | 3 Data Types |

Filter in grid

| ABC CAND_ID | ABC CAND_NAME | ABC CAND_PARTY_AFFILIATION | ⊙ CAND_ELECTION_YEAR | ABC CAND_OFFICE_STATE | ABC CAND_OFFICE |
|---|---|---|---|---|---|
| 4,864 Categories | 4,760 Categories | | 1986 - 2052 | 57 Categories | 3 Categories |
| H0AK00097 | COX, JOHN R. | Rename | 2014 | AK | H |
| H0AL02087 | ROBY, MARTHA | Change type > | 2016 | AL | H |
| H0AL02095 | JOHN, ROBERT E JR | Edit column > | 2016 | AL | H |
| H0AL05049 | CRAMER, ROBERT E "BUD" J | Column Details | 2008 | AL | H |
| H0AL05163 | BROOKS, MO | | 2016 | AL | H |
| H0AL06088 | COOKE, STANLEY KYLE | Find > | 2010 | AL | H |
| H0AL07086 | SEWELL, TERRI A. | Format > | 2016 | AL | H |
| H0AL07094 | HILLIARD, EARL FREDERICK | Filter > | 2010 | AL | H |
| H0AL07177 | CHAMBERLAIN, DON | Clean > | 2012 | AL | H |
| H0AR01083 | CRAWFORD, ERIC ALAN RICK | Formula > | 2016 | AR | H |
| H0AR01091 | GREGORY, JAMES CHRISTOPH | Aggregate > | 2010 | AR | H |
| H0AR01109 | CAUSEY, CHAD | | 2010 | AR | H |
| H0AR01125 | SMITH, PRINCELLA D | Restructure > | 2010 | AR | H |
| H0AR02107 | GRIFFIN, JOHN TIMOTHY | | 2014 | AR | H |
| H0AR02131 | ELLIOTT, JOYCE ANN | Lookup... | 2010 | AR | H |
| H0AR03022 | SKOCH, BERNARD KURT BER | | 2010 | AR | H |
| H0AR03030 | WHITAKER, DAVID JEFFREY | Drop | 2010 | AR | H |
| H0AR03055 | WOMACK, STEVE | REP | 2016 | AR | H |
| H0AS00018 | FALEOMAVAEGA, ENI | DEM | 2014 | AS | H |
| H0AZ01184 | FLAKE, JEFF MR. | REP | 2012 | AZ | H |
| H0AZ01259 | GOSAR, PAUL ANTHONY | REP | 2016 | AZ | H |
| H0AZ01283 | MEHTA, STEVE | REP | 2010 | AZ | H |
| H0AZ01325 | TOBIN, ANDY HON. | REP | 2014 | AZ | H |
| H0AZ01333 | GRESSLEY, FORREST DAYL | REP | 2010 | AZ | H |
| H0AZ03321 | PARKER, VERNON | REP | 2014 | AZ | H |

**New Step** Switch to editor

Cancel    Add to Recipe

**Choose a transformation**

Choose transformation

# Data Wrangling

One often needs to manipulate data prior to analysis. Tasks include reformatting, cleaning, quality assessment, and integration.

*Approaches include:*
Manual manipulation in spreadsheets
Custom code (e.g., dplyr in R, Pandas in Python)
Trifacta Wrangler  http://www.trifacta.com/products/wrangler/
Open Refine  http://openrefine.org/

# Data Quality

"The first sign that a visualization is good is that it shows you a problem in your data…

…every successful visualization that I've been involved with has had this stage where you realize, "Oh my God, this data is not what I thought it would be!" So already, you've discovered something."

Martin Wattenberg

# Graph Viewer

**Roll-up by:**

All

**Visualization:**

Node-Link

**Sort by:**

None

**Edge centrality filters:**

☐ Images
☑ Animate

# Graph Viewer

**Roll-up by:**

All

**Visualization:**

Matrix

**Sort by:**

Linkage

**Edge centrality filters:**

# Graph Viewer

**Roll-up by:**

All

**Visualization:**

Matrix

**Sort by:**

None

**Edge centrality filters:**

# Visualize Friends by School?

| School | Count |
|---|---|
| Berkeley | ||||||||||||||||||||||||| |
| Cornell | |||| |
| Harvard | ||||||||| |
| Harvard University | ||||||| |
| Stanford | ||||||||||||||||||| |
| Stanford University | |||||||||| |
| UC Berkeley | ||||||||||||||||||| |
| UC Davis | |||||||||| |
| University of California at Berkeley | |||||||||||||| |
| University of California, Berkeley | ||||||||||||||||||| |
| University of California, Davis | ||| |

# Data Quality Hurdles

| | |
|---|---|
| Missing Data | no measurements, redacted, …? |
| Erroneous Values | misspelling, outliers, …? |
| Type Conversion | e.g., zip code to lat-lon |
| Entity Resolution | diff. values for the same thing? |
| Data Integration | effort/errors when combining data |

*LESSON*: Anticipate problems with your data. Many research problems around these issues!

# Analysis Example: Motion Pictures Data

# Motion Pictures Data

| Title | String (N) |
| IMDB Rating | Number (Q) |
| Rotten Tomatoes Rating | Number (Q) |
| MPAA Rating | String (O) |
| Release Date | Date (T) |

IMDB Rating (bin)

Rotten Tomatoes Rating (bin)

# Lesson: Exercise Skepticism

Check **data quality** and your **assumptions**.

Start with **univariate summaries**, then start to consider **relationships among variables**.

**Avoid premature fixation!**

# Analysis Example: Antibiotic Effectiveness

## Data Set: Antibiotic Effectiveness

| | |
|---|---|
| Genus of Bacteria | String (N) |
| Species of Bacteria | String (N) |
| Antibiotic Applied | String (N) |
| Gram-Staining? | Pos / Neg (N) |
| Min. Inhibitory Concent. (g) | Number (Q) |

Collected prior to 1951.

# What questions might we ask?

| Table 1: Burtin's data. | Antibiotic | | | |
| --- | --- | --- | --- | --- |
| Bacteria | Penicillin | Streptomycin | Neomycin | Gram Staining |
| Aerobacter *aerogenes* | 870 | 1 | 1.6 | negative |
| Brucella *abortus* | 1 | 2 | 0.02 | negative |
| Brucella *anthracis* | 0.001 | 0.01 | 0.007 | positive |
| Diplococcus *pneumoniae* | 0.005 | 11 | 10 | positive |
| Escherichia *coli* | 100 | 0.4 | 0.1 | negative |
| Klebsiella *pneumoniae* | 850 | 1.2 | 1 | negative |
| Mycobacterium *tuberculosis* | 800 | 5 | 2 | negative |
| Proteus *vulgaris* | 3 | 0.1 | 0.1 | negative |
| Pseudomonas *aeruginosa* | 850 | 2 | 0.4 | negative |
| Salmonella (Eberthella) *typhosa* | 1 | 0.4 | 0.008 | negative |
| Salmonella *schottmuelleri* | 10 | 0.8 | 0.09 | negative |
| Staphylococcus *albus* | 0.007 | 0.1 | 0.001 | positive |
| Staphylococcus *aureus* | 0.03 | 0.03 | 0.001 | positive |
| Streptococcus *fecalis* | 1 | 1 | 0.1 | positive |
| Streptococcus *hemolyticus* | 0.001 | 14 | 10 | positive |
| Streptococcus *viridans* | 0.005 | 10 | 40 | positive |

# How do the drugs compare?



| Bacteria | Penicillin | Antibiotic Streptomycin | Neomycin | Gram stain |
|---|---|---|---|---|
| Aerobacter aerogenes | 870 | 1 | 1.6 | – |
| Brucella abortus | 1 | 2 | 0.02 | – |
| Bacillus anthracis | 0.001 | 0.01 | 0.007 | + |
| Diplococcus pneumoniae | 0.005 | 11 | 10 | + |
| Escherichia coli | 100 | 0.4 | 0.1 | – |
| Klebsiella pneumoniae | 850 | 1.2 | 1 | – |
| Mycobacterium tuberculosis | 800 | 5 | 2 | – |
| Proteus vulgaris | 3 | 0.1 | 0.1 | – |
| Pseudomonas aeruginosa | 850 | 2 | 0.4 | – |
| Salmonella (Eberthella) typhosa | 1 | 0.4 | 0.008 | – |
| Salmonella schottmuelleri | 10 | 0.8 | 0.09 | – |
| Staphylococcus albus | 0.007 | 0.1 | 0.001 | + |
| Staphylococcus aureus | 0.03 | 0.03 | 0.001 | + |
| Streptococcus fecalis | 1 | 1 | 0.1 | + |
| Streptococcus hemolyticus | 0.001 | 14 | 10 | + |
| Streptococcus viridans | 0.005 | 10 | 40 | + |

Original graphic by Will Burtin, 1951

# How do the drugs compare?



| Bacteria | Penicillin | Antibiotic Streptomycin | Neomycin | Gram stain |
|---|---|---|---|---|
| Aerobacter aerogenes | 870 | 1 | 1.6 | − |
| Brucella abortus | 1 | 2 | 0.02 | − |
| Bacillus anthracis | 0.001 | 0.01 | 0.007 | + |
| Diplococcus pneumoniae | 0.005 | 11 | 10 | + |
| Escherichia coli | 100 | 0.4 | 0.1 | − |
| Klebsiella pneumoniae | 850 | 1.2 | 1 | − |
| Mycobacterium tuberculosis | 800 | 5 | 2 | − |
| Proteus vulgaris | 3 | 0.1 | 0.1 | − |
| Pseudomonas aeruginosa | 850 | 2 | 0.4 | − |
| Salmonella (Eberthella) typhosa | 1 | 0.4 | 0.008 | − |
| Salmonella schottmuelleri | 10 | 0.8 | 0.09 | − |
| Staphylococcus albus | 0.007 | 0.1 | 0.001 | + |
| Staphylococcus aureus | 0.03 | 0.03 | 0.001 | + |
| Streptococcus fecalis | 1 | 1 | 0.1 | + |
| Streptococcus hemolyticus | 0.001 | 14 | 10 | + |
| Streptococcus viridans | 0.005 | 10 | 40 | + |

Radius: 1 / log(MIC)
Bar Color: Antibiotic
Background Color: Gram Staining

# How do the drugs compare?



Mike Bostock
Stanford CS448B, Winter 2009

# How do the drugs compare?



X-axis: Antibiotic | log(MIC)
Y-axis: Gram-Staining | Species
Color: Most-Effective?
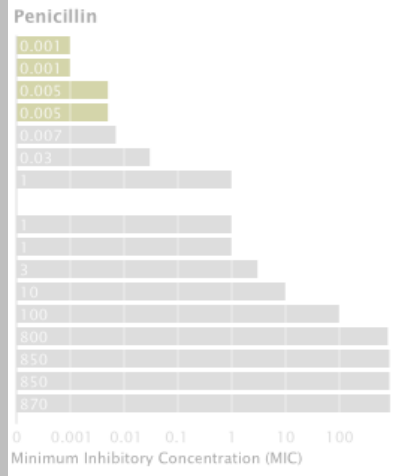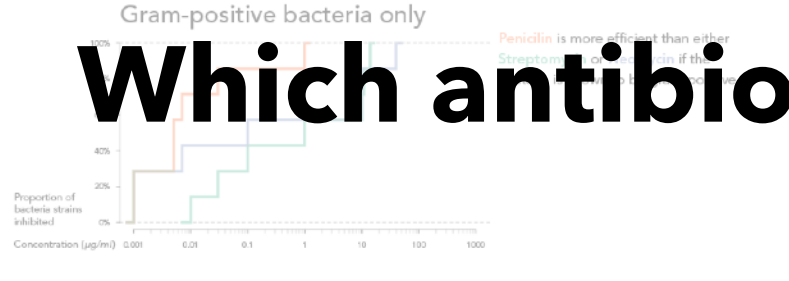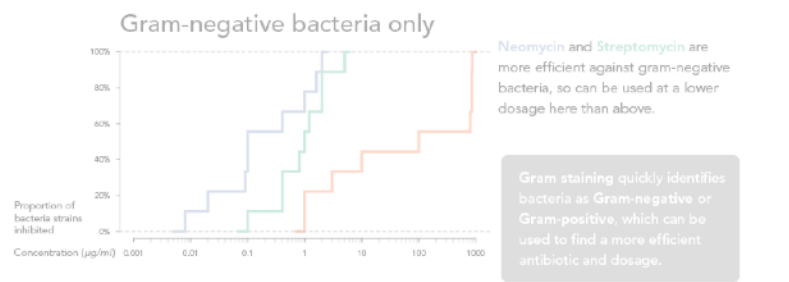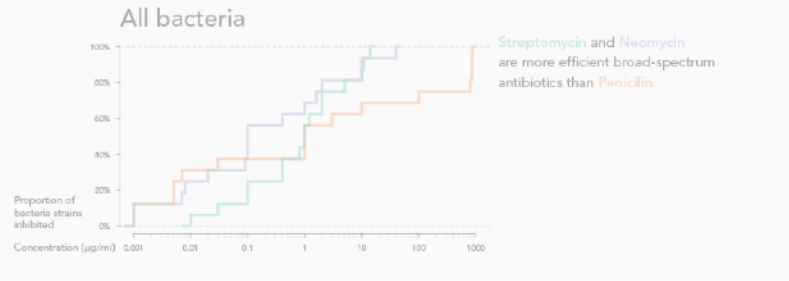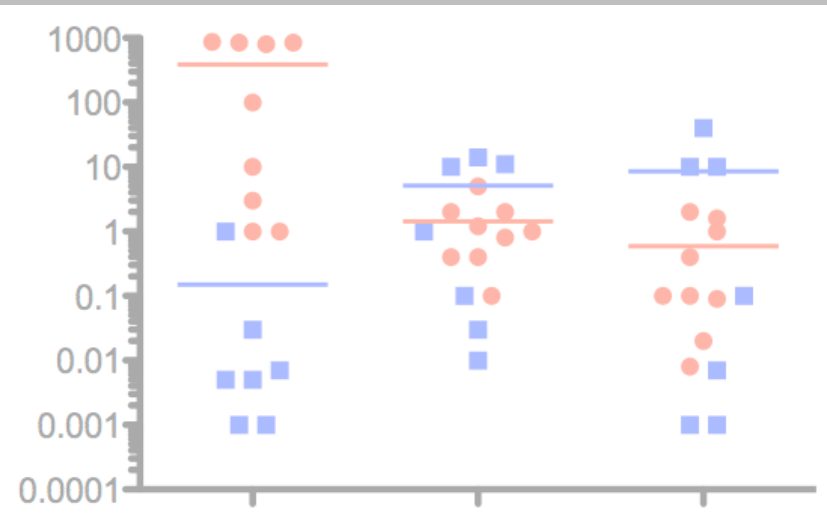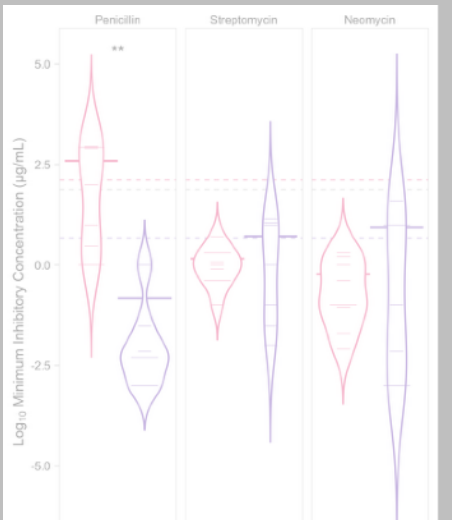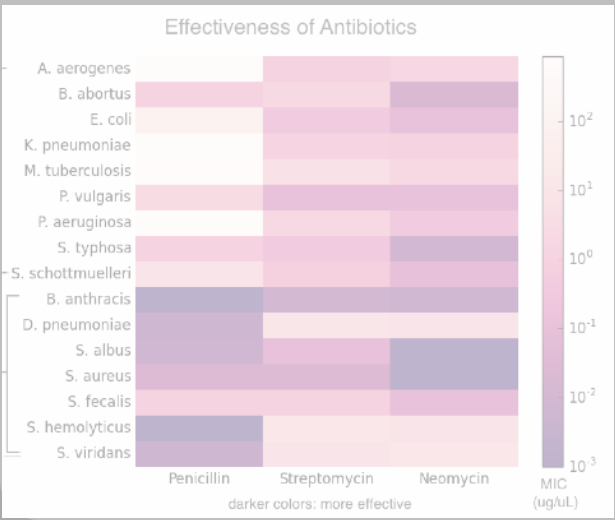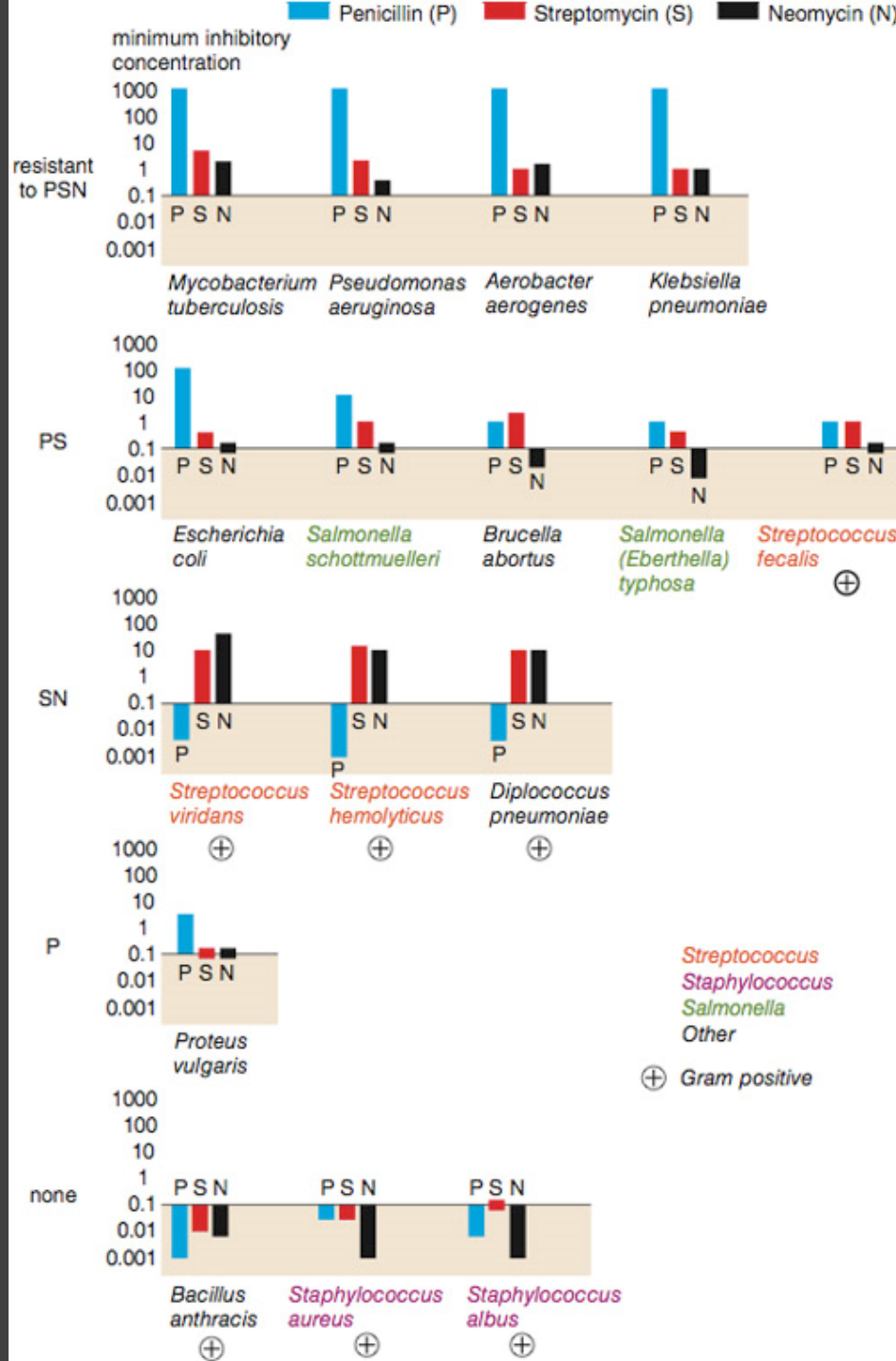
minimum inhibitory concentration of antibiotics

bowen li
cs448b

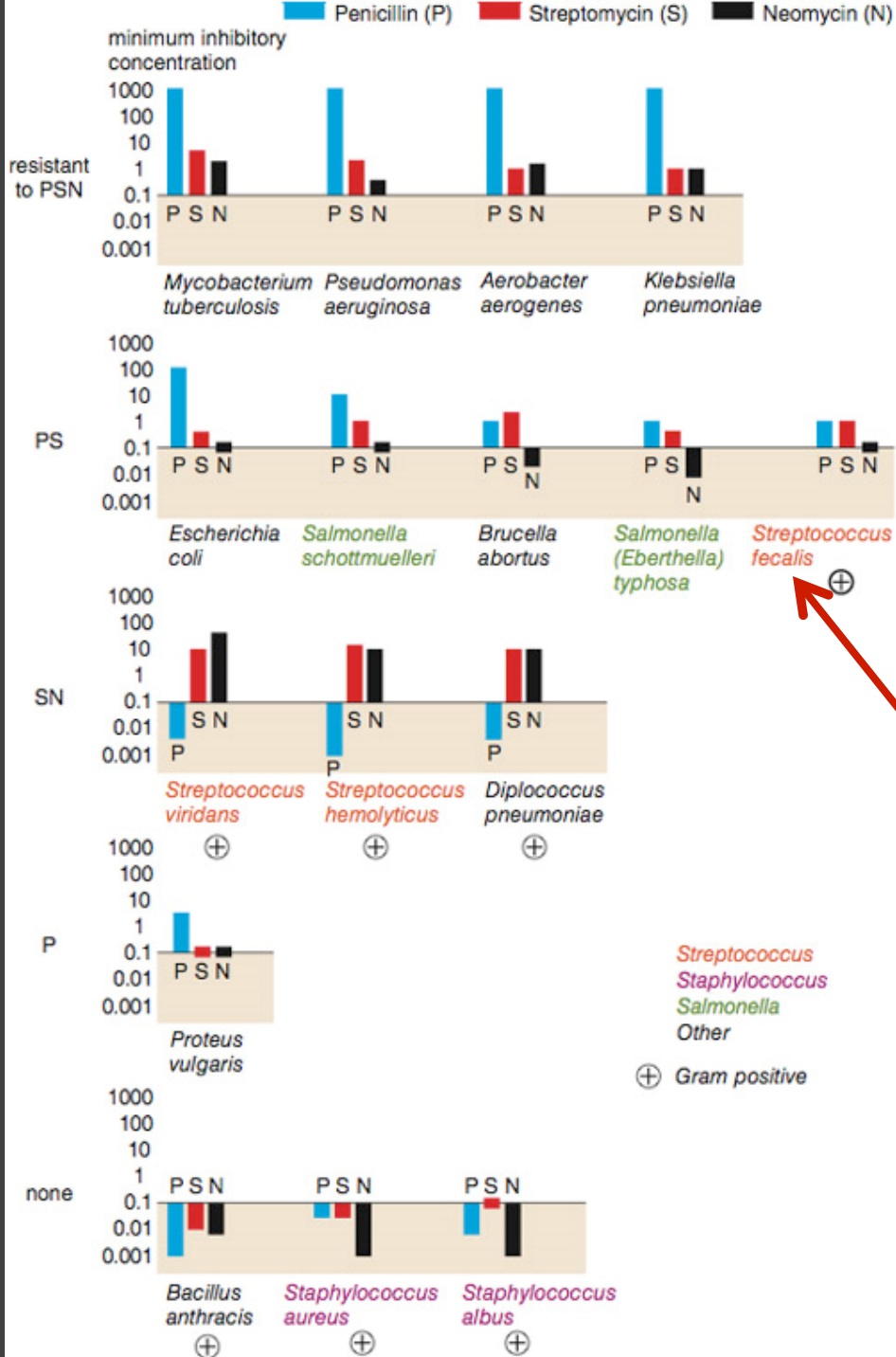# Which antibiotic should one use?

# Do the bacteria group by antibiotic resistance?

# Do the bacteria group by antibiotic resistance?

Wainer & Lysen
*American Scientist*, 2009

# Do the bacteria group by antibiotic resistance?

Not a streptococcus!
(realized ~30 yrs later)
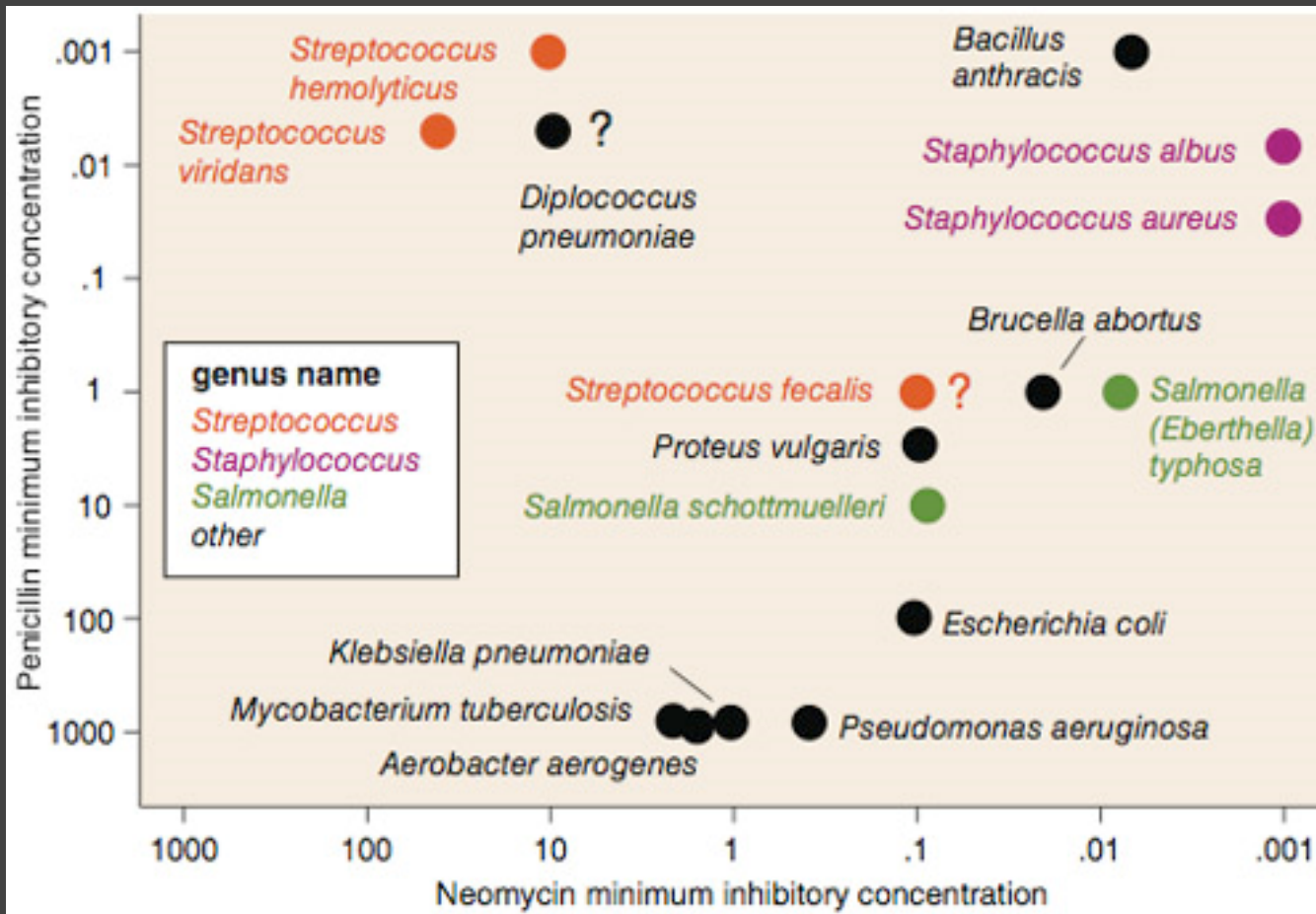
Wainer & Lysen
*American Scientist*, 2009

# Do the bacteria group by antibiotic resistance?

Not a streptococcus!
(realized ~30 yrs later)

Really a streptococcus!
(realized ~20 yrs later)

Do the bacteria group by resistance?
Do different drugs correlate?

**Do the bacteria group by resistance?**
**Do different drugs correlate?**

Wainer & Lysen
*American Scientist,* 2009

# Lesson: Iterative Exploration

**Exploratory Process**
1  Construct graphics to address questions
2  Inspect "answer" and assess new questions
3  Repeat…

**Transform data** appropriately (e.g., invert, log)

**Show data variation, not design variation** [Tufte]

# Administrivia

# A2: Exploratory Data Analysis

Use visualization software to form & answer questions

**First steps:**

Step 1: Pick domain & data
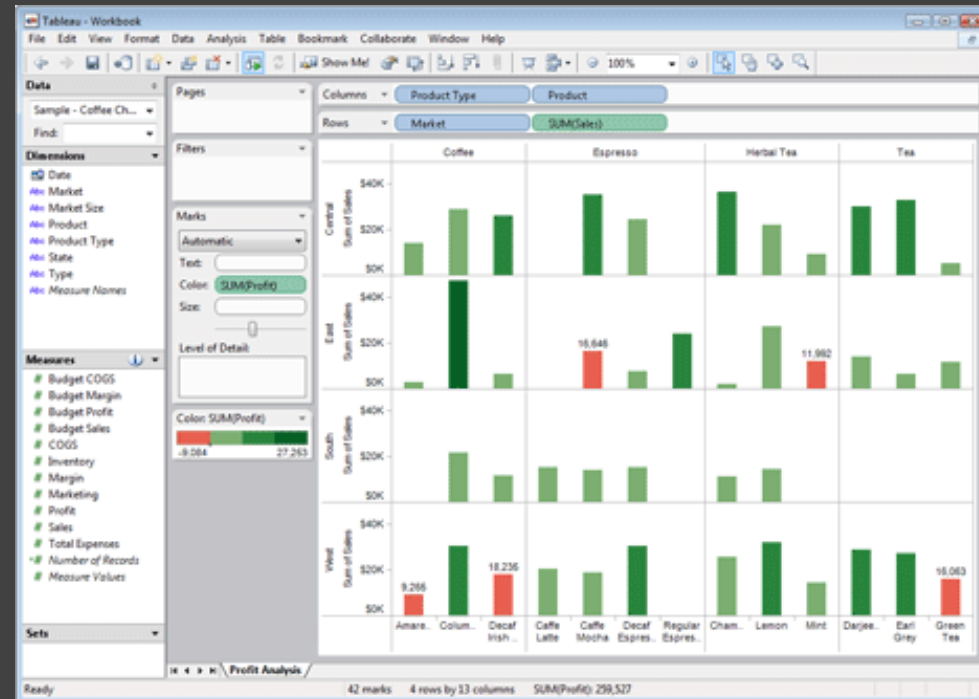
Step 2: Pose questions

Step 3: Profile the data

Iterate as needed

**Create visualizations**

Interact with data

Refine your questions

**Author a report**

Screenshots of most insightful views *(10+)*
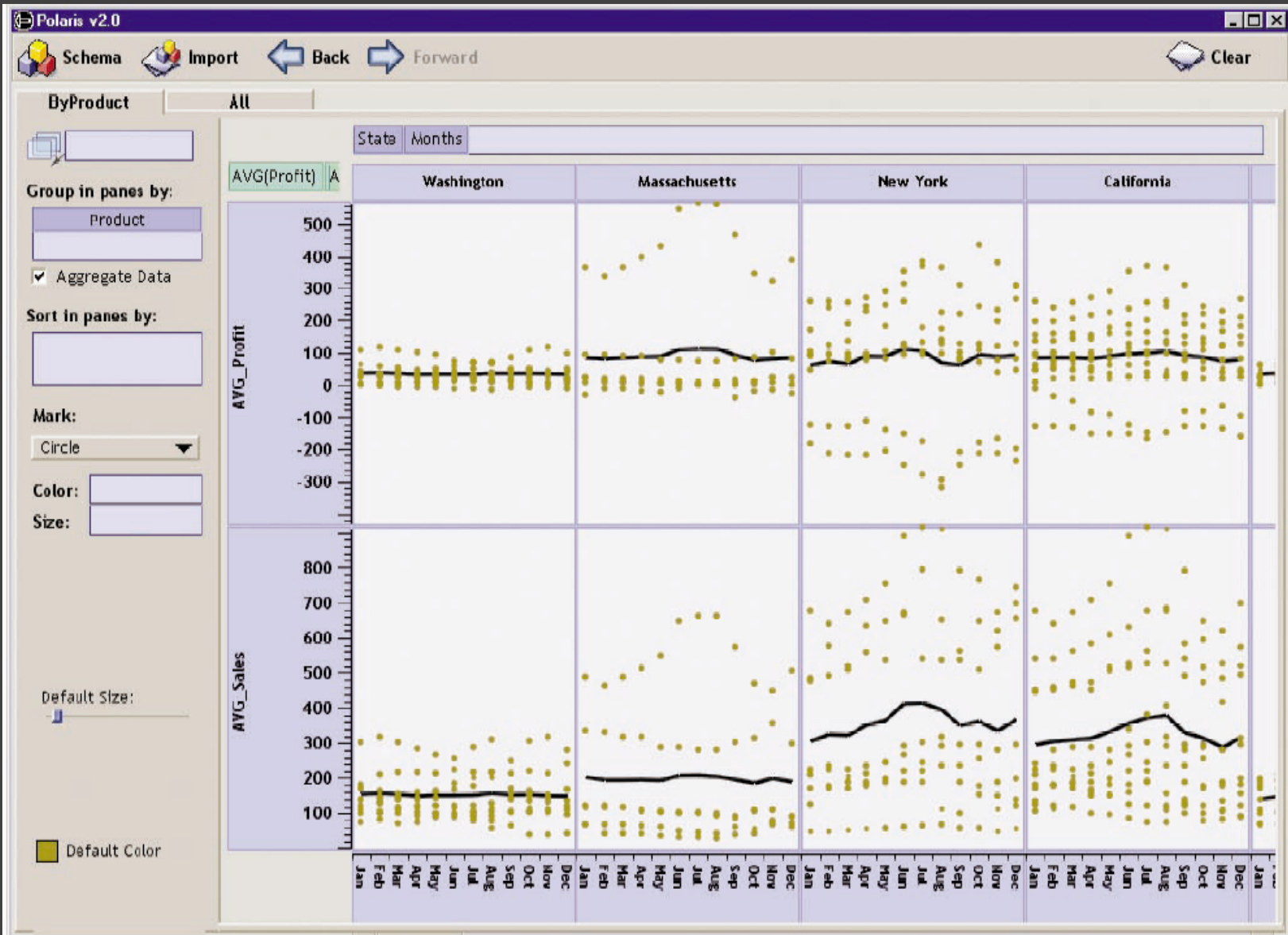
Include titles and captions for each view

Due by 5:00pm
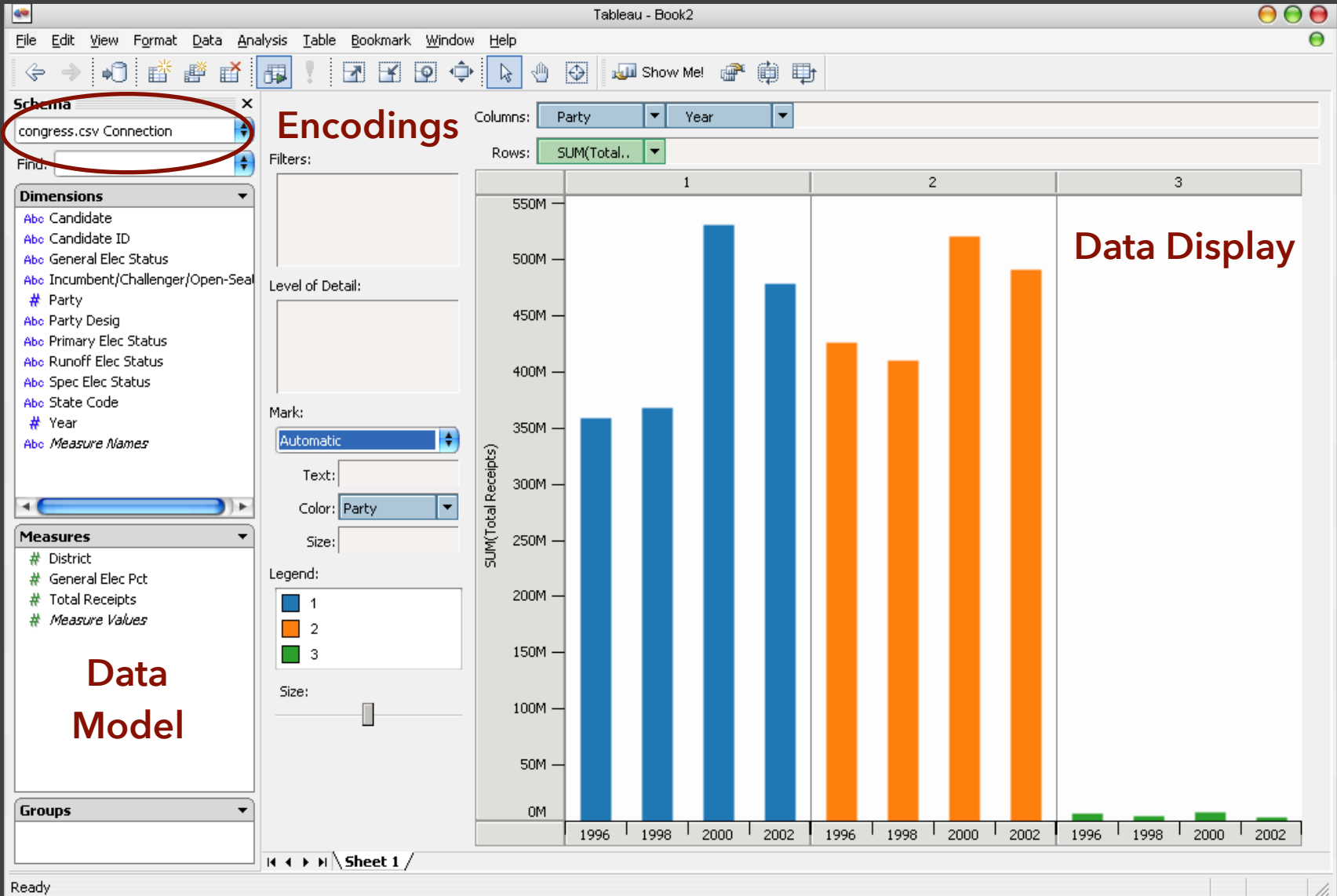
**Monday, Oct 16**

# Tableau / Polaris

# Polaris [Stolte et al.]

# Tableau

# Tableau Demo

**The dataset:**

Federal Elections Commission Receipts

Every Congressional Candidate from 1996 to 2002

4 Election Cycles

9216 Candidacies

# Dataset Schema

Year (Qi)

Candidate Code (N)

Candidate Name (N)

Incumbent / Challenger / Open-Seat (N)

Party Code (N) [1=Dem,2=Rep,3=Other]

Party Name (N)

Total Receipts (Qr)

State (N)

District (N)

This is a subset of the larger data set available from the FEC.

# Hypotheses?

What might we learn from this data?

# Hypotheses?

What might we learn from this data?

Correlation between receipts and winners?

Do receipts increase over time?

Which states spend the most?

Which party spends the most?

Margin of victory vs. amount spent?

Amount spent between competitors?

# Tableau Demo

# Tableau / Polaris Approach

Insight: can simultaneously specify both
   database queries and visualization

Choose data, then visualization, not vice versa

Use smart defaults for visual encodings

Can also suggest encodings upon request

# Specifying Table Configurations

**Operands are the database fields**

Each operand interpreted as a set {…}

Quantitative and Ordinal fields treated differently

**Three operators:**

concatenation (+)

cross product (x)

nest (/)

# Table Algebra: Operands

**Ordinal fields**: interpret domain as a set that partitions table into rows and columns.

Quarter = {(Qtr1),(Qtr2),(Qtr3),(Qtr4)} ->

| Qtr1 | Qtr2 | Qtr3 | Qtr4 |
|------|------|------|------|
| 95892 | 101760 | 105282 | 98225 |

**Quantitative fields**: treat domain as single element set and encode spatially as axes.

Profit = {(Profit[-410,650])} ->

# Concatenation (+) Operator

**Ordered union of set interpretations**
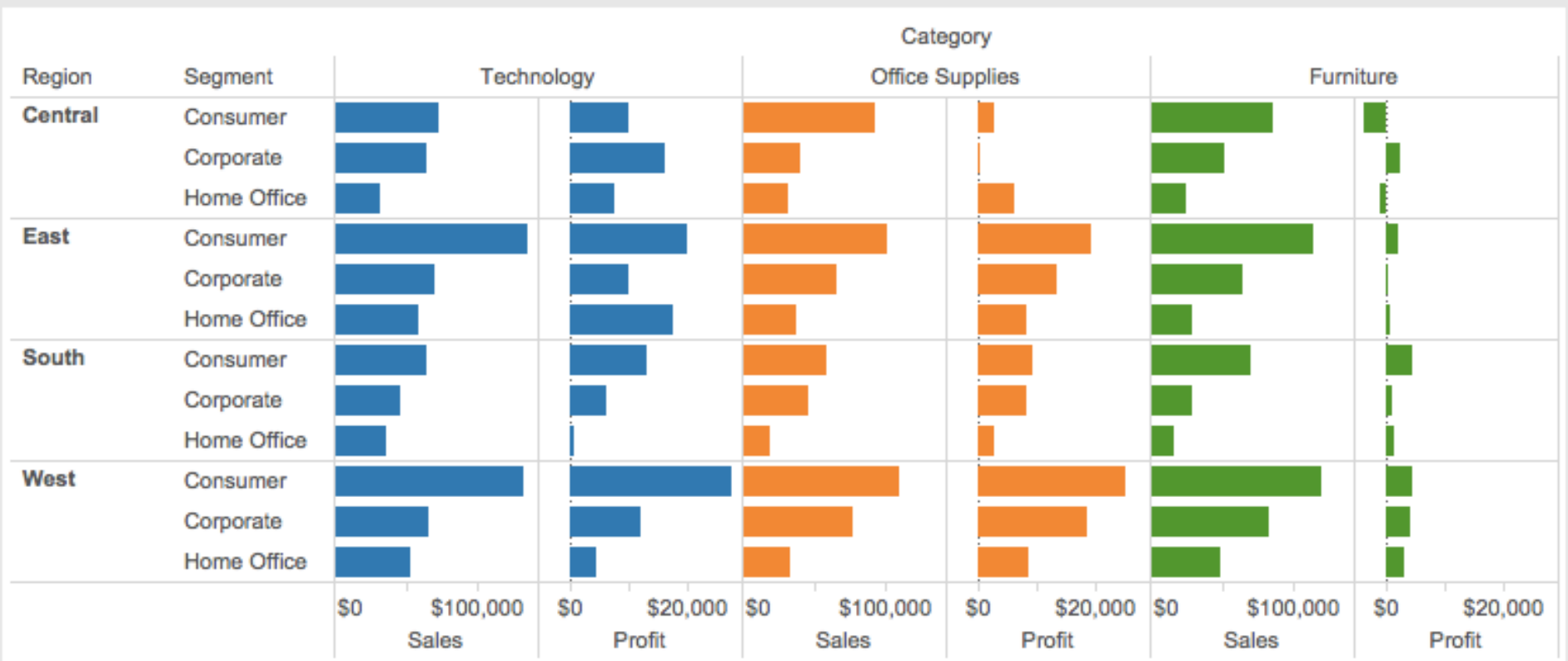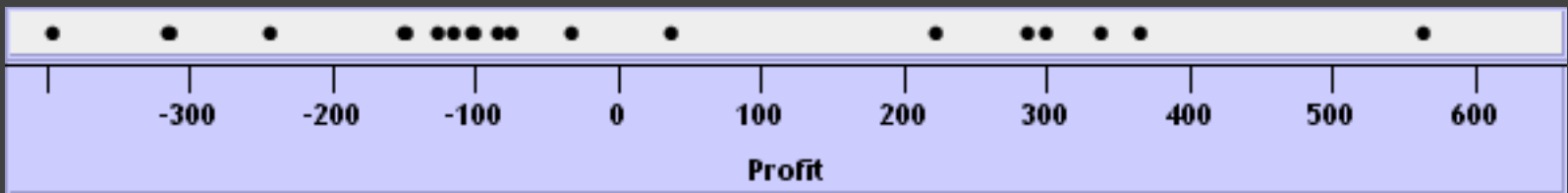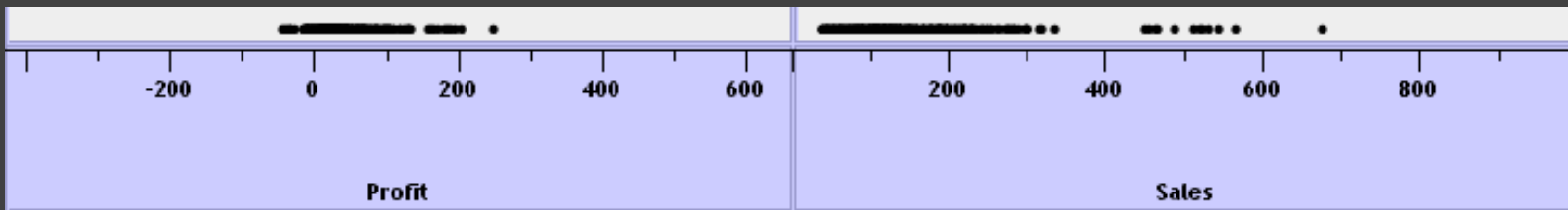
Quarter + Product Type
= {(Qtr1),(Qtr2),(Qtr3),(Qtr4)} + {(Coffee), (Espresso)}
= {(Qtr1),(Qtr2),(Qtr3),(Qtr4),(Coffee),(Espresso)}

| Qtr1 | Qtr2 | Qtr3 | Qtr4 | Coffee | Espresso |
|------|------|------|------|--------|----------|
| 48   | 59   | 57   | 53   | 151    | 21       |

Profit + Sales = {(Profit[-310,620]),(Sales[0,1000])}

# Cross (x) Operator

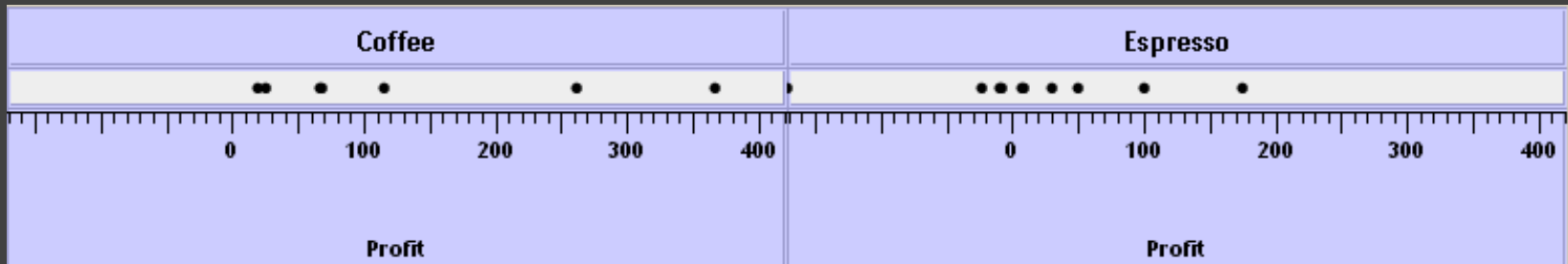## Cross-product of set interpretations

Quarter x Product Type =

    {(Qtr1,Coffee), (Qtr1, Tea), (Qtr2, Coffee), (Qtr2, Tea), (Qtr3, Coffee), (Qtr3, Tea), (Qtr4, Coffee), (Qtr4,Tea)}

| Qtr1 | | Qtr2 | | Qtr3 | | Qtr4 | |
|---|---|---|---|---|---|---|---|
| Coffee | Espresso | Coffee | Espresso | Coffee | Espresso | Coffee | Espresso |
| 131 | 19 | 160 | 20 | 178 | 12 | 134 | 33 |

Product Type x Profit =

# Nest (/) Operator

**Cross-product filtered by existing records**

Quarter x Month ->

creates twelve entries for each quarter. i.e., (Qtr1, December)

Quarter / Month ->

creates three entries per quarter based on tuples in database (not semantics)

# Table Algebra

The operators (+, x, /) and operands (O, Q) provide
  an *algebra* for tabular visualization.

Algebraic statements are then mapped to:
  **Visualizations** - trellis plot partitions, visual encodings
  **Queries** - selection, projection, group-by aggregation

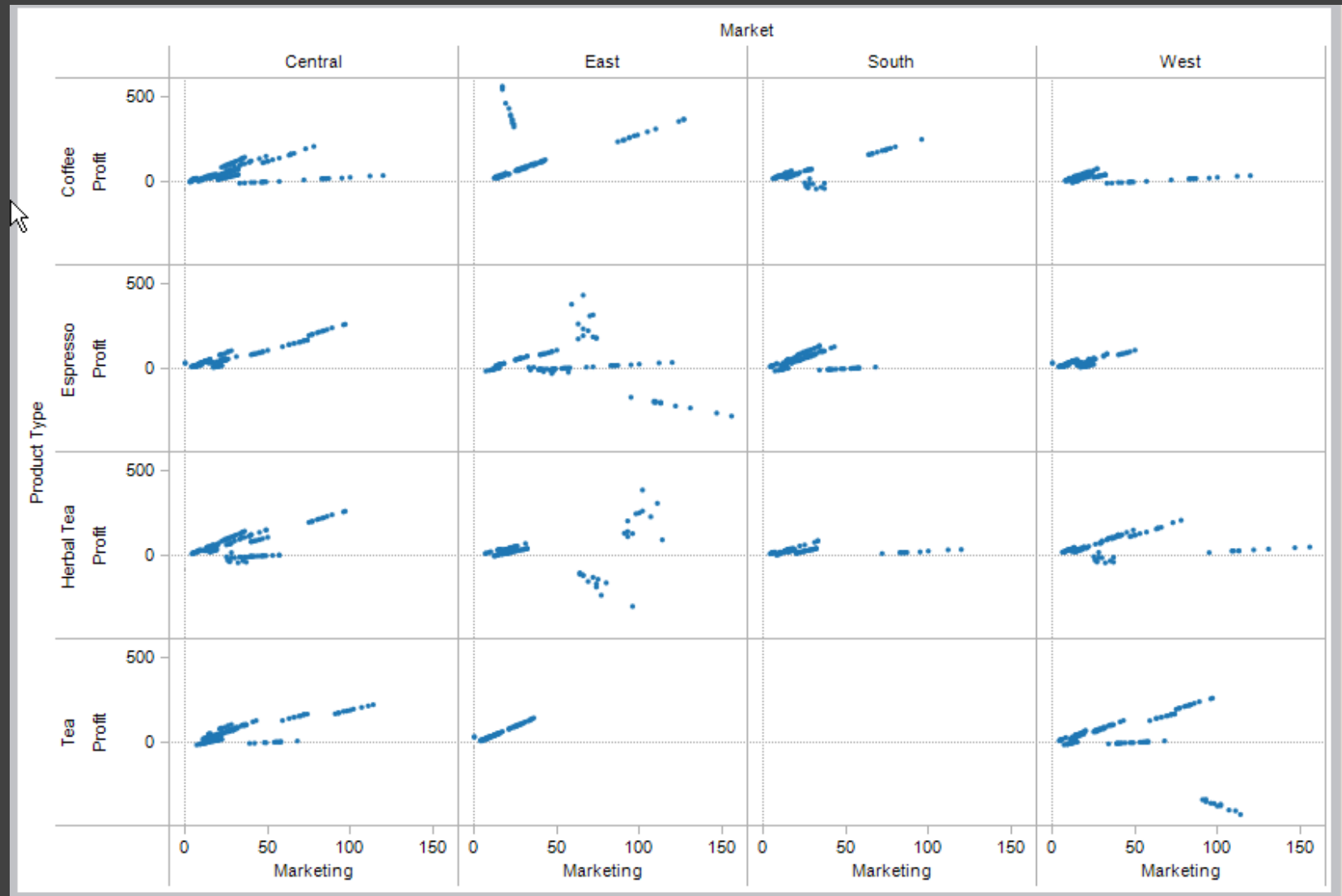In Tableau, users make statements via drag-and-drop
  Note that this specifies operands *NOT* operators!
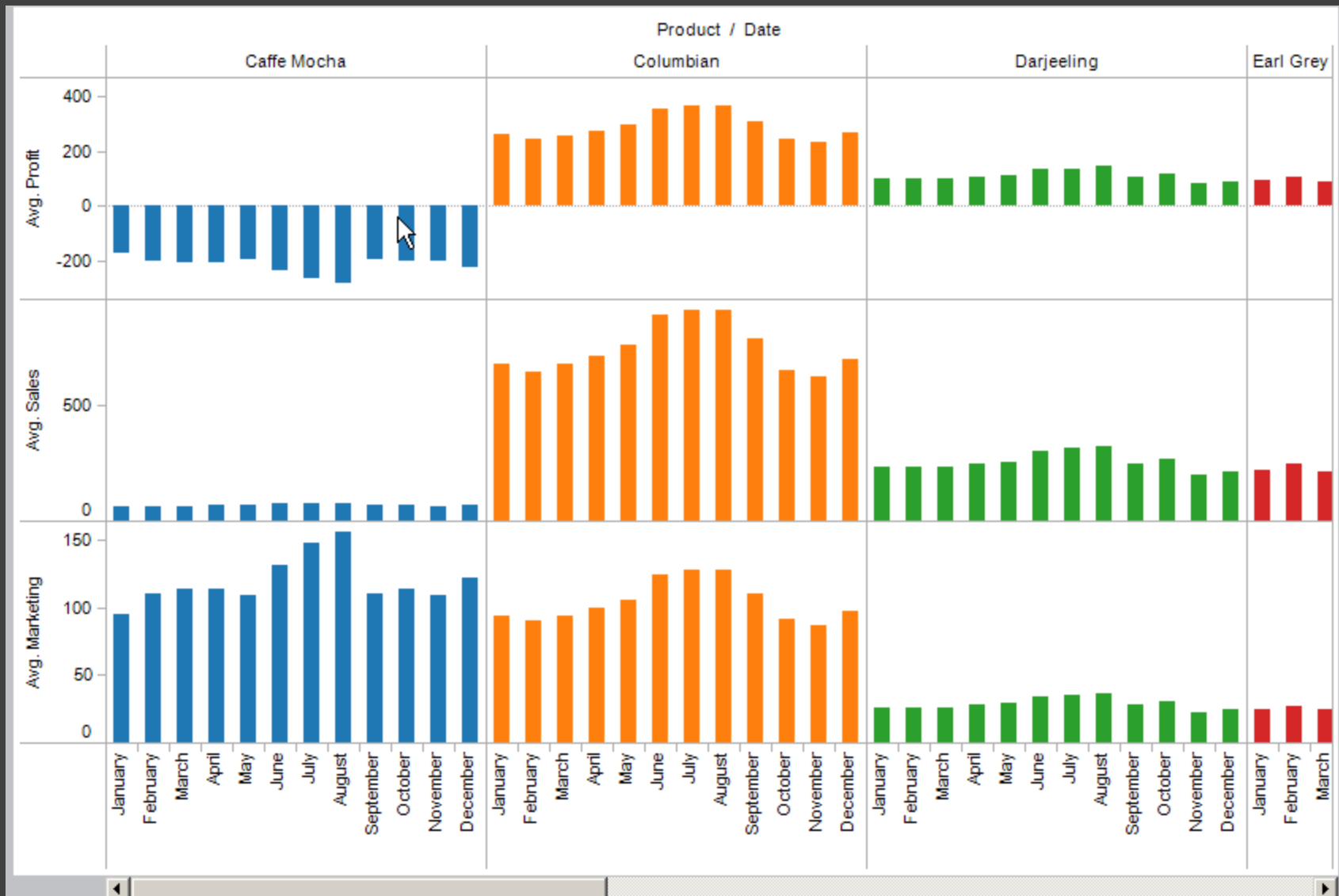  Operators are inferred by data type (O, Q)

# Ordinal-Ordinal

# Quantitative-Quantitative

# Ordinal-Quantitative

# Querying the Database



(1) Select records from the database, filtering by user-defined criteria.

(2) Partition the records into layers and panes. The same record may appear in multiple partitions.

(3) Group, sort, and aggregate the relations within each pane.

(4) Render and compose layers.
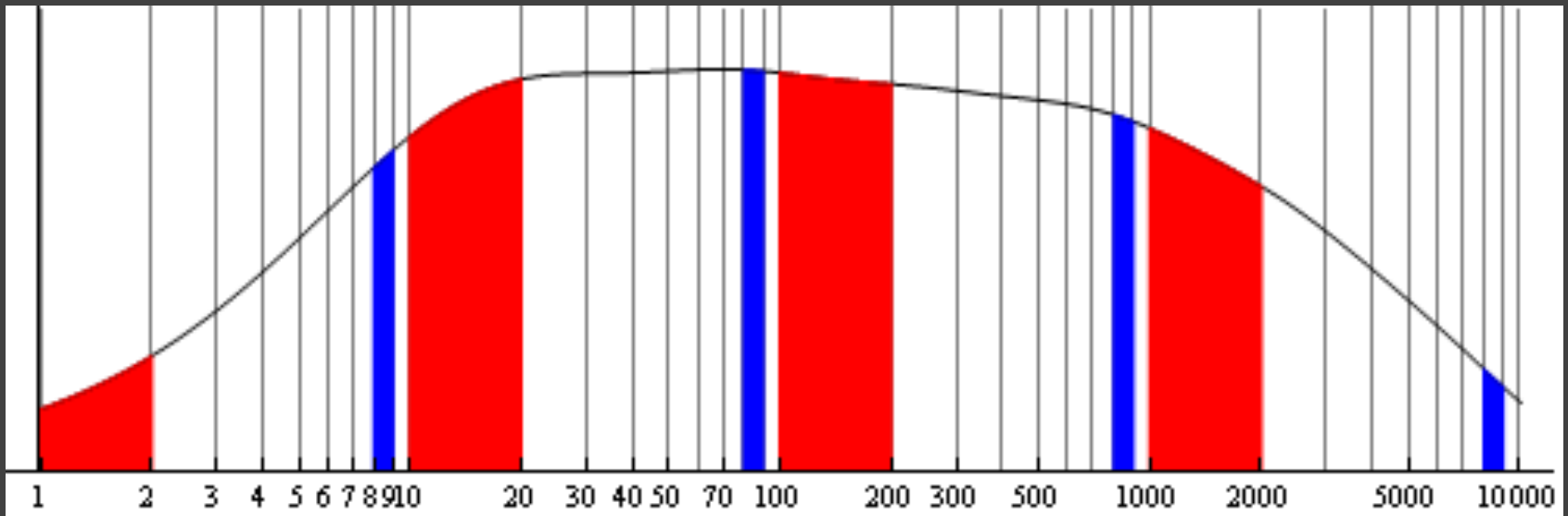
**BONUS TOPIC**

# Data Fraud

# A Detective Story

You have accounting records for two firms that are in dispute. One is lying. *How to tell?*

| *Firm A* | | *Firm B* | LIARS! |
|---:|---:|---:|---:|
| 283.08 | 25.23 | 283.08 | 75.23 |
| 153.86 | 385.62 | 353.86 | 185.25 |
| 1448.97 | 12371.32 | 5322.79 | 9971.42 |
| 18595.91 | 1280.76 | 8795.64 | 4802.43 |
| 21.33 | 257.64 | 61.33 | 57.64 |

Amt. Paid: $34823.72       Amt. Rec'd: $29908.67

# **Benford's Law** (Benford 1938, Newcomb 1881)

The *logarithms* of the values (not the values themselves) are uniformly randomly distributed.



Hence the leading digit **1** has a ~30% likelihood. Larger digits are increasingly less likely.

# **Benford's Law** (Benford 1938, Newcomb 1881)

The *logarithms* of the values (not the values themselves) are uniformly randomly distributed.

Holds for many (but certainly not all) real-life data sets: Addresses, Bank accounts, Building heights, …

Data must span multiple orders of magnitude.

Evidence that records do not follow Benford's Law is admissible in a court of law!