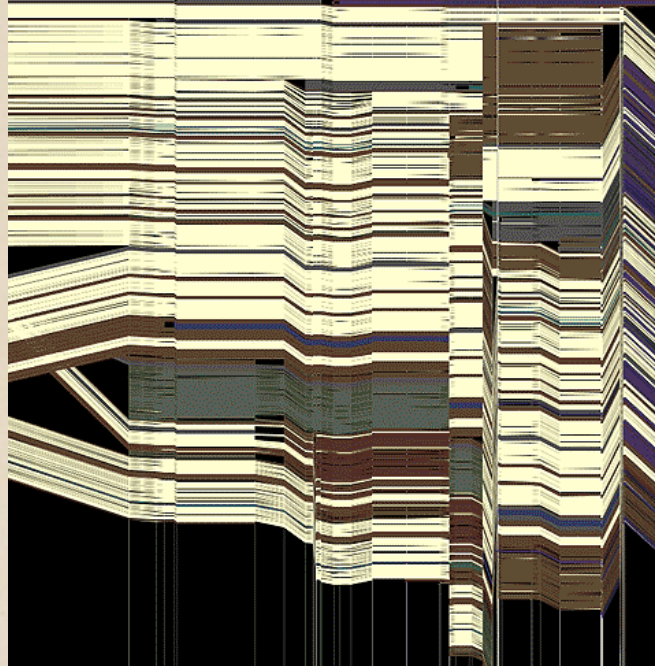# CSE 412 - Intro to Data Visualization
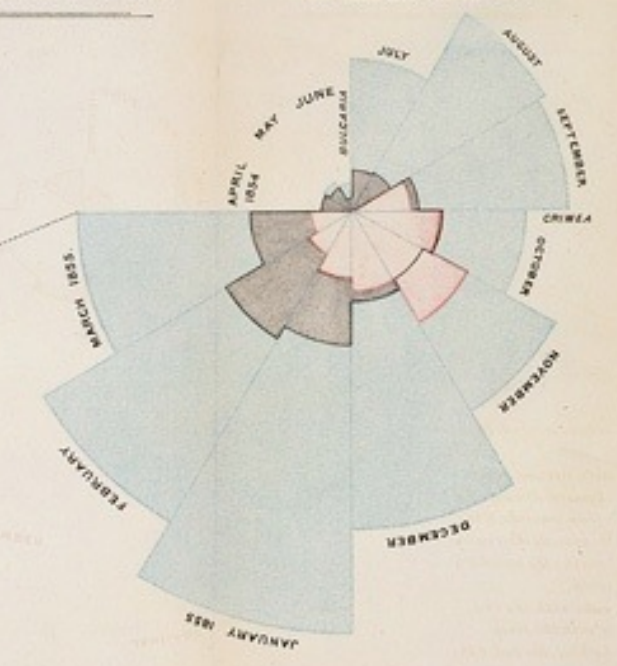
# Exploratory Data Analysis
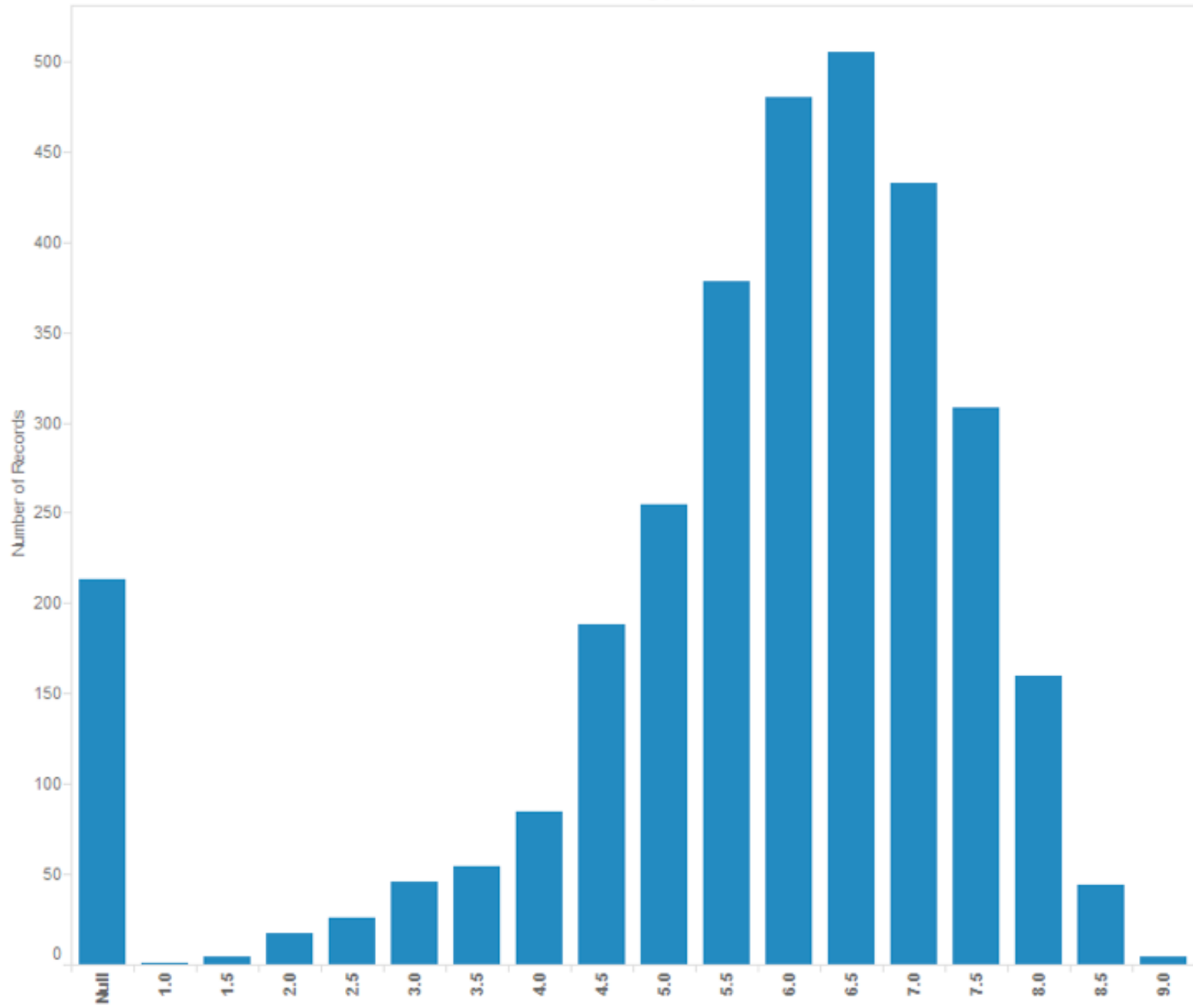


Jane Hoffswell  University of Washington

# Analysis Example: Motion Pictures Data
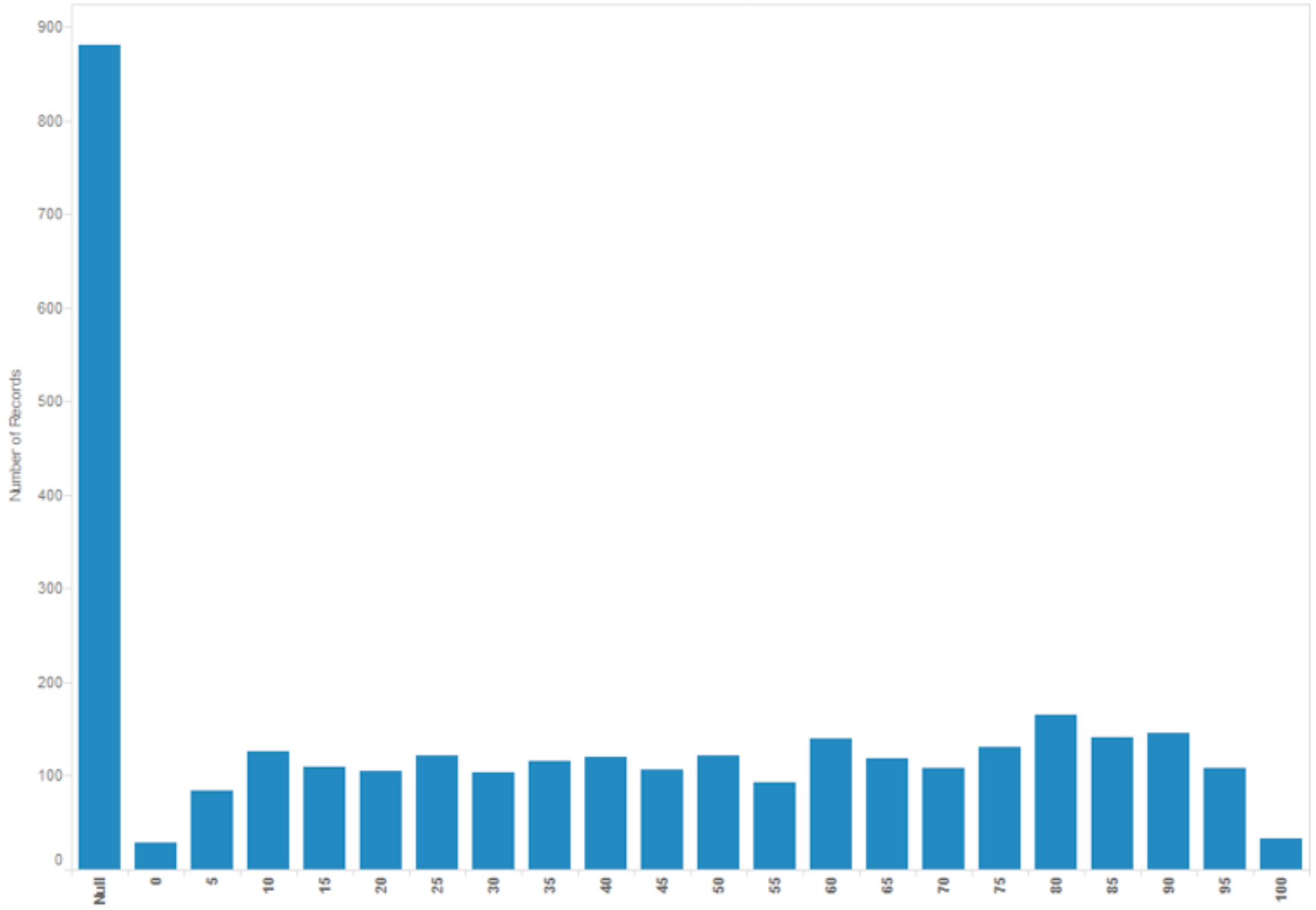
# Motion Pictures Data

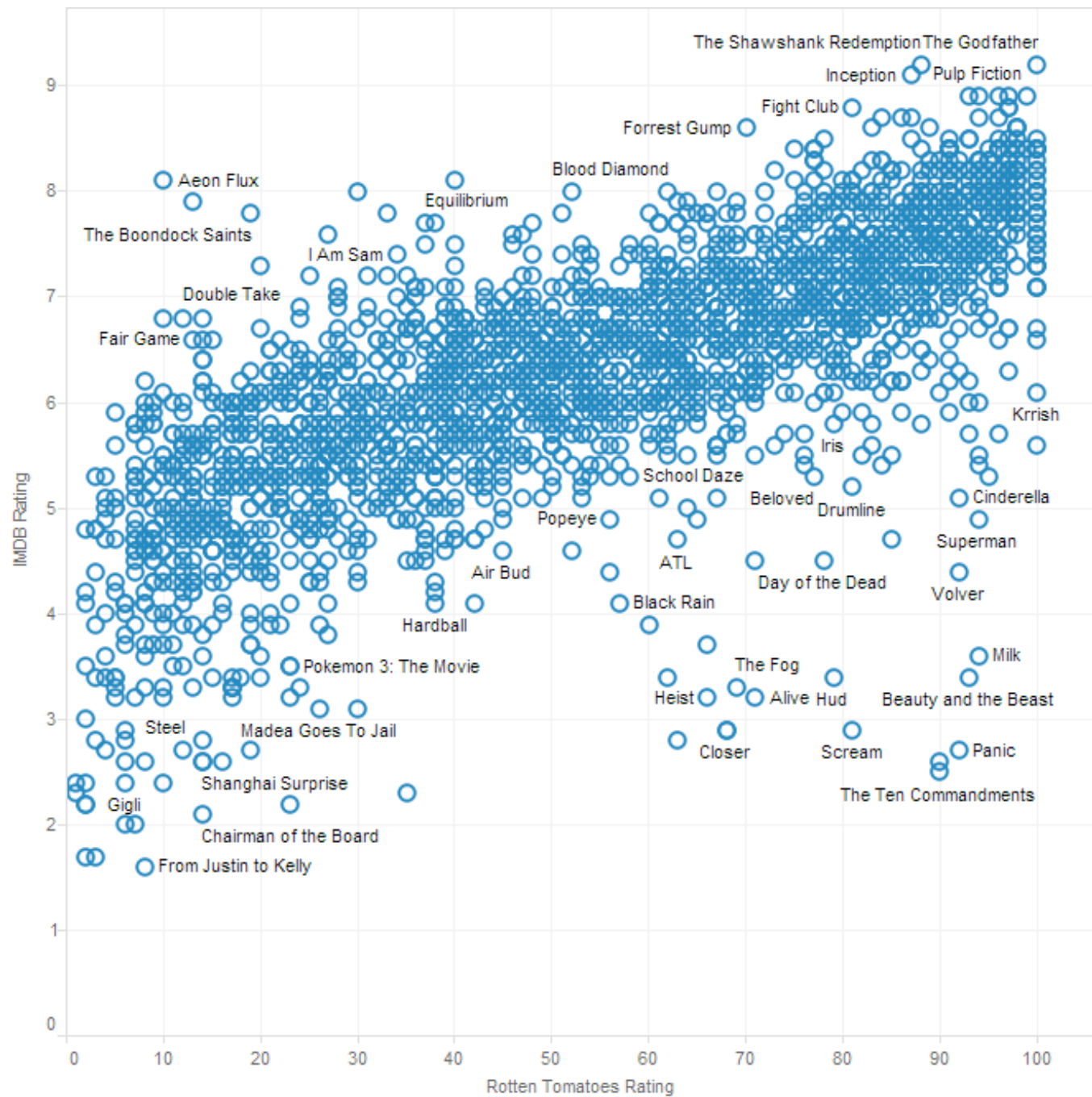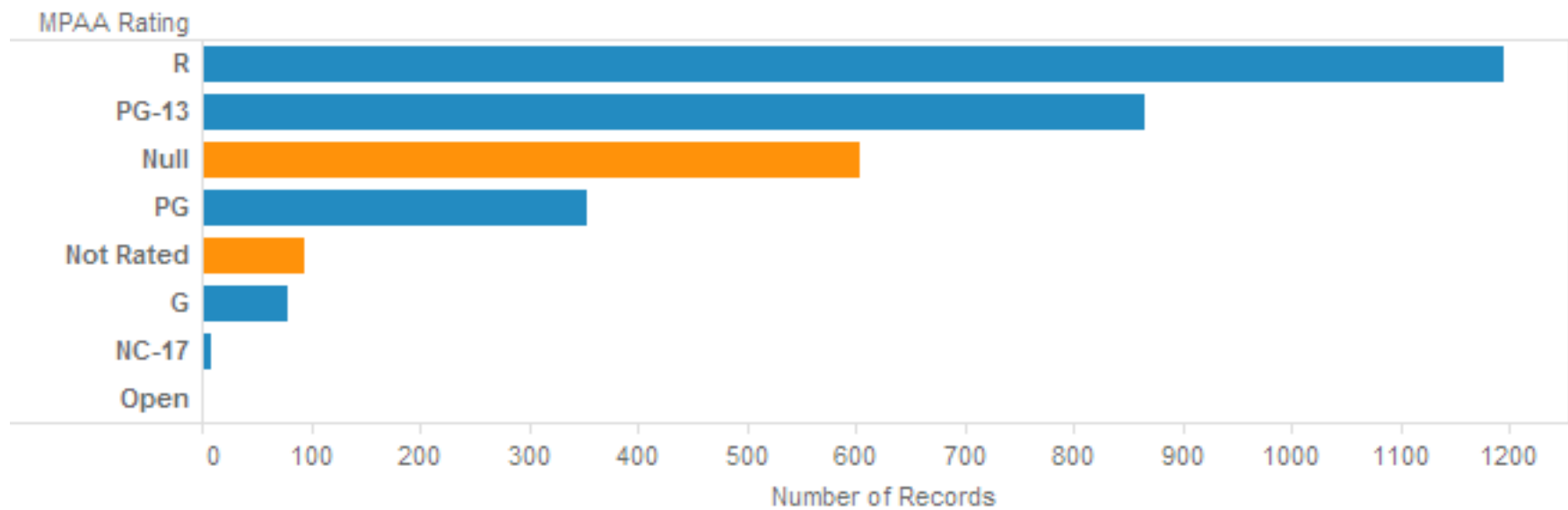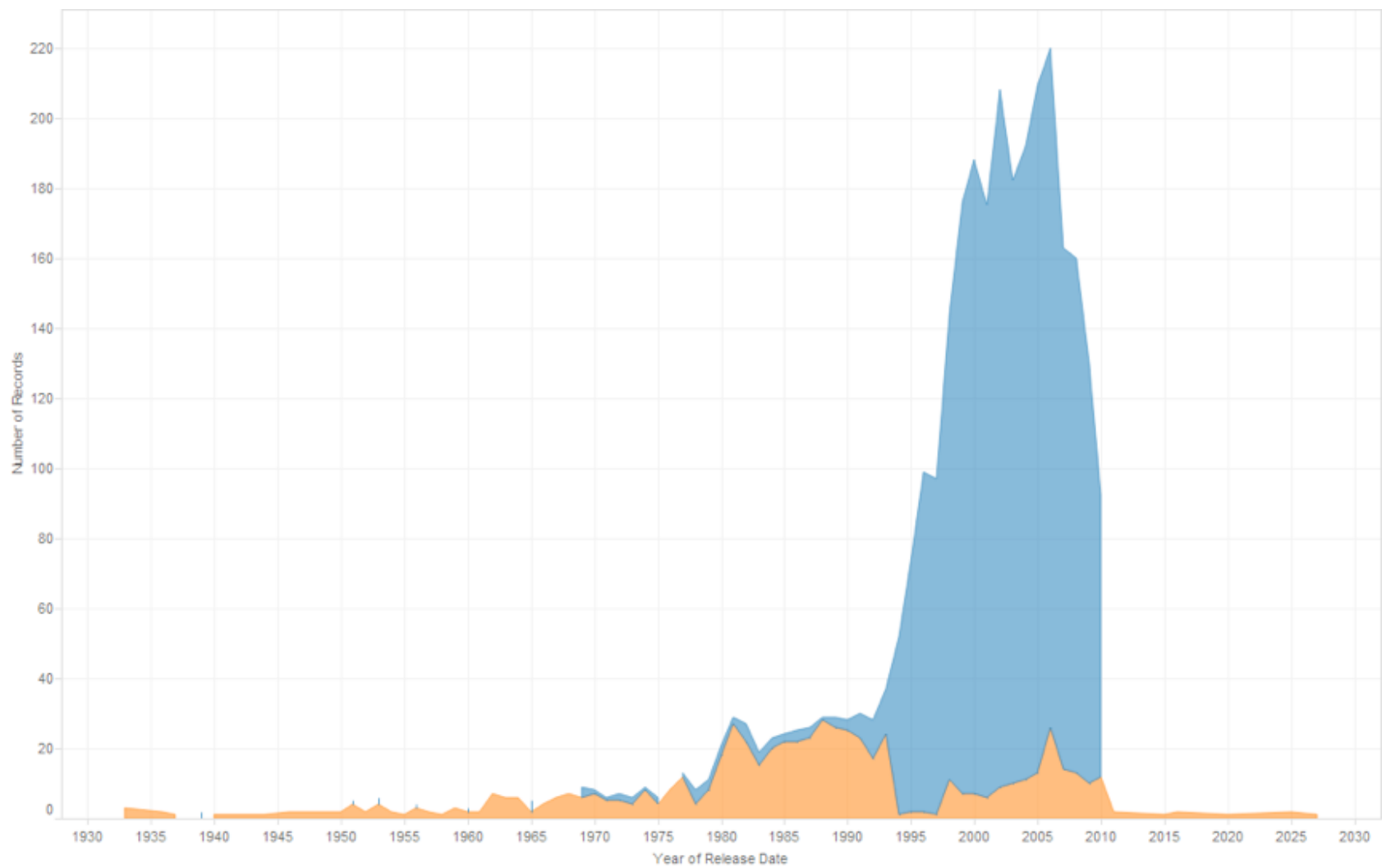| | |
|---|---|
| Title | String (N) |
| IMDB Rating | Number (Q) |
| Rotten Tomatoes Rating | Number (Q) |
| MPAA Rating | String (O) |
| Release Date | Date (T) |

IMDB Rating (bin)

Rotten Tomatoes Rating (bin)

Scatter plot of IMDB Rating (y-axis) versus Rotten Tomatoes Rating (x-axis). Labeled points include:

- The Shawshank Redemption
- The Godfather
- Inception
- Pulp Fiction
- Fight Club
- Forrest Gump
- Blood Diamond
- Aeon Flux
- Equilibrium
- The Boondock Saints
- I Am Sam
- Double Take
- Fair Game
- Krrish
- Iris
- School Daze
- Beloved
- Cinderella
- Drumline
- Popeye
- ATL
- Superman
- Air Bud
- Day of the Dead
- Volver
- Hardball
- Black Rain
- Milk
- Pokemon 3: The Movie
- The Fog
- Heist
- Alive Hud
- Beauty and the Beast
- Steel
- Madea Goes To Jail
- Closer
- Scream
- Panic
- Shanghai Surprise
- The Ten Commandments
- Gigli
- Chairman of the Board
- From Justin to Kelly

The Godfather

Inception

Fight Club

The Godfather: Part II

Forrest Gump

Aeon Flux

Blood Diamond

Casino Royale

Saw

I Am Sam

Double Take

Krrish

Fair Game

Iris

Cinderella

Beloved
Drumline

Superman

Popeye

ATL

Day of the Dead

Volver

Air Bud

Black Rain

Hardball

The Fog

Milk

Pokemon 3: The Movie

Heist
Alive
Hud

Beauty and the Beast

Steel

Madea Goes To Jail

Closer

Scream

Panic

The Ten Commandments

Chairman of the Board

From Justin to Kelly

Premonition  Dude, Where's My Car?

Bad Lieutenant: Port of Call New Orleans

IMDB Rating

Rotten Tomatoes Rating

# 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)
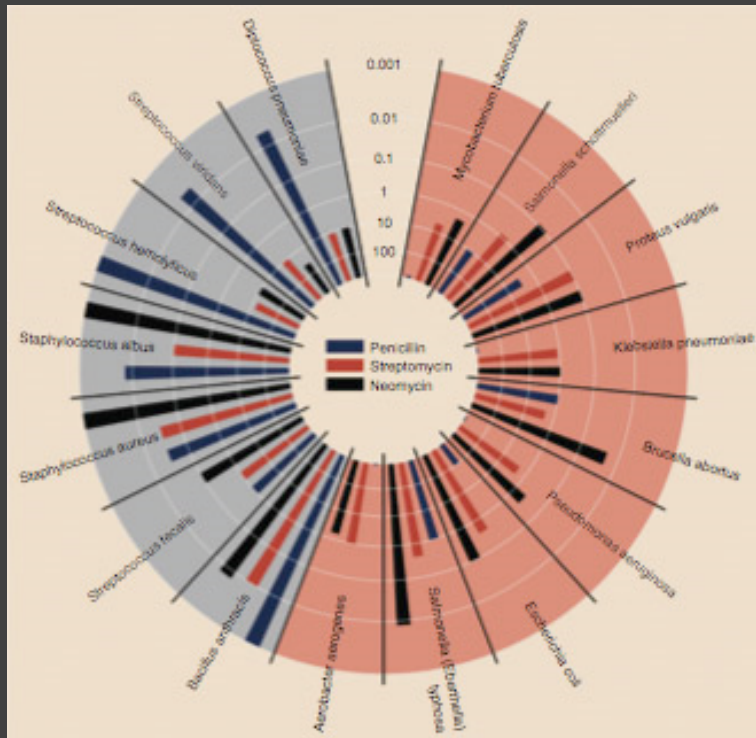
Gram-Staining?                          Pos / Neg (N)

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

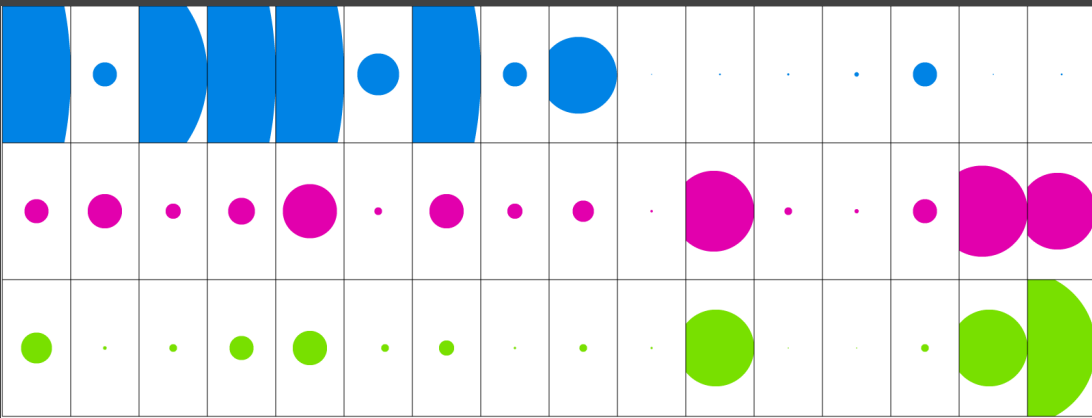# How do the drugs compare?



X-axis: Antibiotic | log(MIC)
Y-axis: Gram-Staining | Species
Color:  Most-Effective?

minimum inhibitory concentration
of antibiotics

bowen li
cs448b

Bowen Li
Stanford CS448B, Fall 2009

## All bacteria

Streptomycin and Neomycin are more efficient broad-spectrum antibiotics than Penicilin.

Proportion of bacteria strains inhibited

Concentration (µg/ml)

## Gram-negative bacteria only

Neomycin and Streptomycin are more efficient against gram-negative bacteria, so can be used at a lower dosage here than above.

Gram staining quickly identifies bacteria as Gram-negative or Gram-positive, which can be used to find a more efficient antibiotic and dosage.

## Gram-positive bacteria only

Penicilin is more efficient than either Streptomycin or Neomycin if the bacteria is known to be gram-positive.

Proportion of bacteria strains inhibited

Concentration (µg/ml)

### Penicillin
0.001
0.001
0.005
0.005
0.007
0.03
1

1
1
3
10
100
800
850
850
870

### Streptomycin
0.01
14
11
10
0.1
0.03
1

2
0.4
0.1
0.8
0.4
5
1.2
2
1

### Neomycin
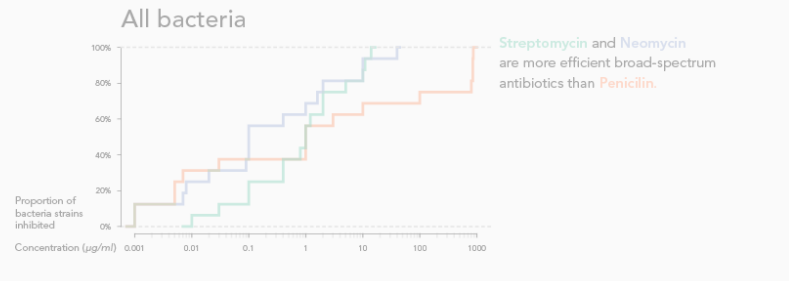0.007
10
10
40
0.001
0.001
0.1

0.02
0.008
0.1
0.09
0.1
2
1
0.4
1.6

Minimum Inhibitory Concentration (MIC)

### Effectiveness of Antibiotics

Rows: A. aerogenes, B. abortus, E. coli, K. pneumoniae, M. tuberculosis, P. vulgaris, P. aeruginosa, S. typhosa, S. schottmuelleri, B. anthracis, D. pneumoniae, S. albus, S. aureus, S. fecalis, S. hemolyticus, S. viridans

Columns: Penicillin, Streptomycin, Neomycin

darker colors: more effective

MIC (ug/uL)

Penicillin     Streptomycin     Neomycin

$Log_{10}$ Minimum Inhibitory Concentration (µg/mL)
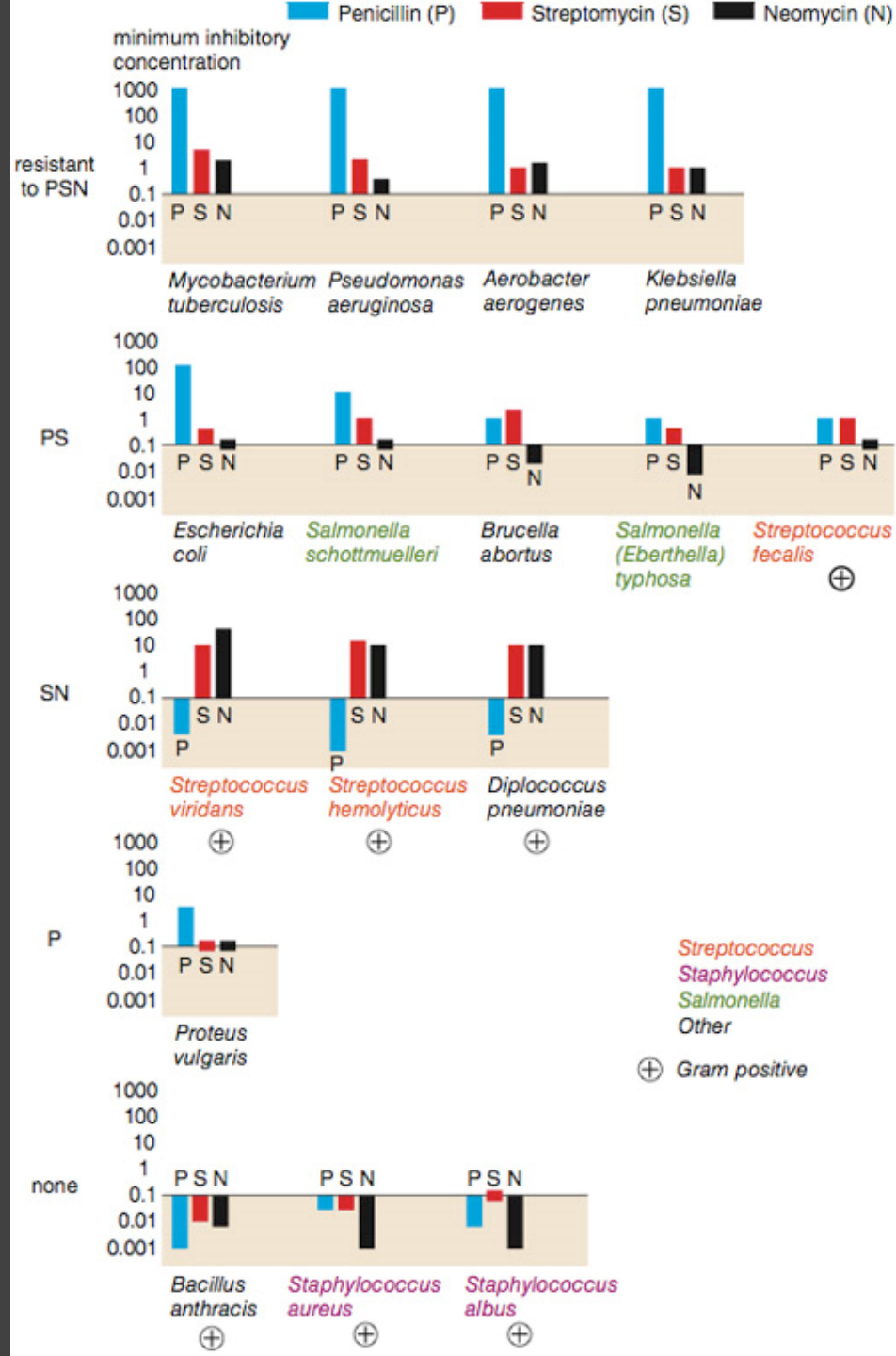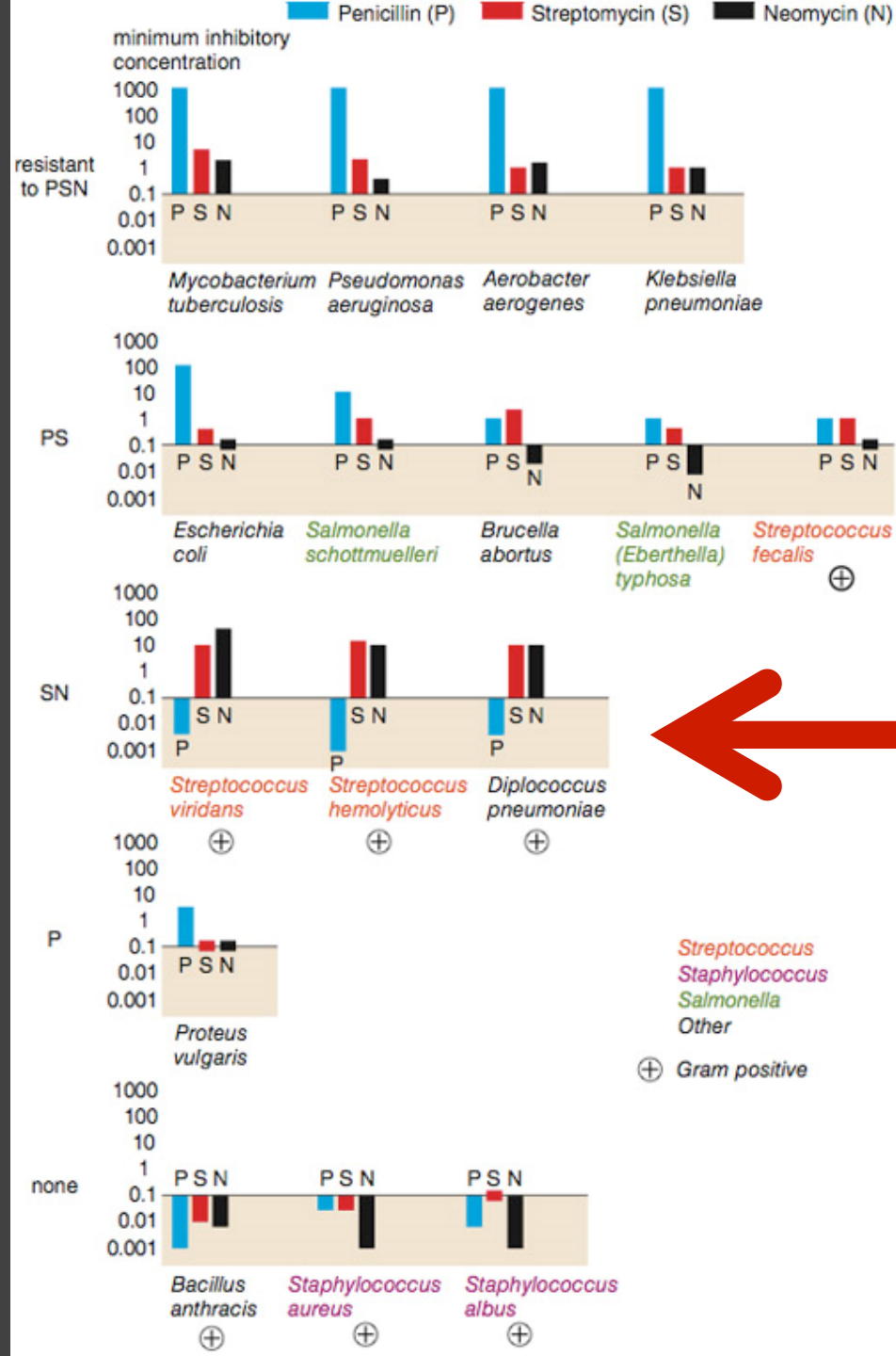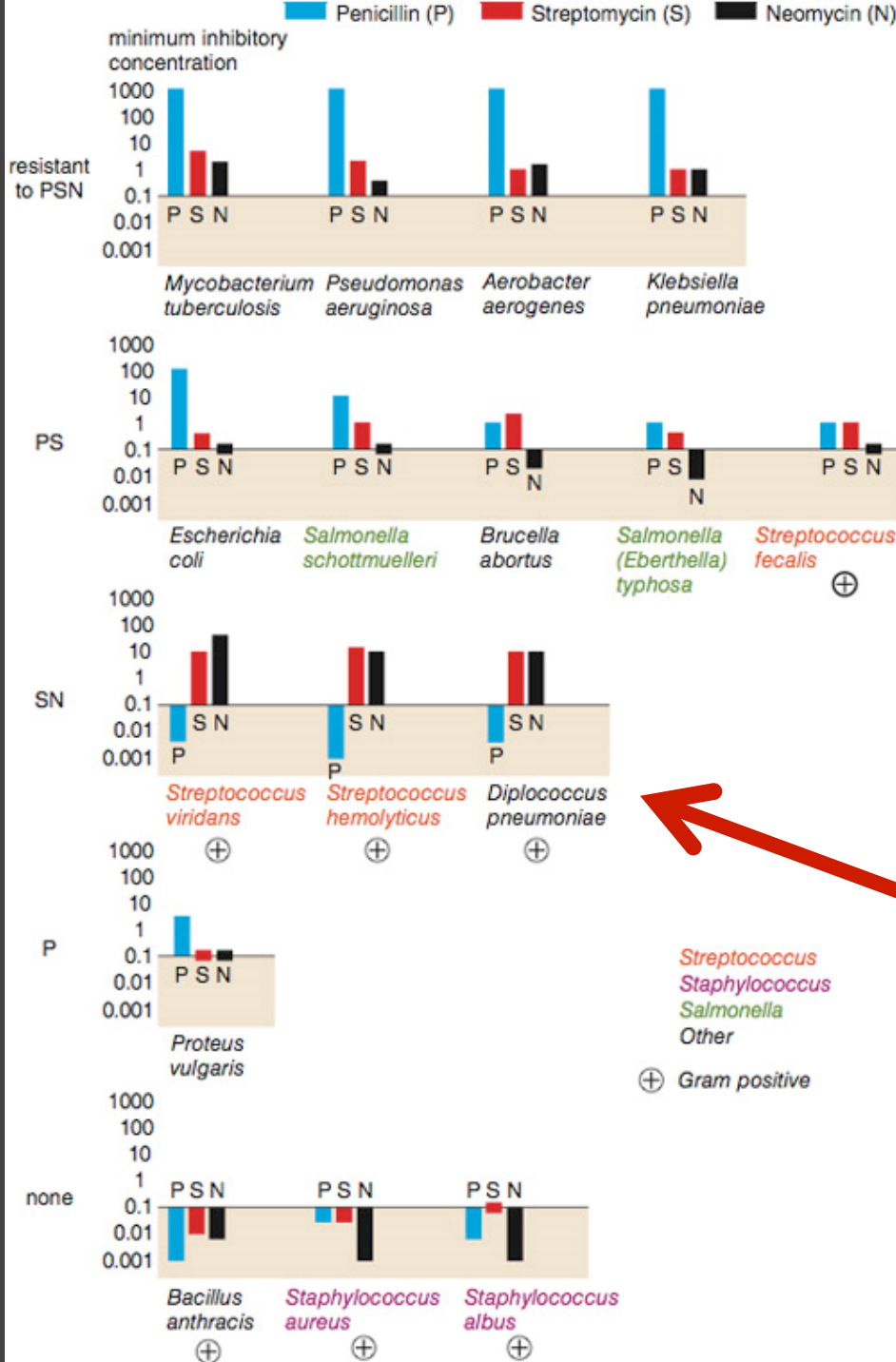
1000
100
10
1
0.1
0.01
0.001
0.0001

# Do the bacteria group by antibiotic resistance?

# Do the bacteria group by antibiotic resistance?

Do the bacteria group by antibiotic resistance?

Wainer & Lysen
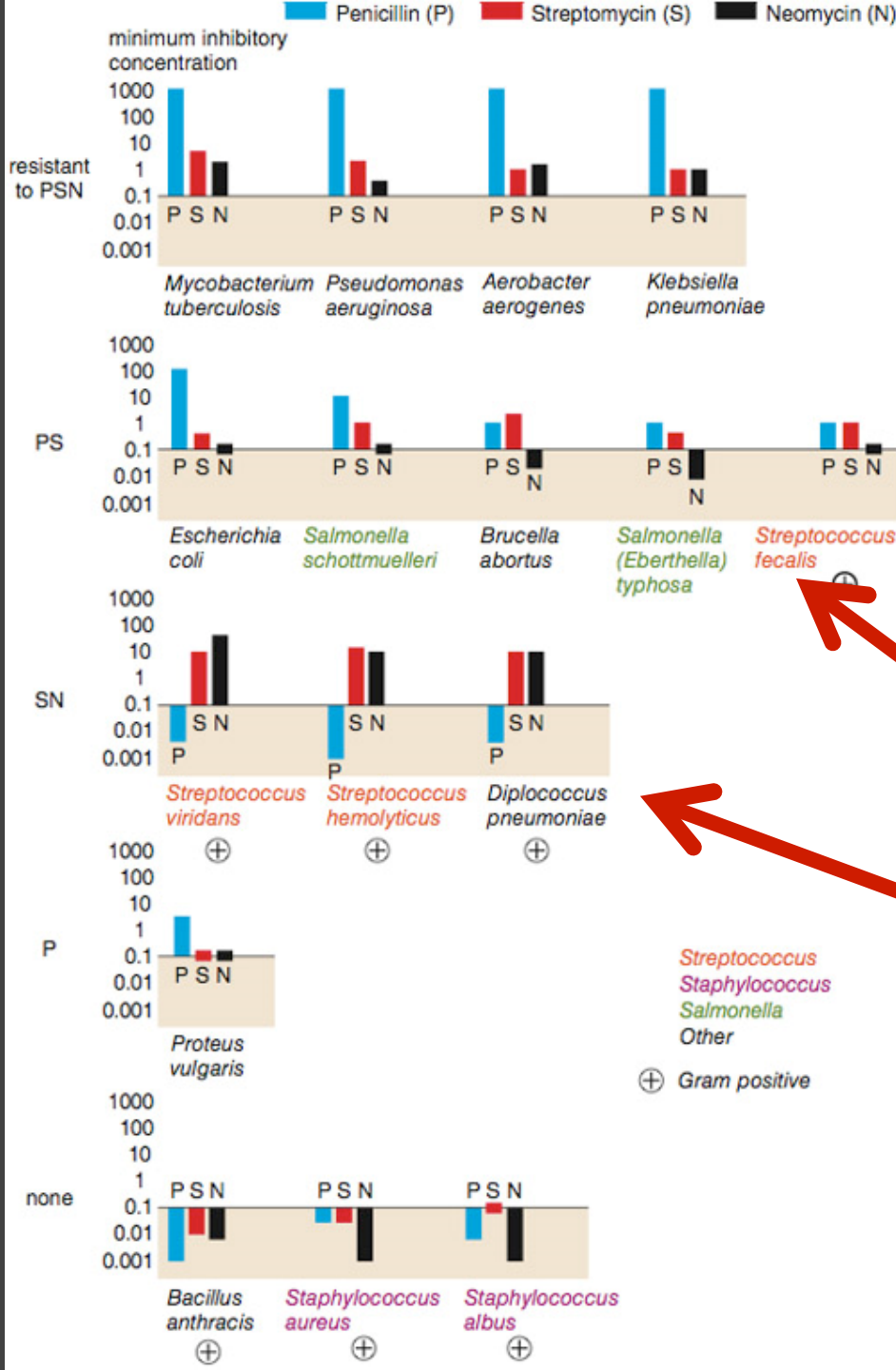*American Scientist*, 2009

**Do the bacteria group by antibiotic resistance?**

Really a streptococcus! (realized ~20 yrs later)

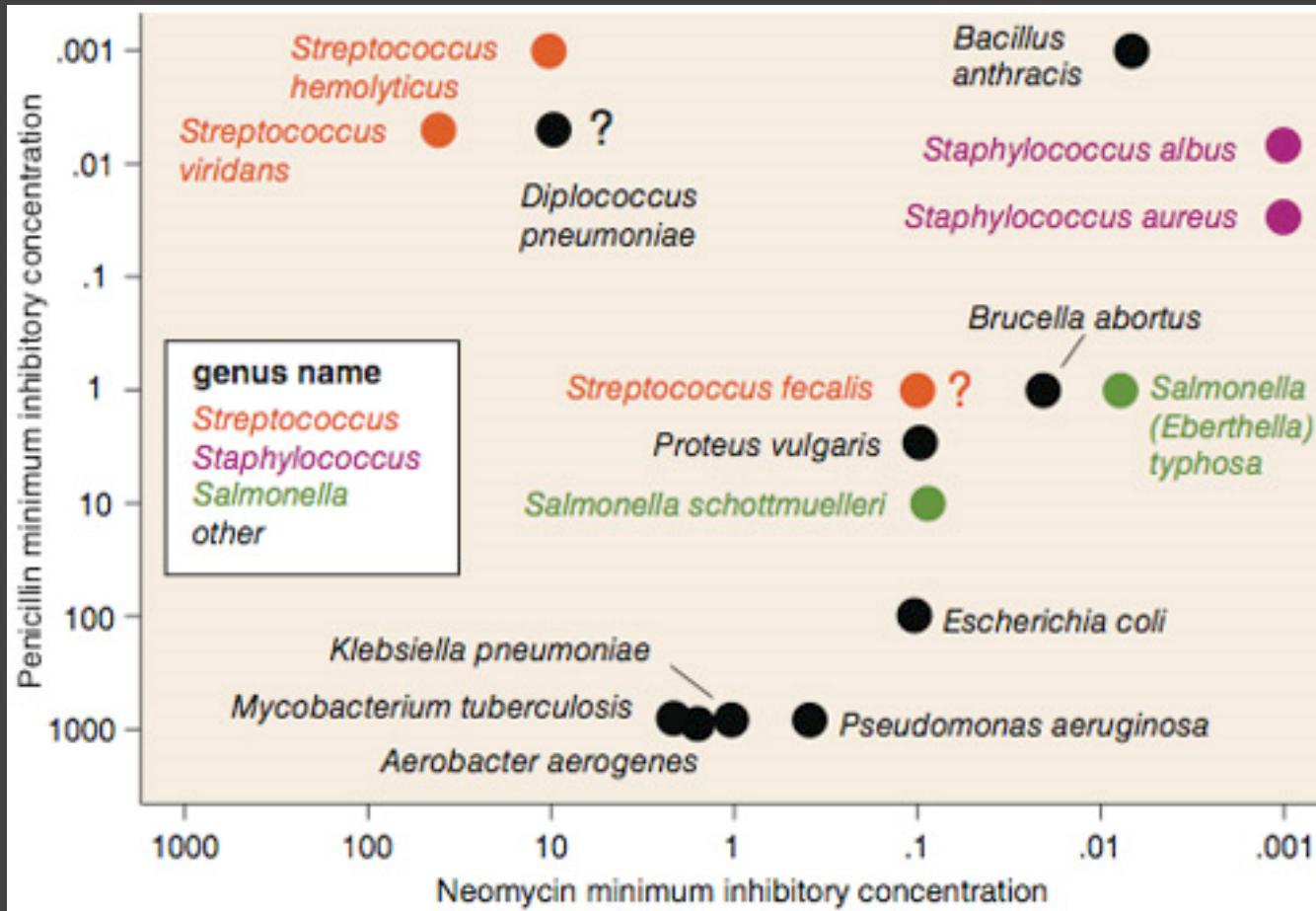**Do the bacteria group by antibiotic resistance?**

Not a streptococcus! (realized ~30 yrs later)

Really a streptococcus! (realized ~20 yrs later)

Wainer & Lysen
*American Scientist*, 2009

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

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

# 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

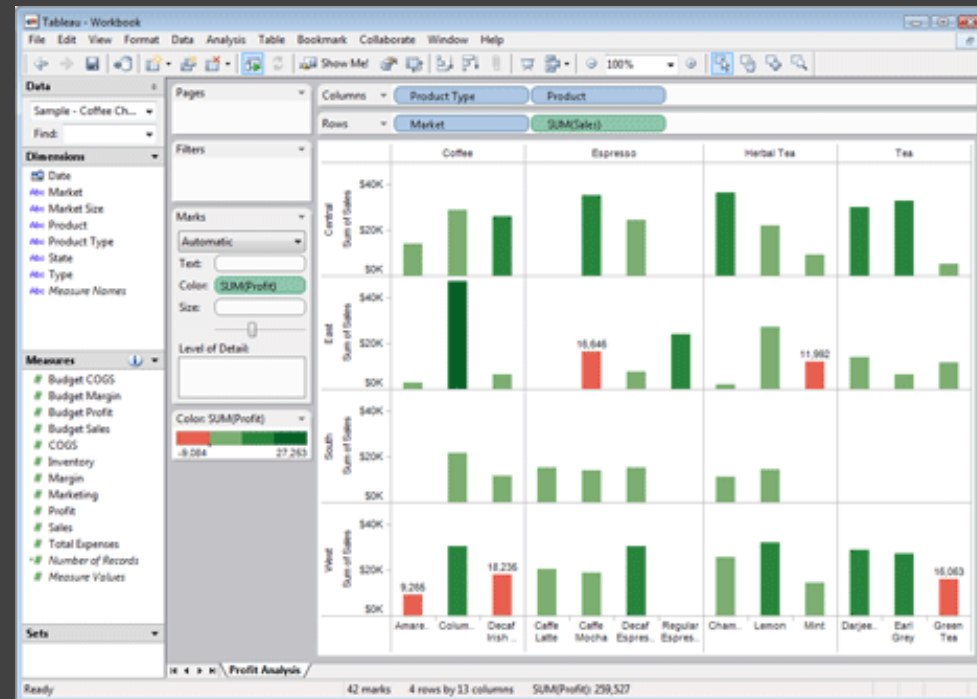Iterate as needed

**Create visualizations**

Interact with data

Refine your questions

**Author a report**

Screenshots of most insightful views *(8+)*
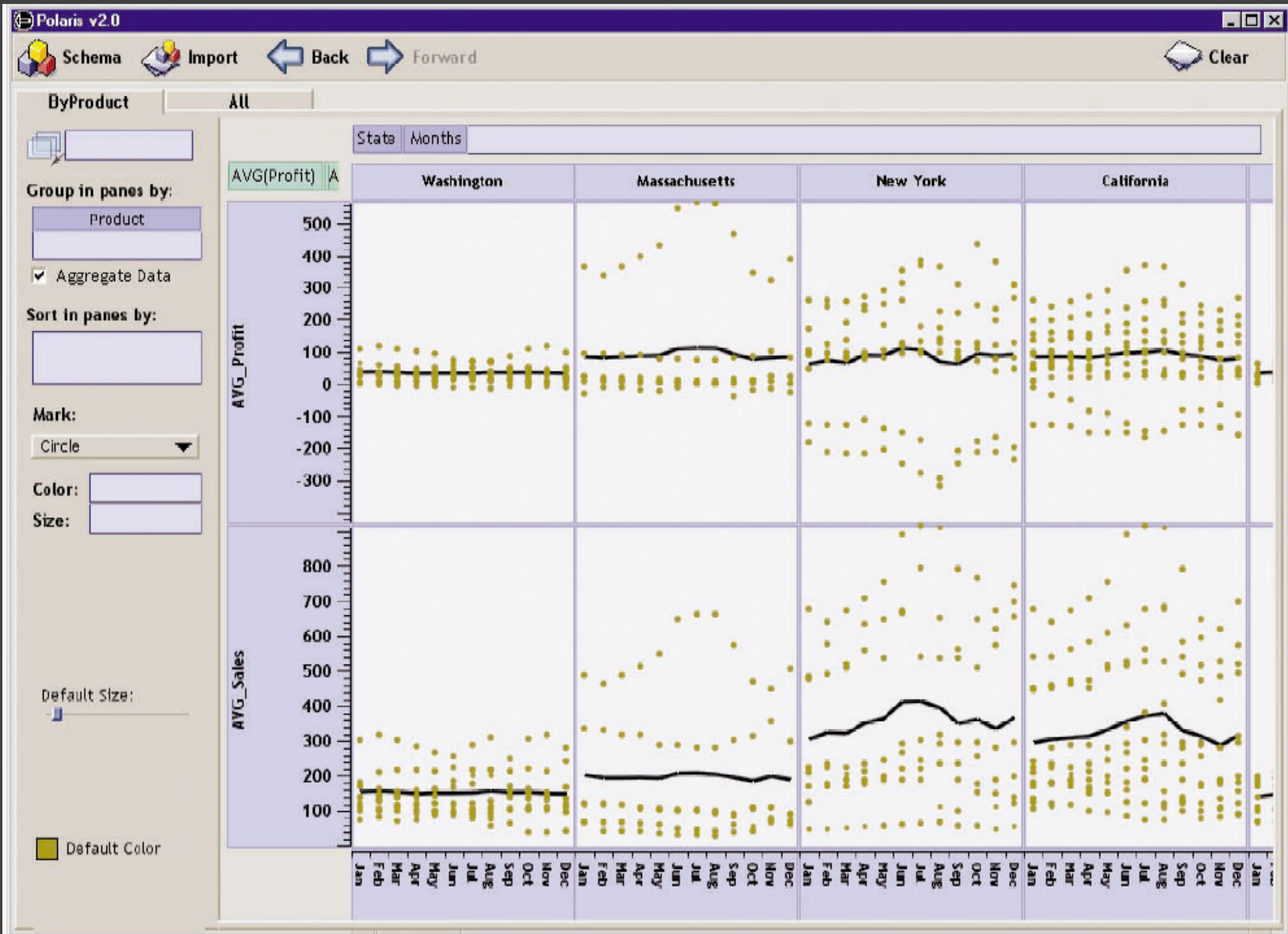
Include titles and captions for each view

Due by 11:59pm
**Monday, Jan 25**

# Tableau / Polaris

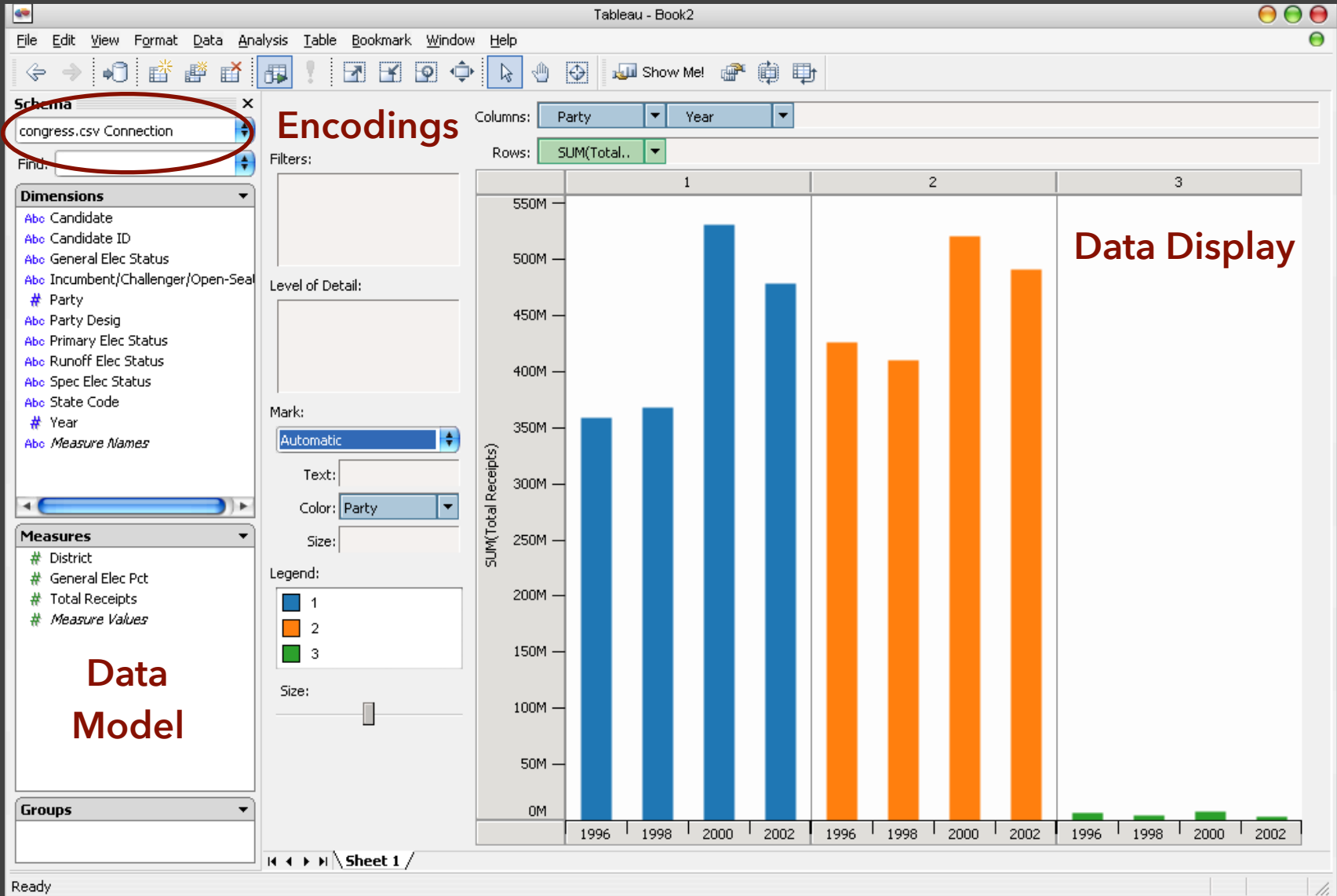# Polaris [Stolte et al.]

# Tableau

# 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

# Tableau Demo

**The dataset:**

Federal Elections Commission Receipts (2012)

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

# 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 (/)

Tableau - Book1

Normal

Show Me

Data | Analytics

Sample – Superstore

Dimensions

- Customer
  - Abc Customer Name
  - Abc Segment
- Order
- Location
- Product
  - Abc Category
  - Abc Sub-Category
  - Manufacturer
  - Abc Product Name
  - Profit (bin)
  - Abc Region
  - Abc Measure Names

Measures

- # Discount
- # Profit
- # Profit Ratio
- # Quantity
- # Sales
- Latitude (generated)
- Longitude (generated)
- Number of Records
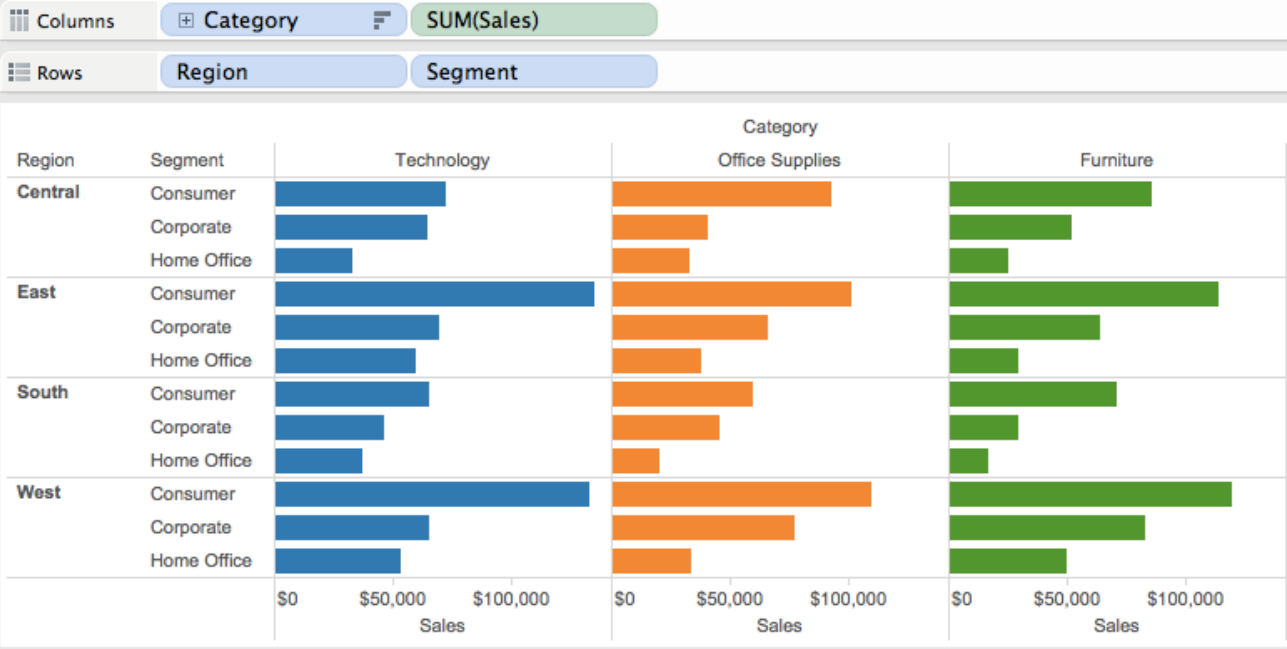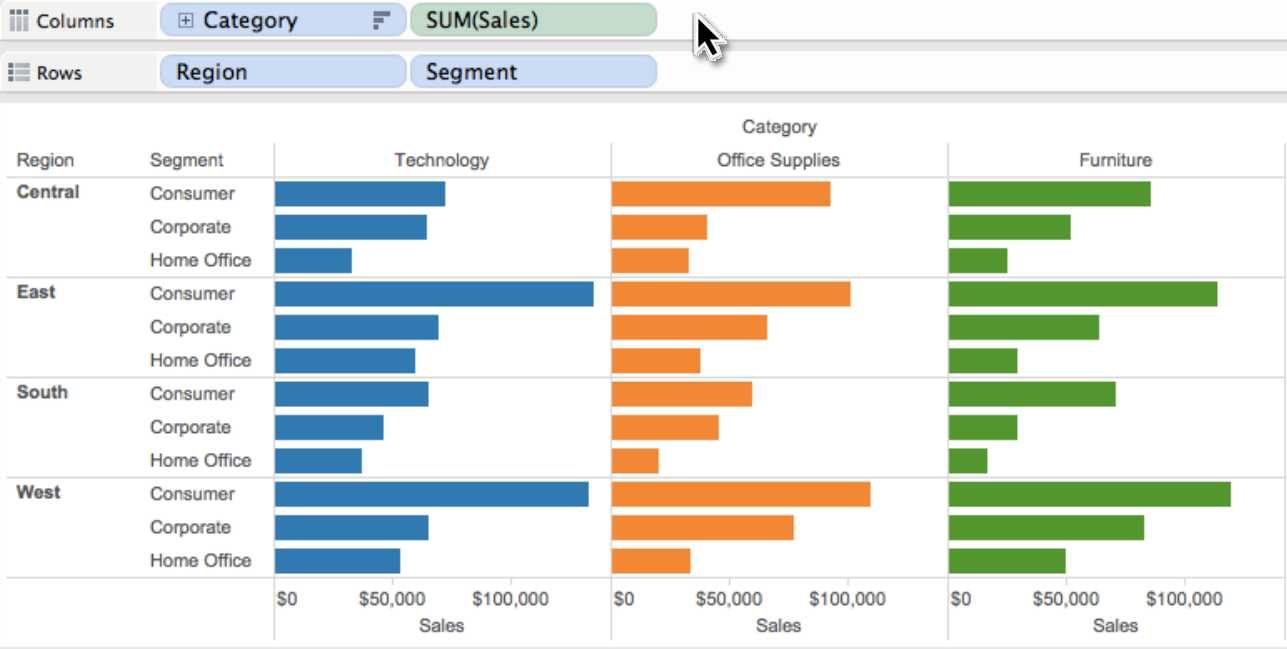- # Measure Values

Pages

Filters

Marks

Automatic

Color | Size | Label

Detail | Tooltip

Category

Category

- Technology
- Office Supplies
- Furniture

Columns: Category | SUM(Sales)
Rows: Region | Segment

Category

| Region | Segment | Technology | Office Supplies | Furniture |
|---|---|---|---|---|
| Central | Consumer | | | |
| | Corporate | | | |
| | Home Office | | | |
| East | Consumer | | | |
| | Corporate | | | |
| | Home Office | | | |
| South | Consumer | | | |
| | Corporate | | | |
| | Home Office | | | |
| West | Consumer | | | |
| | Corporate | | | |
| | Home Office | | | |

$0   $50,000   $100,000   Sales

Data Source | Sheet 1

36 marks    12 rows by 3 columns    SUM(Sales): $2,297,201

# Table Algebra

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

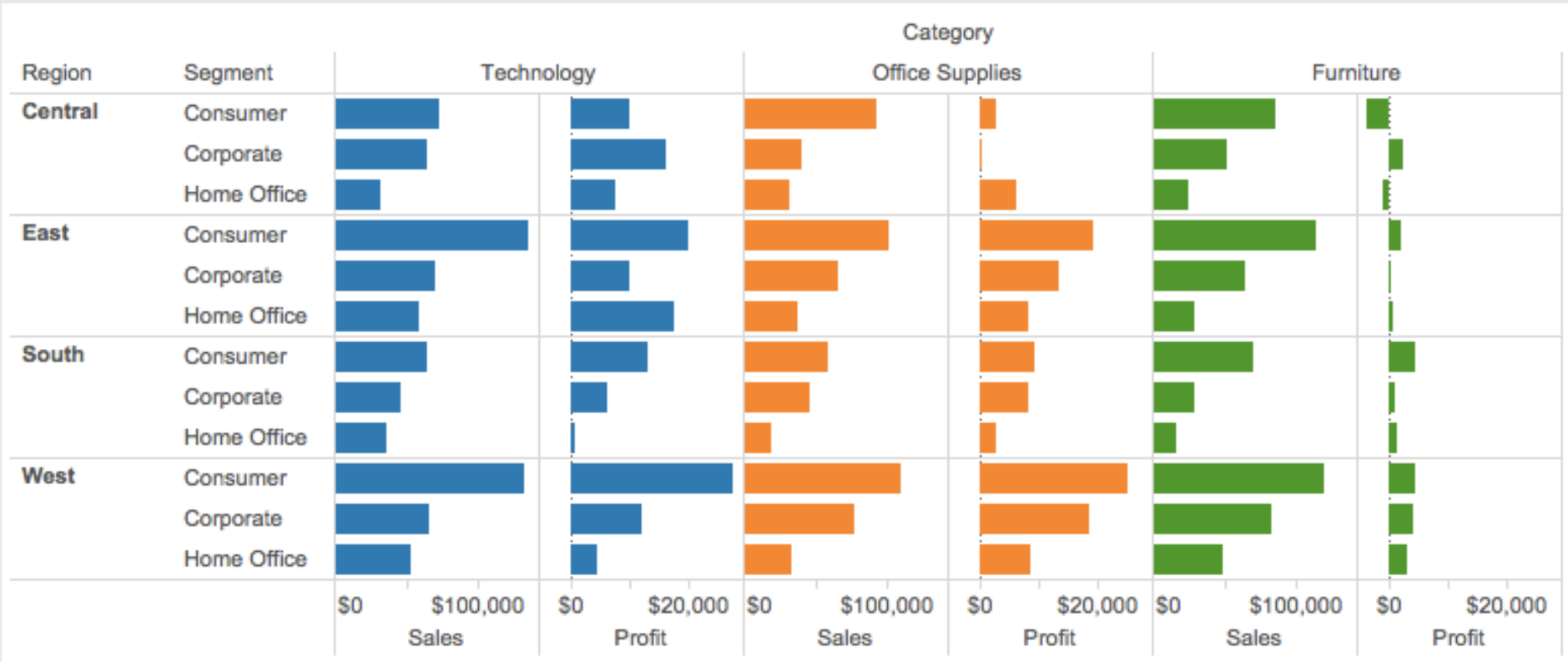Algebraic statements are then mapped to:
**Visualizations** - trellis plot partitions, visual encodings
**Queries** - selection, projection, group-by aggregation

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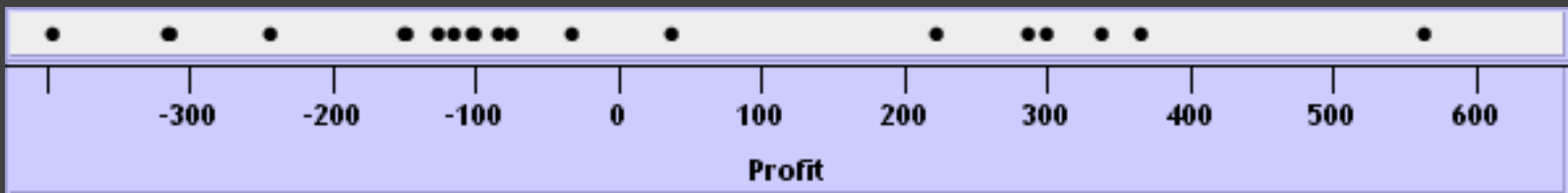
# Table Algebra: Operands

**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
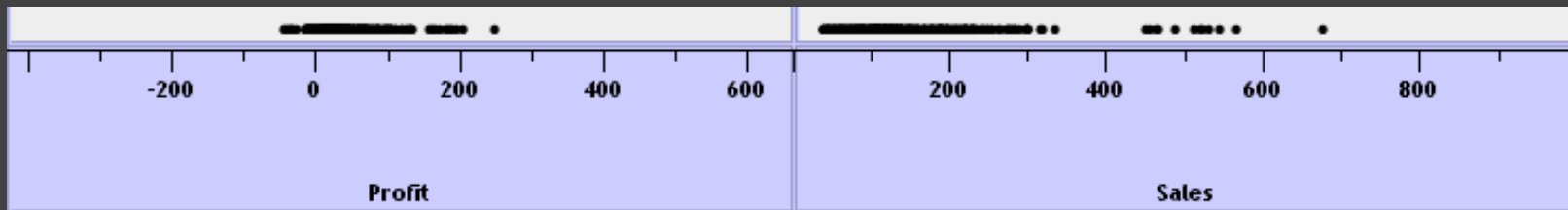
**Ordered union of set interpretations**

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       |

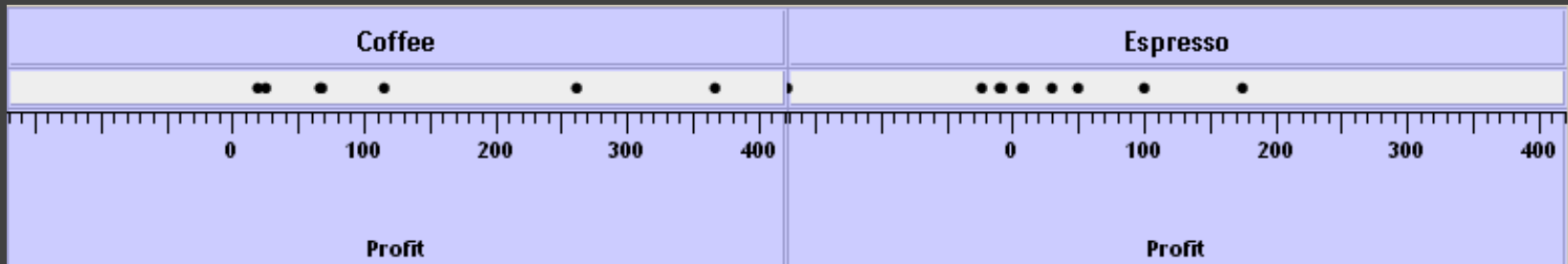Profit + Sales = {(Profit[-310,620]),(Sales[0,1000])}

# 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 ->

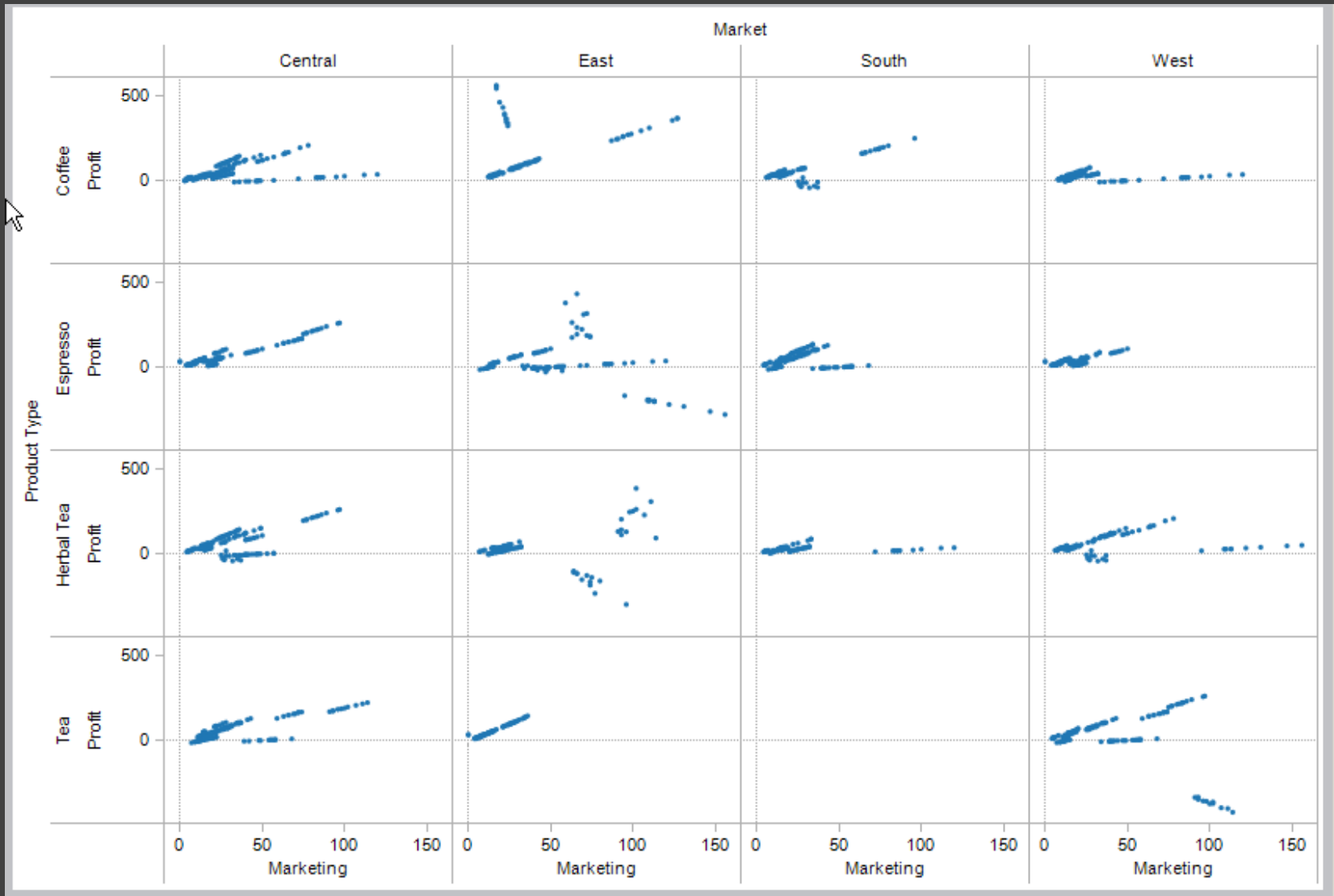creates twelve entries for each quarter. i.e., (Qtr1, December)

Quarter / Month ->

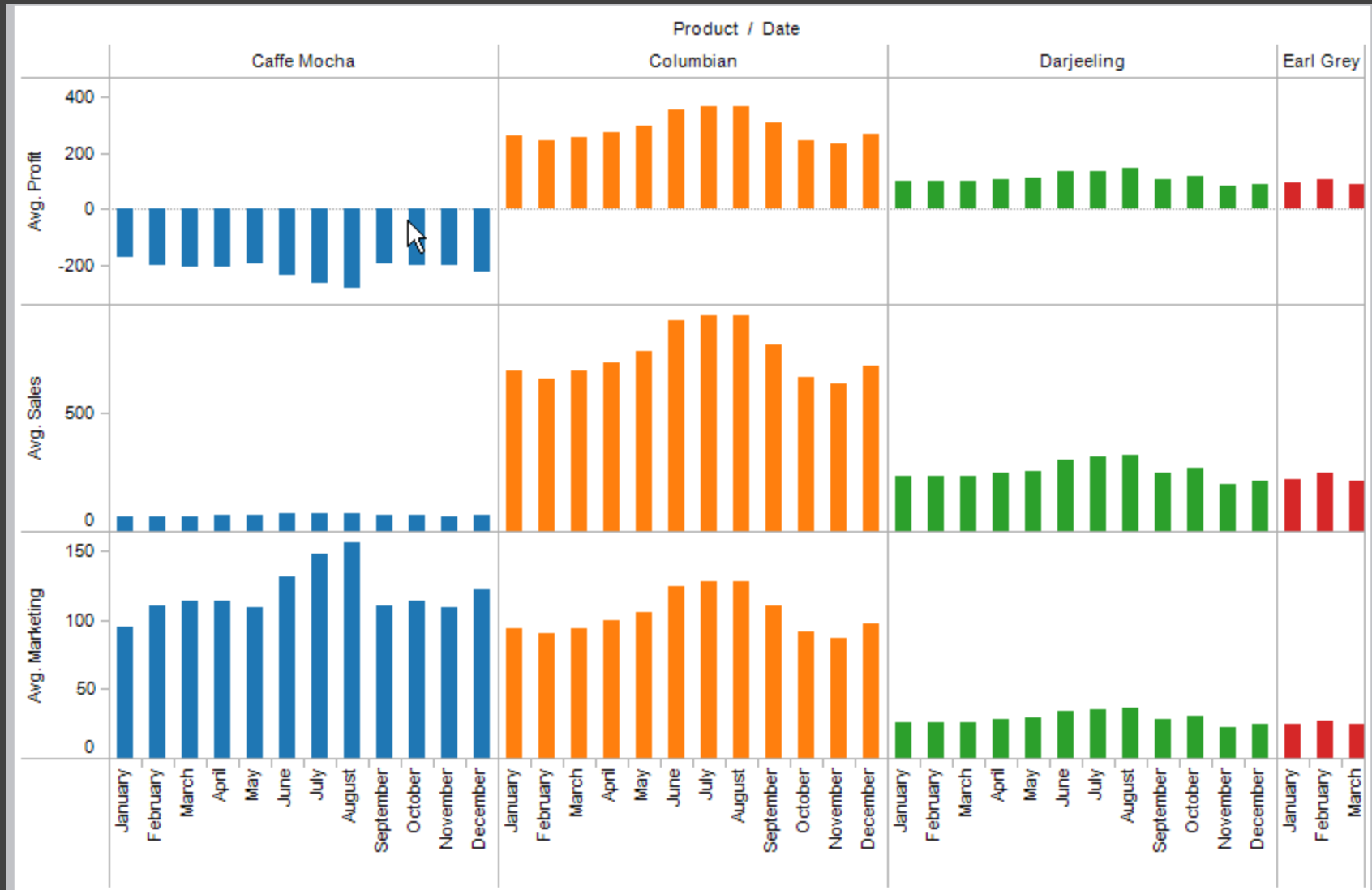creates three entries per quarter based on tuples in database (not semantics)

# 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.

# Quiz Section: Tableau

Tomorrow, Thursday January 14th

Introduction and hands-on experience in Tableau
Come prepared with Tableau installed
See announcement on Ed for instructions

**Up Next:** Jane's Office Hour (link on Canvas)