

# Data Exploration

Dr. Christan Grant

Dr. Laura Melissa Cruz Castro

CAP5771 – Introduction to Data Science

University of Florida



Data Objects and  
Attribute Types



Basic Statistical  
Descriptions of Data



Data Visualization

# Types of Data Sets

## Record (items with attributes)

Relational records

Data matrix, e.g., numerical matrix, crosstabs

Document data: text documents: term-frequency vector

Transaction data

## Graph and network (nodes with relationships)

World Wide Web

Social or information networks

Molecular Structures

## Ordered (sequential records)

Video data: sequence of images

Temporal data: time-series

Sequential Data: transaction sequences

Genetic sequence data

## Spatial, image and multimedia (locations)

Spatial data: maps

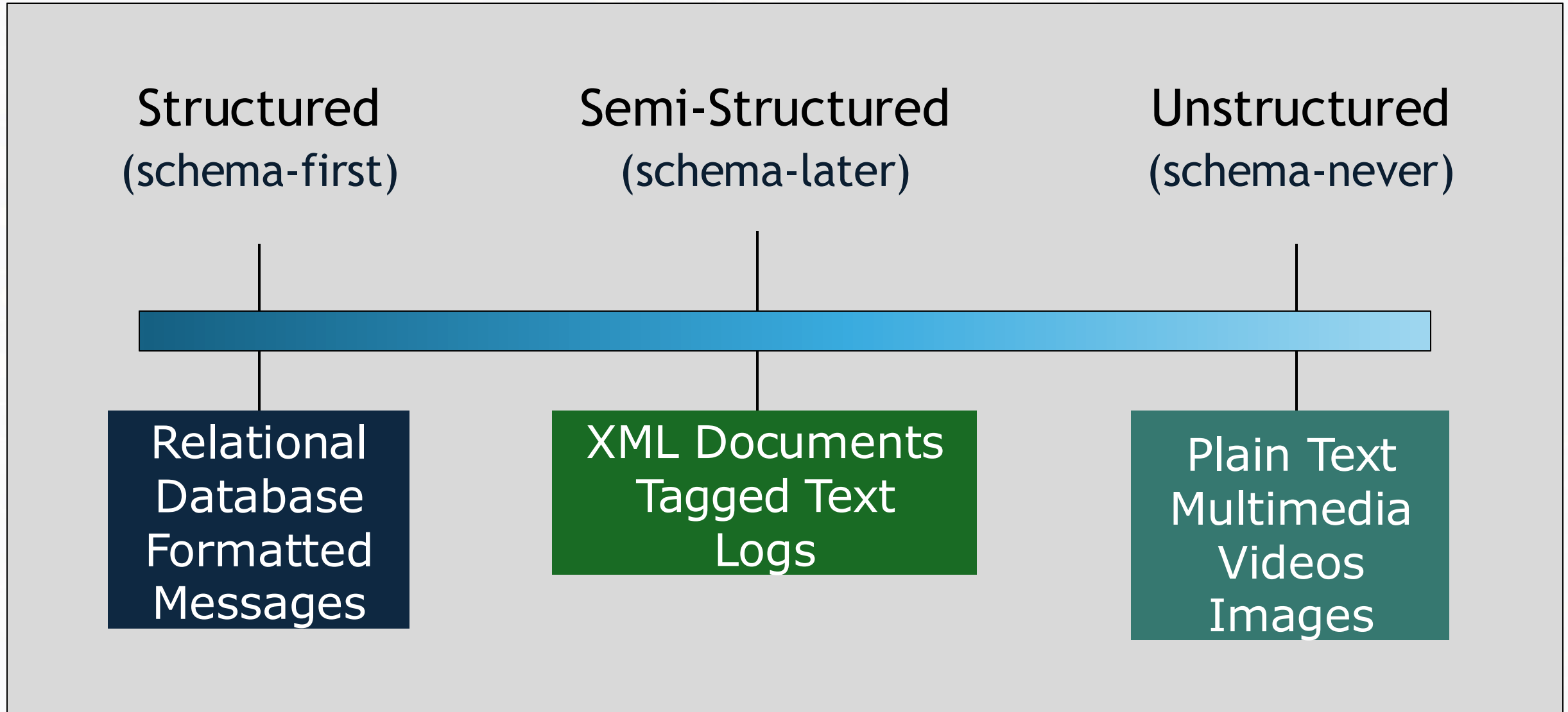
Image data

Video data

	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

# The Structure Spectrum



# Data Sources at Web Companies

## Examples from Facebook

### Structured Data

- Application databases
- Wikipedia (and other knowledge bases)

### Semi-Structured Data

- Web server logs
- Event logs
- API server logs
- Ad server logs
- Search server logs

### Unstructured Data

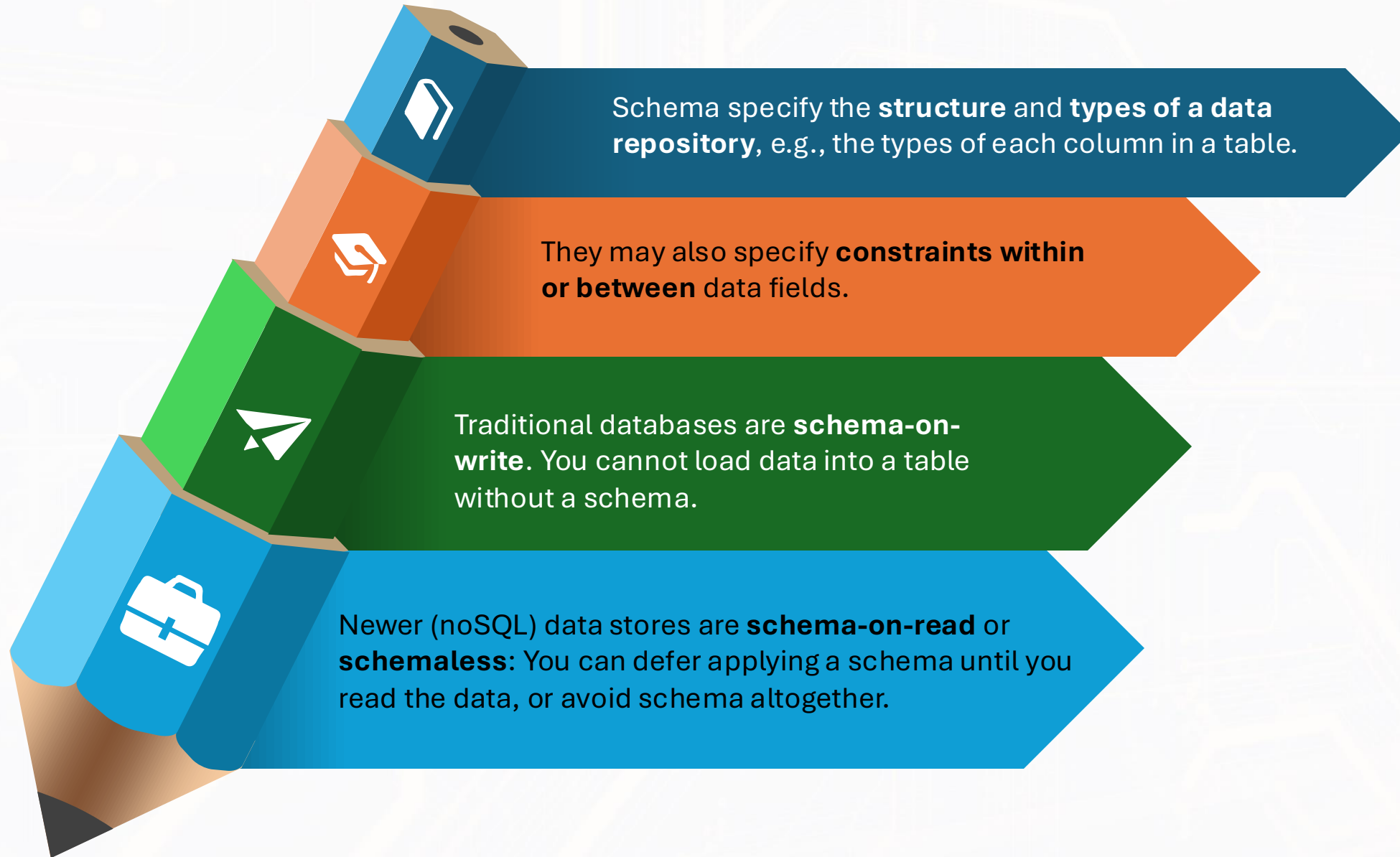
- Advertisement landing page content
- Images and video

# What is a Schema

A **named** collection of tables, views, functions, constraints, indexes, sequences, etc.

- Similar to a namespaces
- The items in the schema define how the data will be treaded.

# Changing Role of Schema



# Schema-on-Write

Predefined data types and structures.

Note: sqlite3 allows you to create a database as a file:

```
>> sqlite3 movies.db
```

```
CREATE SCHEMA hollywood
CREATE TABLE films (title text, release date, awards text[])
CREATE VIEW winners AS
    SELECT title, release FROM films WHERE awards IS NOT NULL;
```

Create

**INSERT INTO** films (text, date, awards) **VALUES**

```
('Spirited Away', 2002, '{"Academy Award", "National Board of Review"}'),
('Major Payne', 1995, NULL);
```

Write



# Schema-on-Read

**XML:** XML schema can be applied later to interpret XML data and specify data types. Here is some XML-encoded data:

```
<location>  
<latitude>37.78333</latitude>  
<longitude>122.4167</longitude>  
</location>
```

# Schema-on-Write vs. Schema-on-Read

## Schema-on-Write

Traditional Approach

## Schema-on-Read

Data is simply copied to the file store, no transformation is needed.

A SerDe (Serializer/Deserializer) is applied during read time to extract the required columns (late binding)

New data can start flowing anytime and will appear retroactively once the SerDe is updated to parse it.

## Pros and Cons

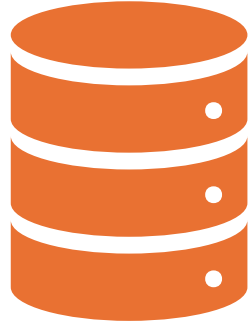
Read is fast

Standards/Governance

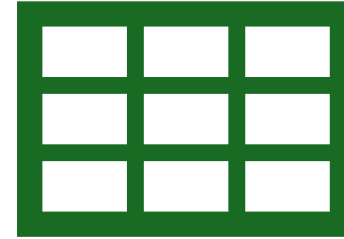
Load is fast

Flexibility/Agility

# Data Model and Schema



A data model is a collection of concepts for describing data.



A schema is a description of a particular collection of data, using a given data model.

# Some Common Data Models

- ⌚ **The relational model of data is the most widely used for record keeping.**

Main concepts: relations, columns/attributes, values

Machine friendly

- ⌚ **Semi-structured models in increasing use (e.g. XML)**

Main concepts: self-describing documents representing tree of labeled values or free text documents

Human friendly

- ⌚ **Others:** RDF Triple, Graph, Streaming, Probabilistic Data, Key-Value, Array/Matrix, Column Stores, Text/Audio/Video

# The Relational Model

## **A Data Model based on Set/Bag Theory**

Support Relational Algebra

## **The Relational Model is Ubiquitous:**

MySQL, PostgreSQL, Oracle, DB2, SQLServ

Foundational work done at

IBM - System R

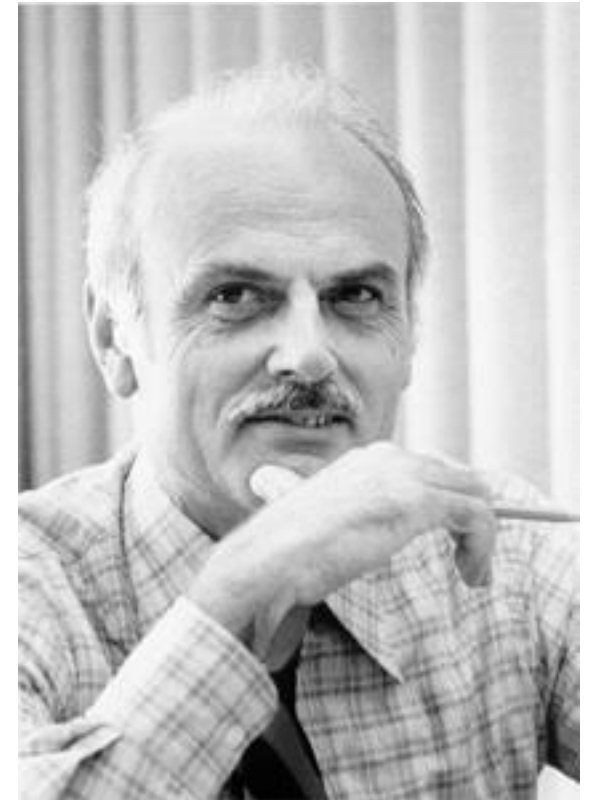
UC Berkeley - Ingres

## **Object-oriented concepts have been merged in**

Early work: POSTGRES research project at Berkeley

Informix, IBM DB2, Oracle 8i

Codd, E. F. (1970). A relational model of data for large shared data banks. *Communications of the ACM*, 13(6), 377-387.



Tedd Codd

# Instance of Students Relation Example

sid	name	login	age	gpa
53666	Jones	jones@cs	18	3.4
53688	Smith	smith@eecs	18	3.2
53650	Smith	smith @math	19	3.8

Cardinality = **3**,  
Arity = **5**,  
all rows distinct

The relation is true for these tuples and false for others  
(*a.k.a, the closed world assumption*)

Arity => The number of distinct attributes in a relation.  
Cardinality => The Number of records in a relation.

# Other Table-Like Data Models: Pandas/Python

**Series:** is a one-dimensional labeled array capable of holding any data type (integers, strings, floating point numbers, Python objects, etc.).  
The axis labels are collectively referred to as the index.

```
>>> s = pd.Series(data, index=index)
```

```
In [3]: s = pd.Series(np.random.randn(5), index=["a", "b", "c", "d", "e"])
```

```
In [4]: s
```

```
Out[4]:
```

```
a    0.469112
b   -0.282863
c   -1.509059
d   -1.135632
e    1.212112
dtype: float64
```

**DataFrame:** a table with named columns  
Represented as a map Dict (col\_name -> series)  
Each Series object represents a column

```
In [37]: d = {
.....:     "one": pd.Series([1.0, 2.0, 3.0], index=["a", "b", "c"]),
.....:     "two": pd.Series([1.0, 2.0, 3.0, 4.0], index=["a", "b", "c", "d"]),
.....: }
.....:
```

```
In [38]: df = pd.DataFrame(d)
```

# Tabular Data in Excel, Google Sheets, Airtable, etc.

	A	B	C	D	E	F	G	H	I
1	rank	company	cik	ticker	sic	state_location	state_of_incorporation	revenues	profits
2	1	Wal-Mart Stores	104169	WMT	5331	AR	DE	421849	16389
3	2	Exxon Mobil	34088	XOM	2911	TX	NJ	354674	30460
4	3	Chevron	93410	CVX	2911	CA	DE	196337	19024
5	4	ConocoPhillips	1163165	COP	2911	TX	DE	184966	11358
6	5	Fannie Mae	310522	FNM	6111	DC	DC	153825	-14014
7	6	General Electric	40545	GE	3600	CT	NY	151628	11644
8	7	Berkshire Hathaway	1067983	BRKA	6331	NE	DE	136185	12967
9	8	General Motors	1467858	GM	3711	MI	MI	135592	6172
10	9	Bank of America Corp.	70858	BAC	6021	NC	DE	134194	-2238
11	10	Ford Motor	37996	F	3711	MI	DE	128954	6561
12	11	Hewlett-Packard	47217	HPQ	3570	CA	DE	126033	8761
13	12	AT&T	732717	T	4813	TX	DE	124629	19864
14	13	J.P. Morgan Chase & Co.	19617	JPM	6021	NY	DE	115475	17370
15	14	Citigroup	831001	C	6021	NY	DE	111055	10602
16	15	McKesson	927653	MCK	5122	CA	DE	108702	1263
17	16	Verizon Communications	732712	VZ	4813	NY	DE	106565	2549
18	17	American International Group	5272	AIG	6331	NY	DE	104417	7786
19	18	International Business Machines	51143	IBM	3570	NY	NY	99870	14833
20	19	Cardinal Health	721371	CAH	5122	OH	OH	98601.9	642.2
21	20	Freddie Mac	37785	FMC	2800	PA	DE	98368	-14025



# Log Files – Example Apache Web Log

Processes, usually daemons, create logs

e.g., httpd, mysqld, syslogd

```
66.249.65.107 - - [08/Oct/2007:04:54:20 -0400] "GET /support.html
HTTP/1.1" 200 11179 "-" "Mozilla/5.0 (compatible; Googlebot/2.1;
+http://www.google.com/bot.html)"
```

```
111.111.111.111 - - [08/Oct/2007:11:17:55 -0400] "GET / HTTP/1.1"
200 10801 "http://www.google.com/search?q=log+analyzer&ie=utf-8&oe=utf-8
&aq=t&rls=org.mozilla:en-US:official&client=firefox-a" "Mozilla/5.0
(Windows; U; Windows NT 5.2; en-US; rv:1.8.1.7) Gecko/20070914
Firefox/2.0.0.7"
```

```
111.111.111.111 - - [08/Oct/2007:11:17:55 -0400] "GET /style.css
HTTP/1.1" 200 3225 ""http://www.loganalyzer.net/" "Mozilla/5.0 (Windows;
U; Windows NT 5.2; en-US; rv:1.8.1.7) Gecko/20070914
Firefox/2.0.0.7"
```

# Well-Formed XML Example

<?xml version = "1.0" standalone = "yes" ?>

<BARS>

<BAR> <NAME>Joe's Bar</NAME>

<BEER> <NAME>Bud</NAME>

<PRICE>2.50</PRICE> </BEER>

<BEER> <NAME>Miller</NAME>

<PRICE>3.00</PRICE> </BEER>

</BAR>

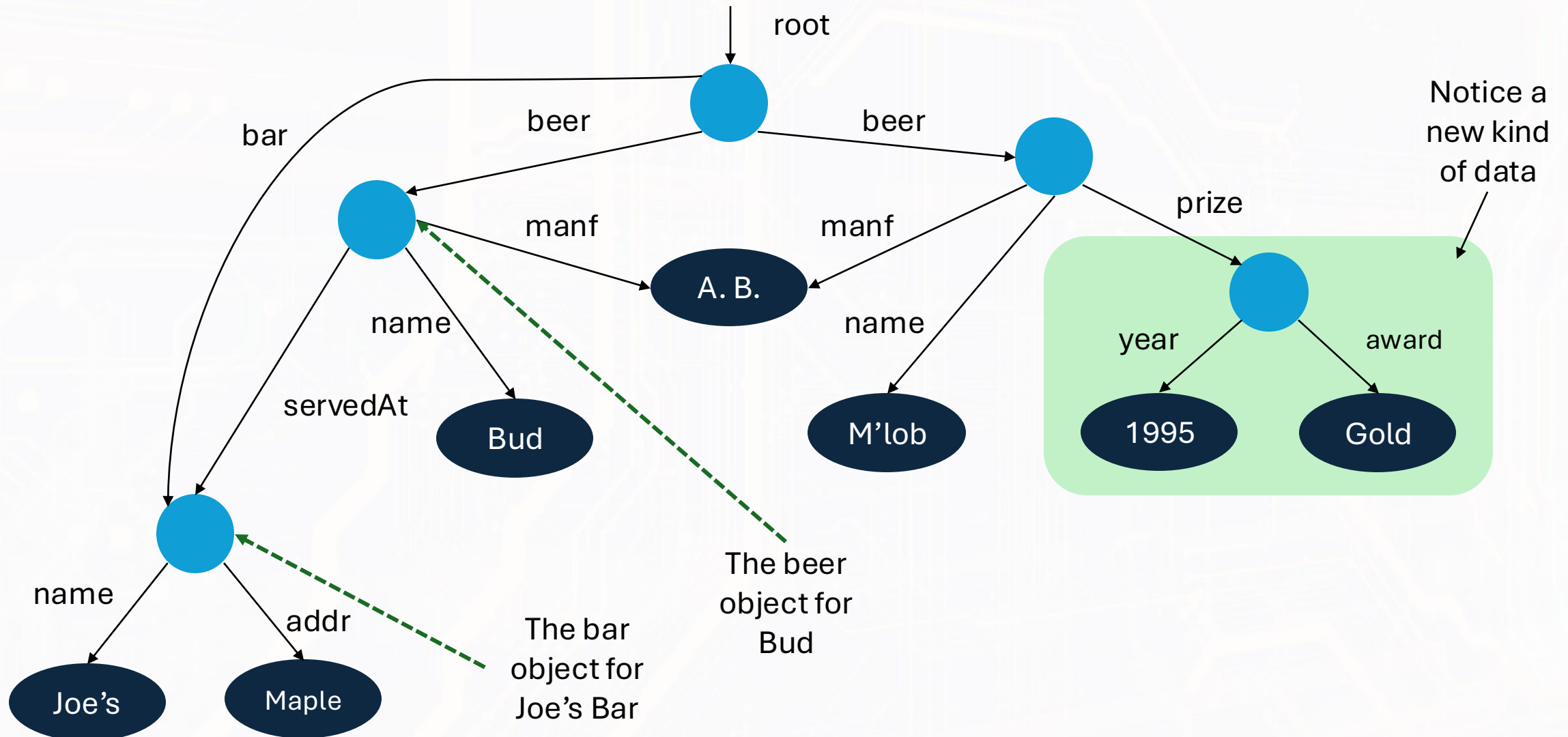
<BAR> ...

</BARS>

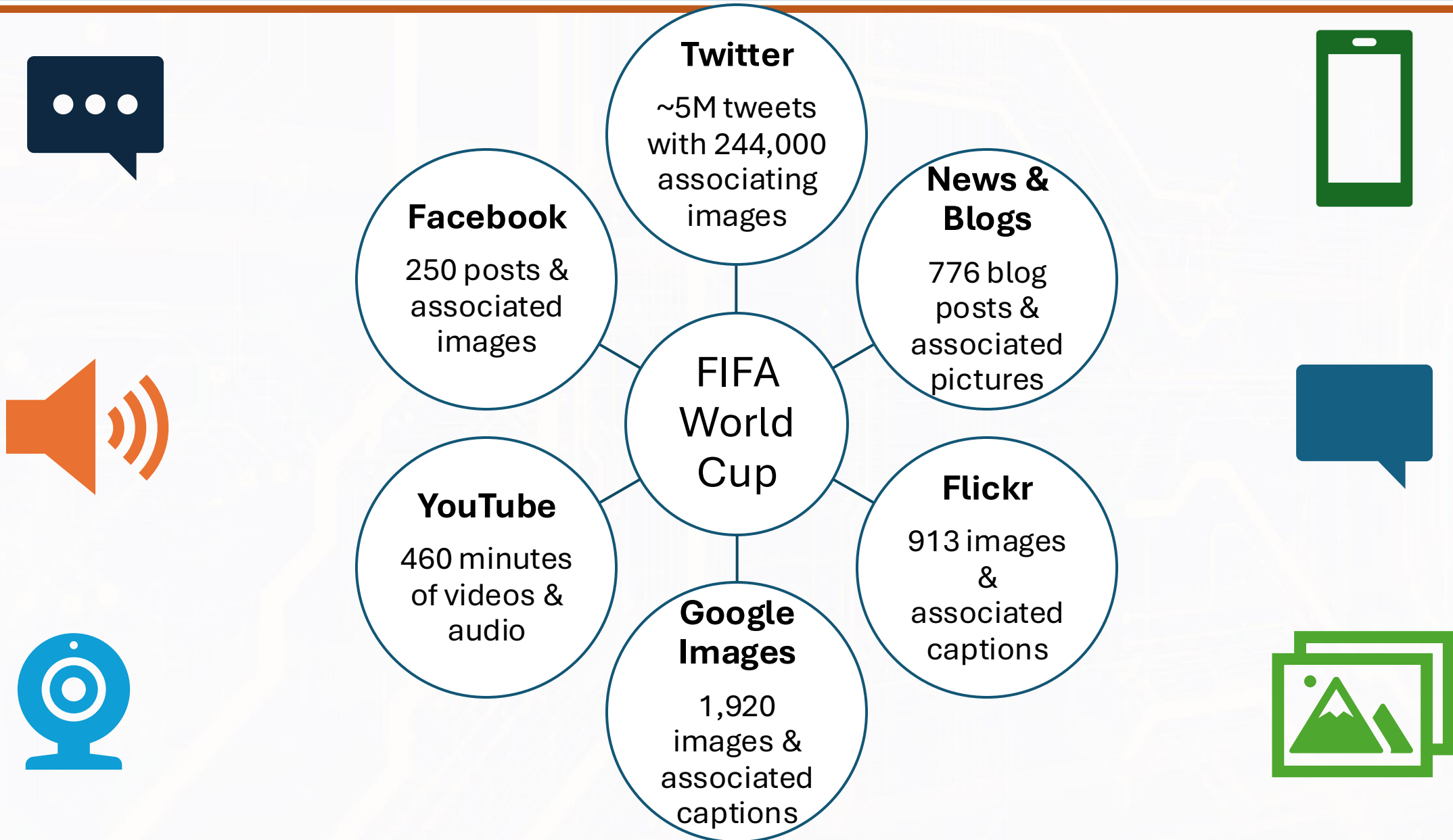
A NAME  
subobject

A BEER  
subobject

# Data Tree Example



# Multimodal Data Sources on World Cup Common

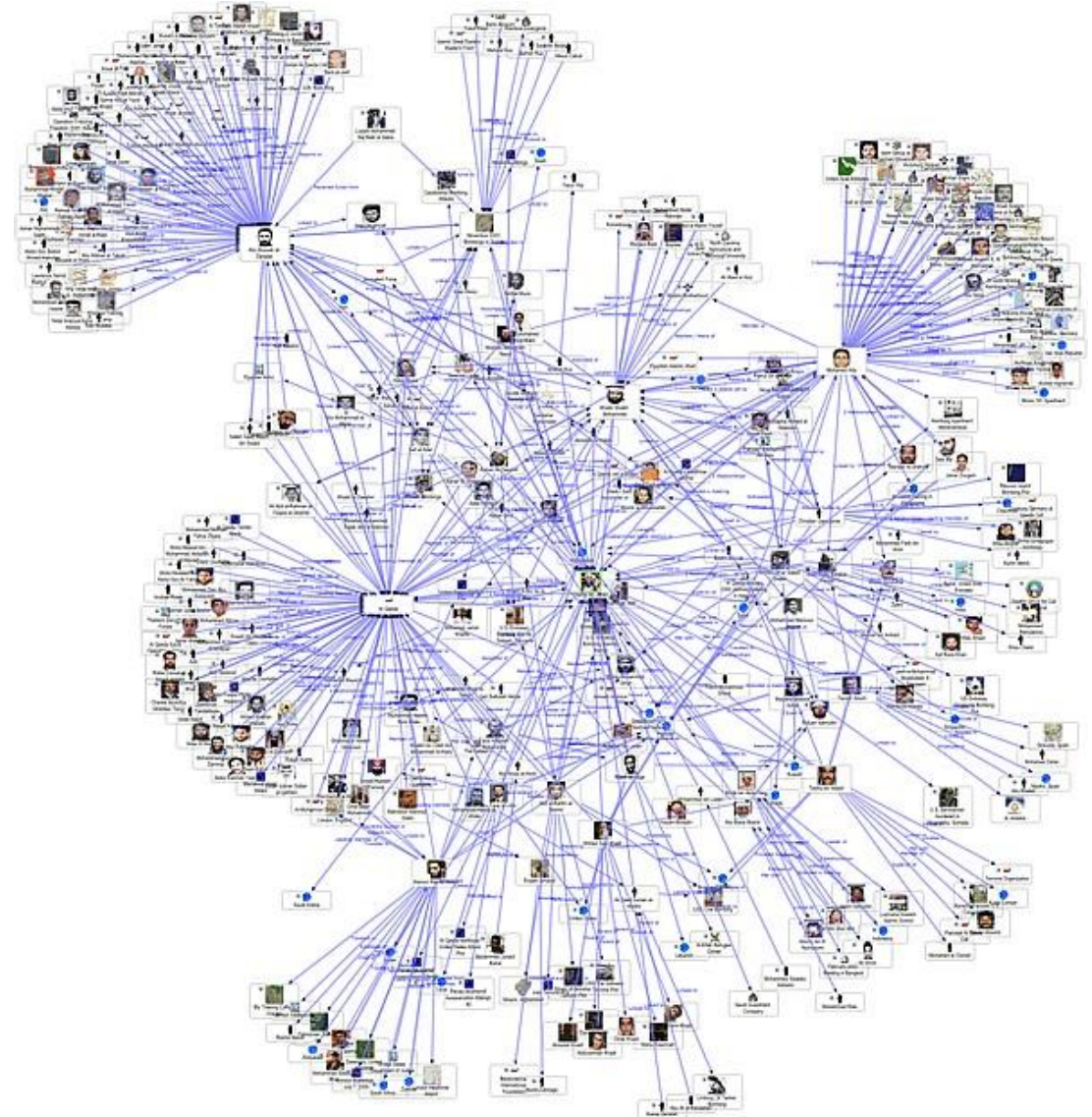


# Graph Data

Lots of interesting data has a graph structure:

- ❖ Social networks
- ❖ Communication networks
- ❖ Computer Networks
- ❖ Road networks
- ❖ Citations
- ❖ Collaborations/Relationships
- ❖ ...

Some of these graphs can get quite large (e.g., Facebook\* user graph)



# Data Objects

Data sets are made up of data objects.

A **data object** represents an entity.

Examples:  
sales database:  
customers, store items, sales  
medical database:  
patients, treatments  
university database:  
students, professors, courses

Database rows → data objects  
Database columns → attributes

Also called  
samples,  
examples,  
instances, data  
points, objects,  
tuples.

Data objects  
are described  
by attributes.

# Attributes

**Attribute** (or dimensions, features, variables): a data field, representing a characteristic or feature of a data object.

E.g., customer\_ID, name, address

Types:

- Nominal

- Binary

- Numeric: quantitative

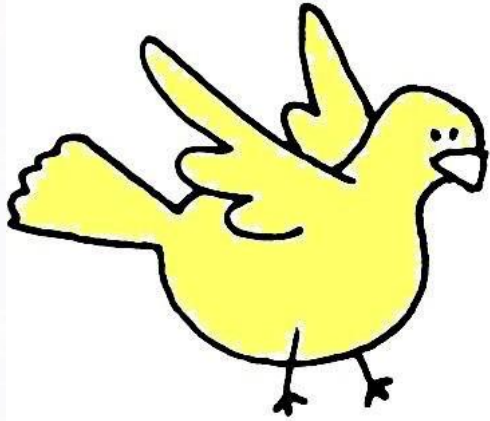
- Interval-scaled

- Ratio-scaled

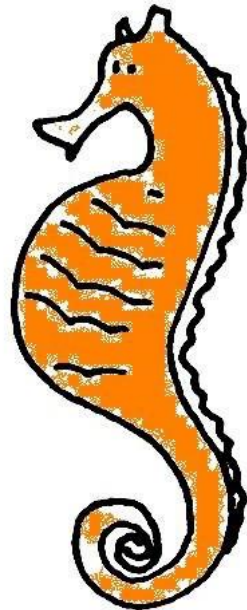


# Descriptive Data

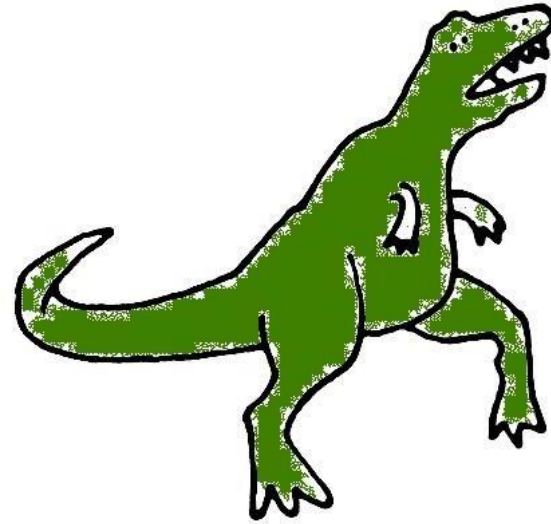
## CATEGORICAL DATA:



I am a bird.  
I am yellow.  
I am awesome.



I am a seahorse.  
I am orange.  
I am super awesome.



I am a T-rex.  
I am green.  
I am extinct.



# Attribute Types

## Nominal

- Categories, states, or “names of things”
- Hair\_color = {auburn, black, blond, brown, grey, red, white}
- marital status, occupation, ID numbers, zip codes

## Binary

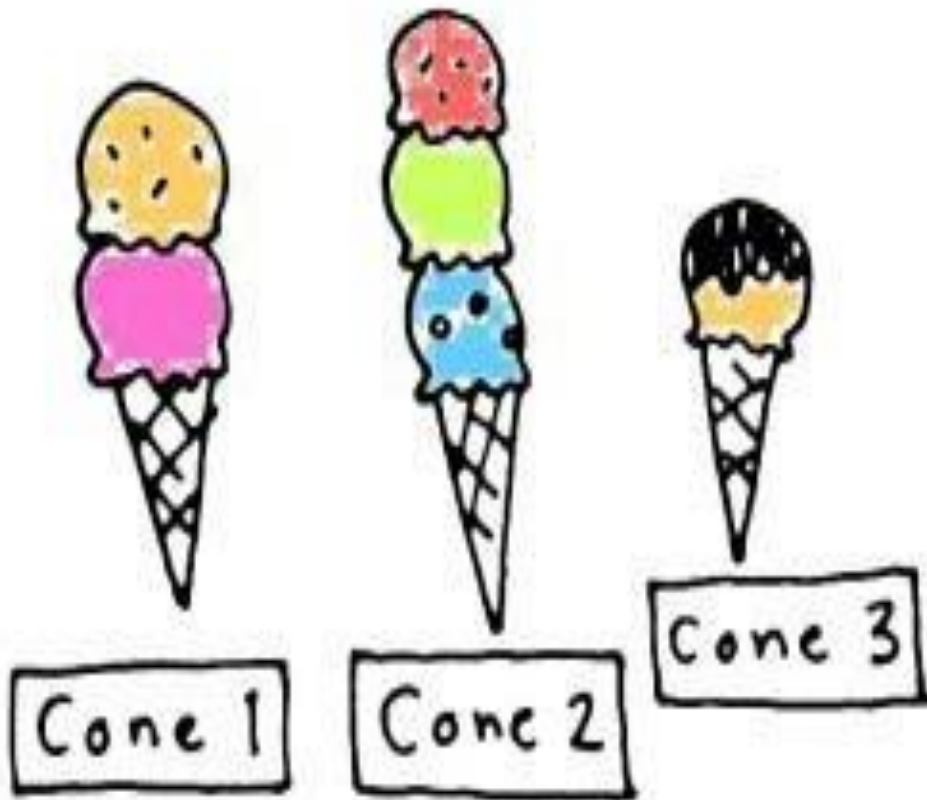
- Nominal attribute with only 2 states
- Symmetric binary: both outcomes equally important
  - e.g., gender
- Asymmetric binary: outcomes not equally important.
  - e.g., medical test (positive vs. negative)
- Convention: assign 1 to most important outcome (e.g., HIV positive)

## Ordinal

- Values have a meaningful order (ranking) but magnitude between successive values is not known.
- Size = {small, medium, large}, grades, army rankings

Data – It's Numeric

# QUANTITATIVE DATA:



Discrete data:

- There are 3 cones
- Cone 1 has 2 scoops

Continuous data:

- Cone 3 weighs 79.4 grams
- cone 2 ice cream is at 8.3°F

# Discrete vs. Continuous Attributes

## Discrete Attribute

- Has only a finite or countably infinite set of values
- E.g., zip codes, profession, or the set of words in a collection of documents
- Sometimes, represented as integer variables
- Note: Binary attributes with 0/1 values are a special case of discrete attributes

## Continuous Attribute

- Has real numbers as attribute values
- E.g., temperature, height, or weight
- Practically, real values can only be measured and represented using a finite number of digits
- Continuous attributes are typically represented as floating-point variables

# Numeric Attribute Types

**Quantity** (integer or real-valued)

## Interval

- Measured on a scale of equal-sized units
- Values have order
- E.g., temperature in  $^{\circ}\text{C}$  or  $^{\circ}\text{F}$ , calendar dates
- No true zero-point

## Ratio

- Inherent zero-point
- We can speak of values as being an order of magnitude larger than the unit of measurement ( $10\text{ K}^{\circ}$  is twice as high as  $5\text{ K}^{\circ}$ ).
- E.g., temperature in Kelvin, length, counts, monetary quantities



Data Objects and  
Attribute Types



Basic Statistical  
Descriptions of Data



Data Visualization

# Basic Statistical Descriptions of Data

## Motivation

To better understand the data/distribution

## Measuring the central tendencies

Mean, median, mode, etc.

## Data dispersion characteristics

Max, min, quantiles, outliers, variance, etc.

# Measuring the Central Tendency



## Mean

- Sample vs. Population. Note: n is sample size and N is population size.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \mu = \frac{\sum x}{N}$$

- Weighted arithmetic mean:
$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$
- Trimmed mean
  - (Remove a %age)

## Median

- Middle value if odd number of values, or average of the middle two values otherwise
- Estimated by interpolation (for grouped data)

## Mode

- Value that occurs most frequently in the data
- Unimodal, bimodal, trimodal

age	frequency
1-5	200
6-15	450
16-20	300
21-50	1500
51-80	700
81-110	44

# Measuring the Central Tendency: (2) Median

## Median:

Middle value if odd number of values, or average of the middle two values otherwise

- Estimated by interpolation (for *grouped data*):

<i>age</i>	<i>frequency</i>
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44

Approximate median

Sum before the median interval

Interval width ( $L_2 - L_1$ )

$$median = L_1 + \left( \frac{n/2 - (\sum freq)_l}{freq_{median}} \right) width$$

Low interval limit

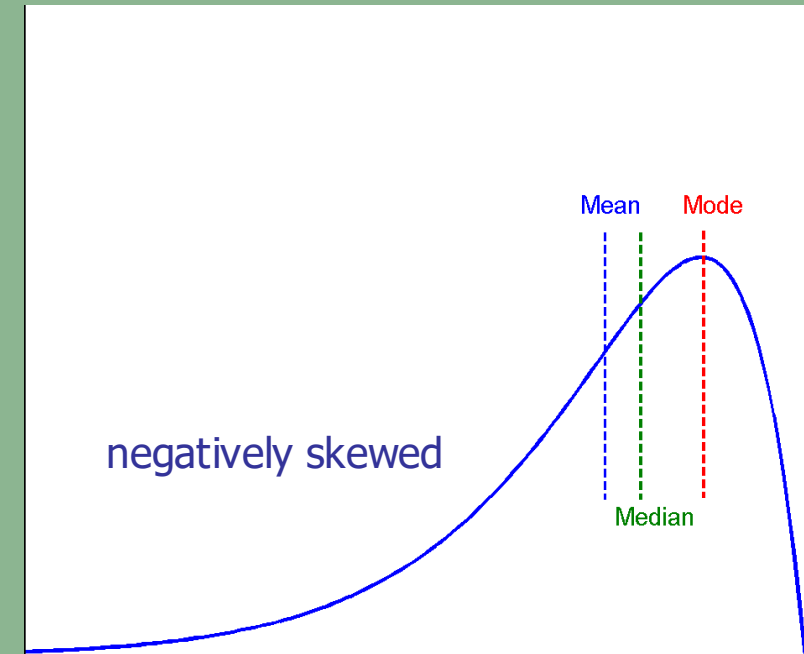
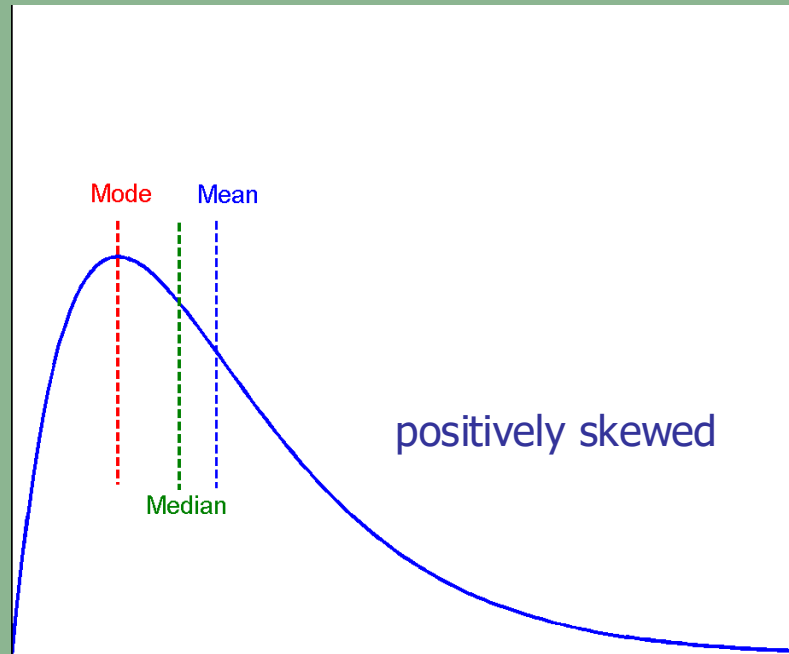
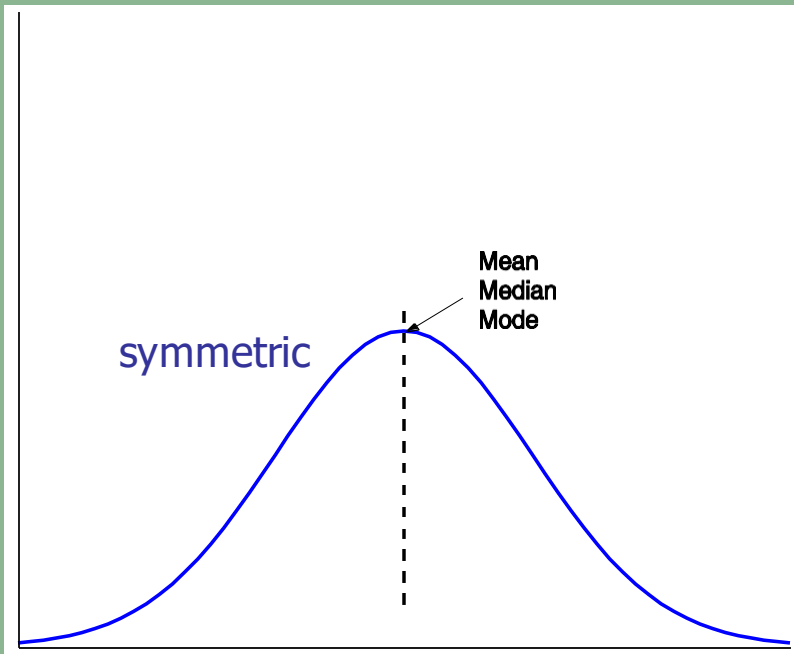


# Statistics on Different Types of Variables

OK to compute....	Nominal	Ordinal	Interval	Ratio
frequency distribution	Yes	Yes	Yes	Yes
median and percentiles	No	Yes	Yes	Yes
add or subtract	No	No	Yes	Yes
mean, standard deviation, standard error of the mean	No	No	Yes	Yes
ratio, or coefficient of variation	No	No	No	Yes

# Symmetric vs. Skewed Data

## Median, Mean and Mode



# Measuring the Dispersion of Data

## Quartiles and outliers

Quartiles:  $Q_1$  (25th percentile),  $Q_3$  (75th percentile)

Inter-quartile range:  $IQR = Q_3 - Q_1$

Five number summary: min,  $Q_1$ , median,  $Q_3$ , max

Outlier: usually, a value higher/lower than  $1.5 \times IQR$

## Variance and standard deviation -

Sample vs. population standard deviation  $s$  vs.  $\sigma$

### Variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right] \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

Standard deviation  $s$  (or  $\sigma$ ) is the square root of variance  $s^2$  (or  $\sigma^2$ )

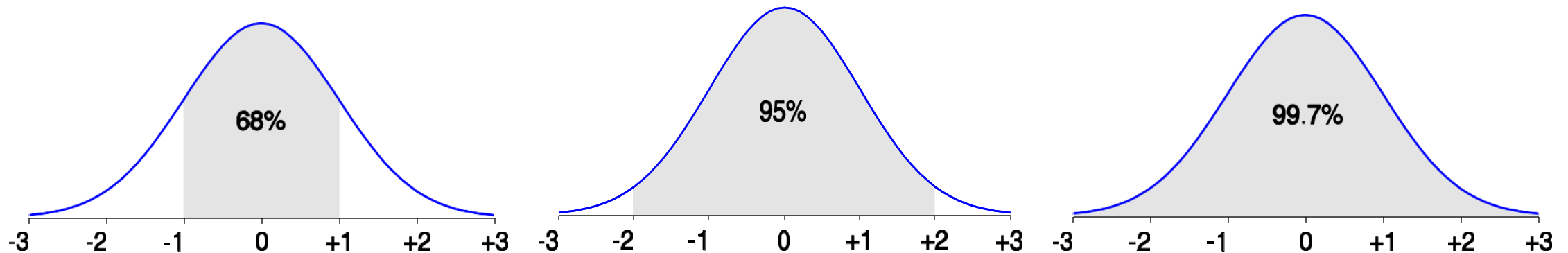
# Properties of Normal Distribution Curve

## The normal (distribution) curve

From  $\mu - \sigma$  to  $\mu + \sigma$ : contains about 68% of the measurements ( $\mu$ : mean,  $\sigma$ : standard deviation)

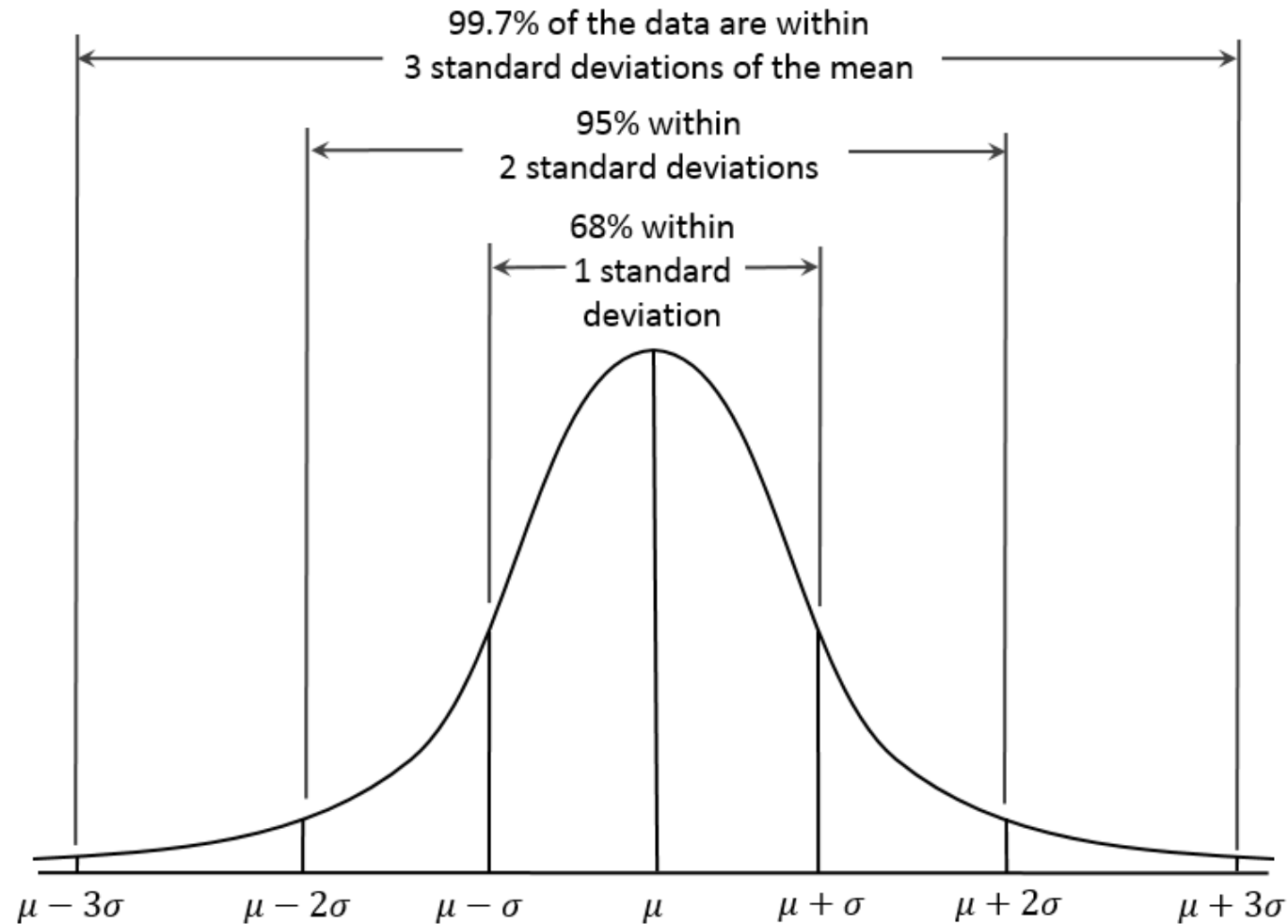
From  $\mu - 2\sigma$  to  $\mu + 2\sigma$ : contains about 95% of it

From  $\mu - 3\sigma$  to  $\mu + 3\sigma$ : contains about 99.7% of it



# Properties of Normal Distribution Curve

← — — — — — Represent data dispersion, spread — — — — — →

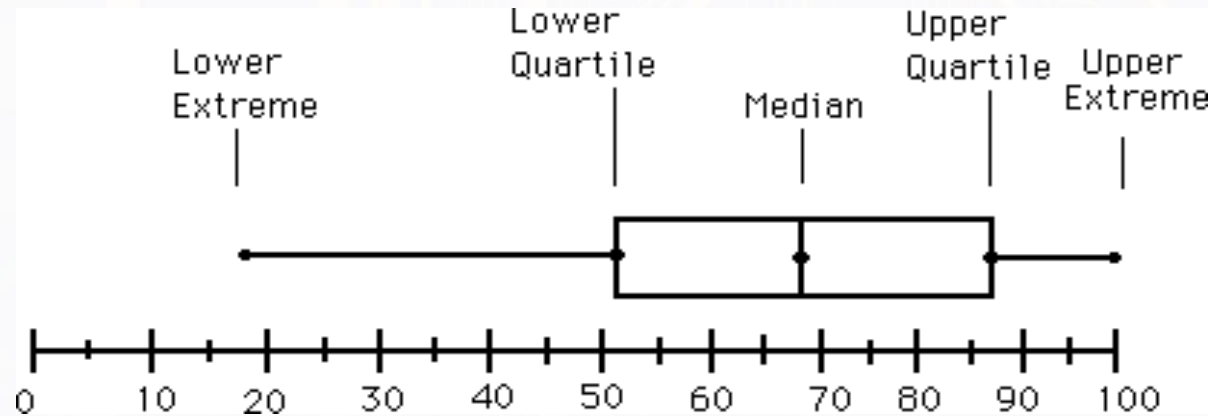


Represent central tendency

# Boxplot Analysis

**Five-number summary of a distribution**

Minimum, Q1, Median, Q3, Maximum



**Boxplot**

Data is represented with a box

The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR

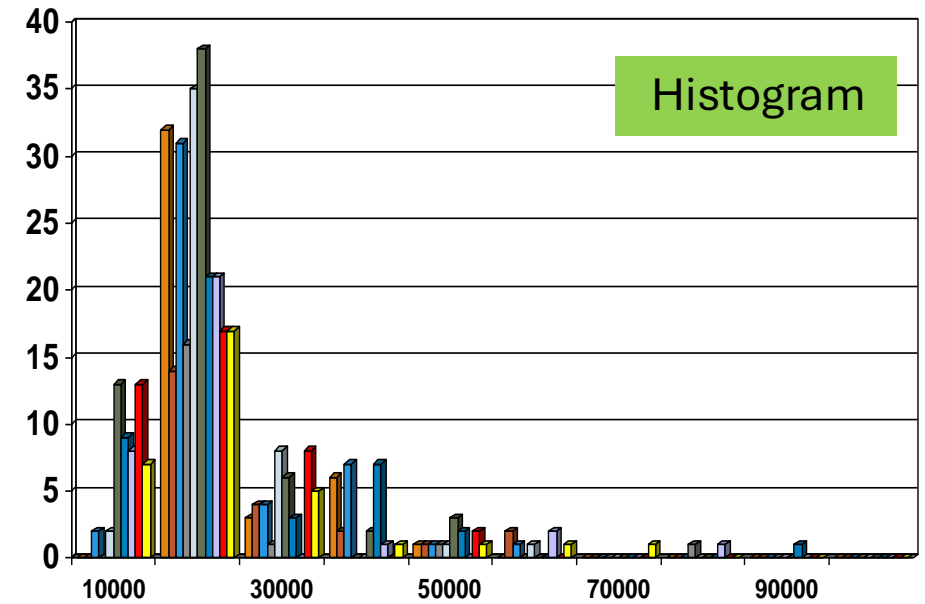
The median is marked by a line within the box

Whiskers: two lines outside the box extended to Minimum and Maximum

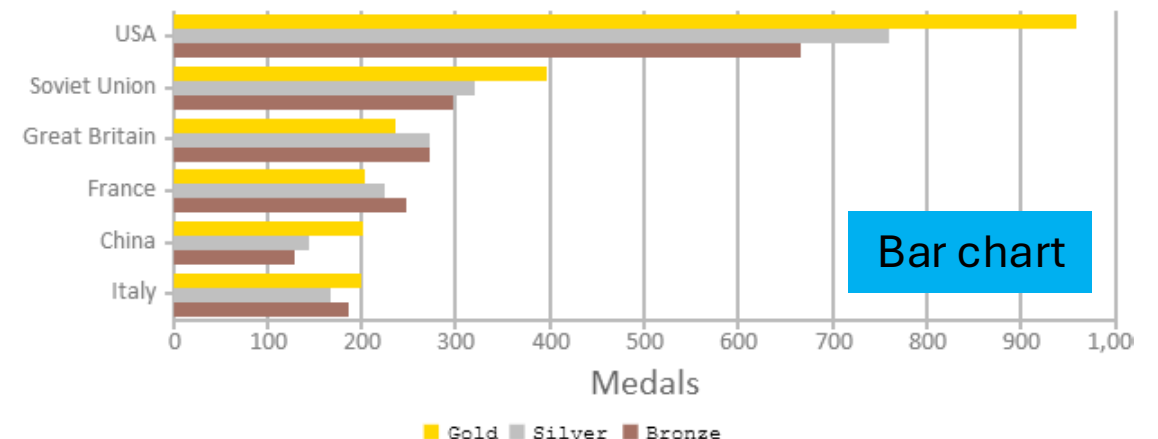
Outliers: points beyond a specified outlier threshold, plotted individually

# Histogram Analysis

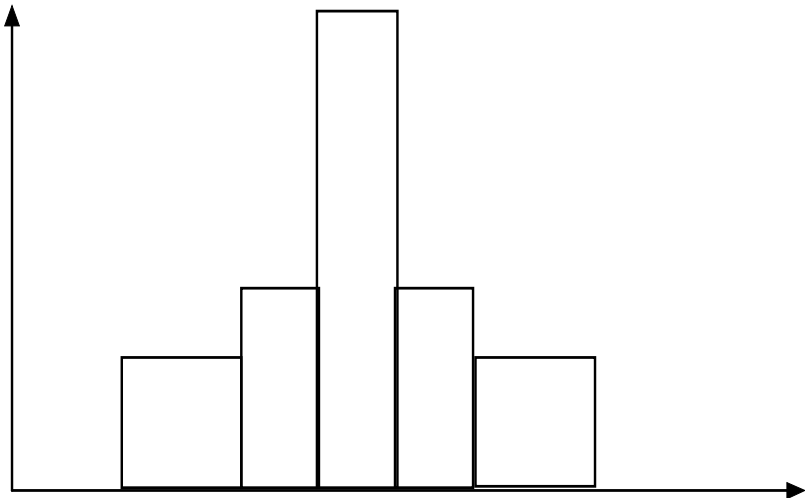
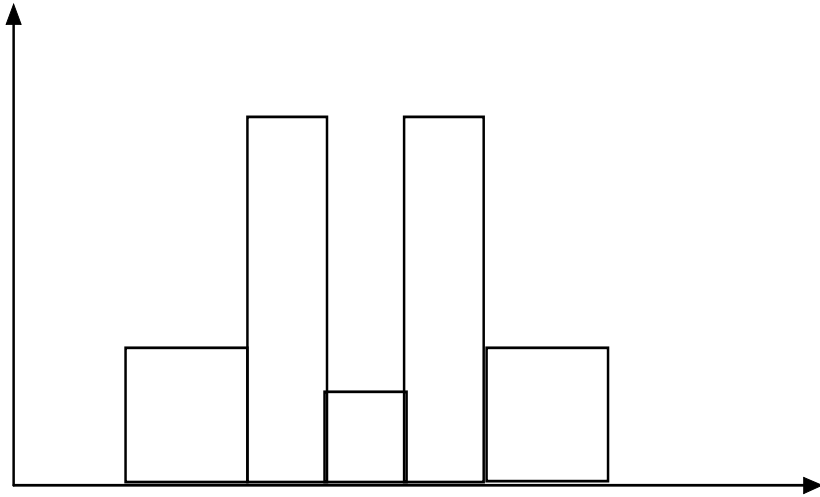
- Histogram: Graph display of tabulated frequencies, shown as bars
- Differences between histograms and bar charts
  - Histograms are used to show **distributions of variables** while bar charts are used to **compare variables**
  - Histograms **plot binned quantitative data** while bar charts **plot categorical data**
  - Bars can be **reordered** in bar charts but not in histograms
  - Differs from a bar chart in that it is **the area of the bar that denotes the value**, not the height as in bar charts, a crucial distinction when the categories are not of uniform width



Olympic Medals of all Times (till 2012 Olympics)



# Histograms Often Tell More than Boxplots



- The two histograms shown in the left may have the same boxplot representation
- The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions

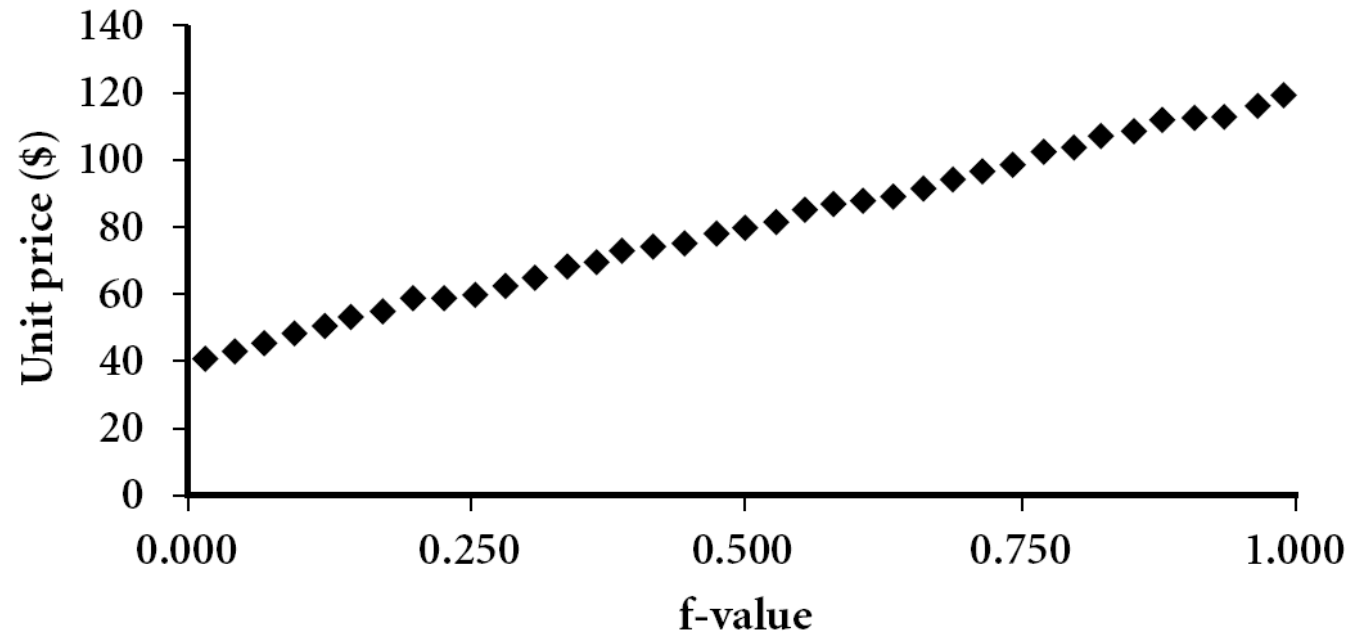


# Quantile Plot

▮ Displays all of the data  
(allowing the user to assess  
both the overall behavior  
and unusual occurrences)

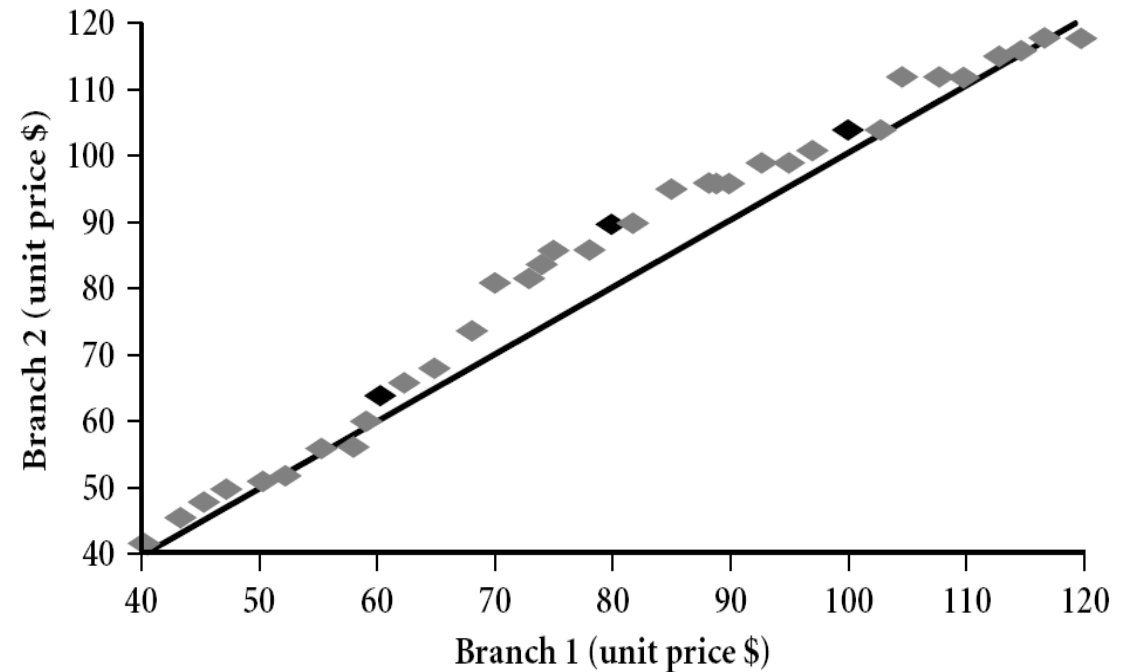
▮ Plots quantile information

▮ For a data  $x_i$  data **sorted in increasing order**,  $f_i$  indicates that approximately 100  $f_i$ % of the data are below or equal to the value  $x_i$



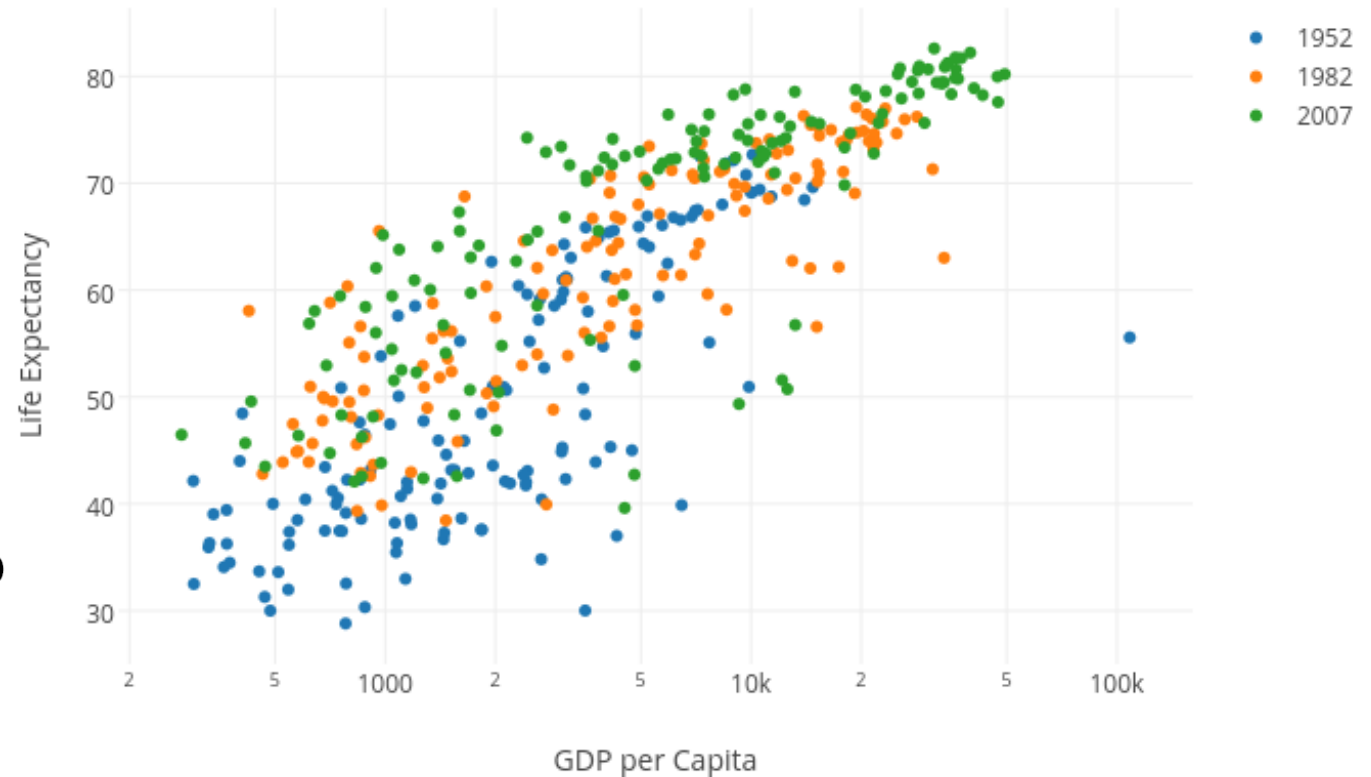
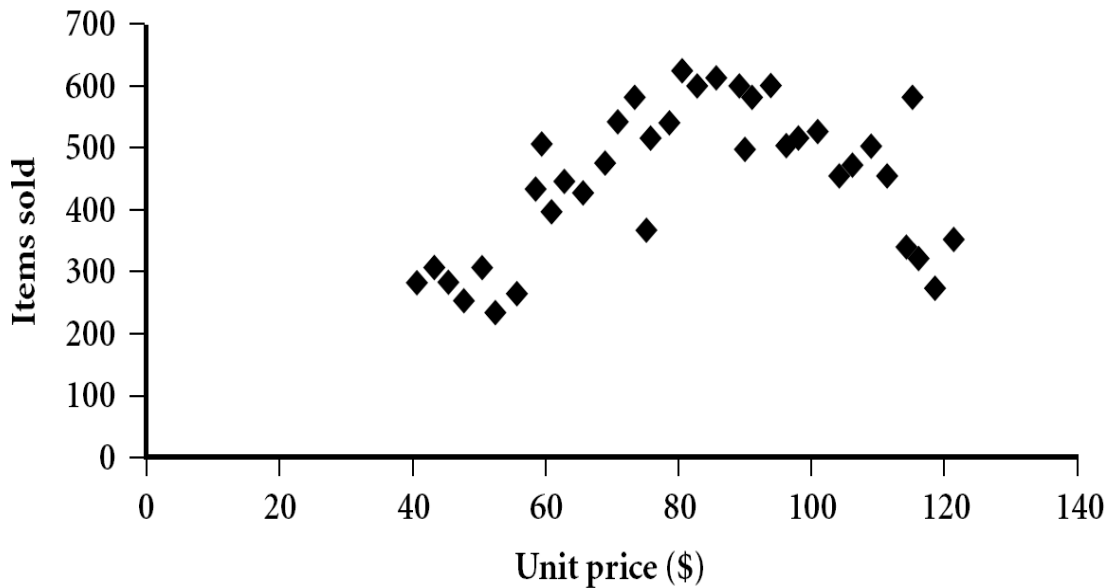
# Quantile-Quantile (Q-Q) Plot

- 📊 Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- 📊 View: Is there is a shift in going from one distribution to another?
- 📊 Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.



# Scatter plot

- Provides a first look at bivariate data to see clusters of points, outliers, etc.
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane





Data Objects and  
Attribute Types



Basic Statistical  
Descriptions of Data



Data Visualization

# Trouble with Summary Stats

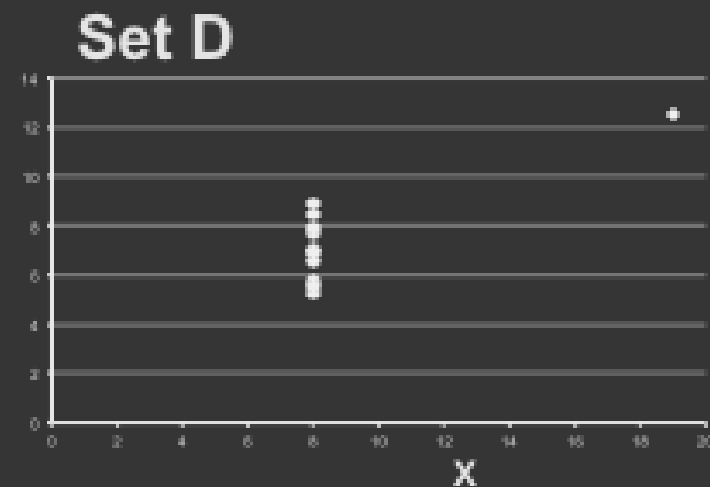
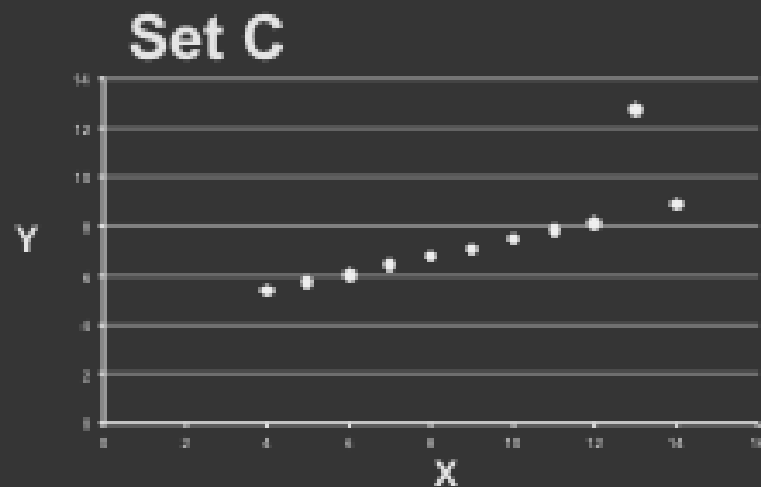
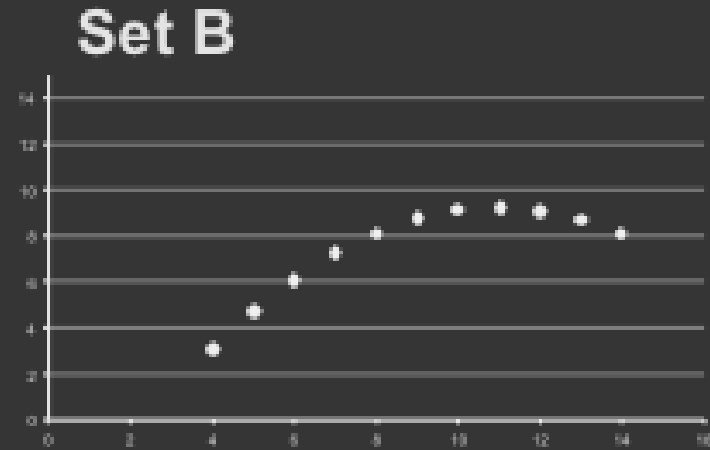
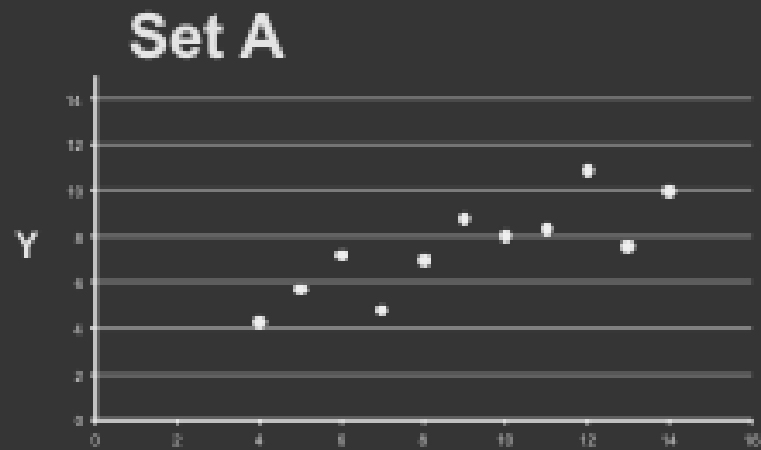
Set A		Set B		Set C		Set D	
X	Y	X	Y	X	Y	X	Y
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.1	14	8.84	8	7.04
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.1	4	5.39	19	12.5
12	10.84	12	9.11	12	8.15	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

$$\begin{aligned}u_x &= 9.0 & \sigma_x &= 3.317 & Y &= 3 + 0.5 X \\u_y &= 7.5 & \sigma_y &= 2.03 & R^2 &= 0.67\end{aligned}$$

[Anscombe 73]

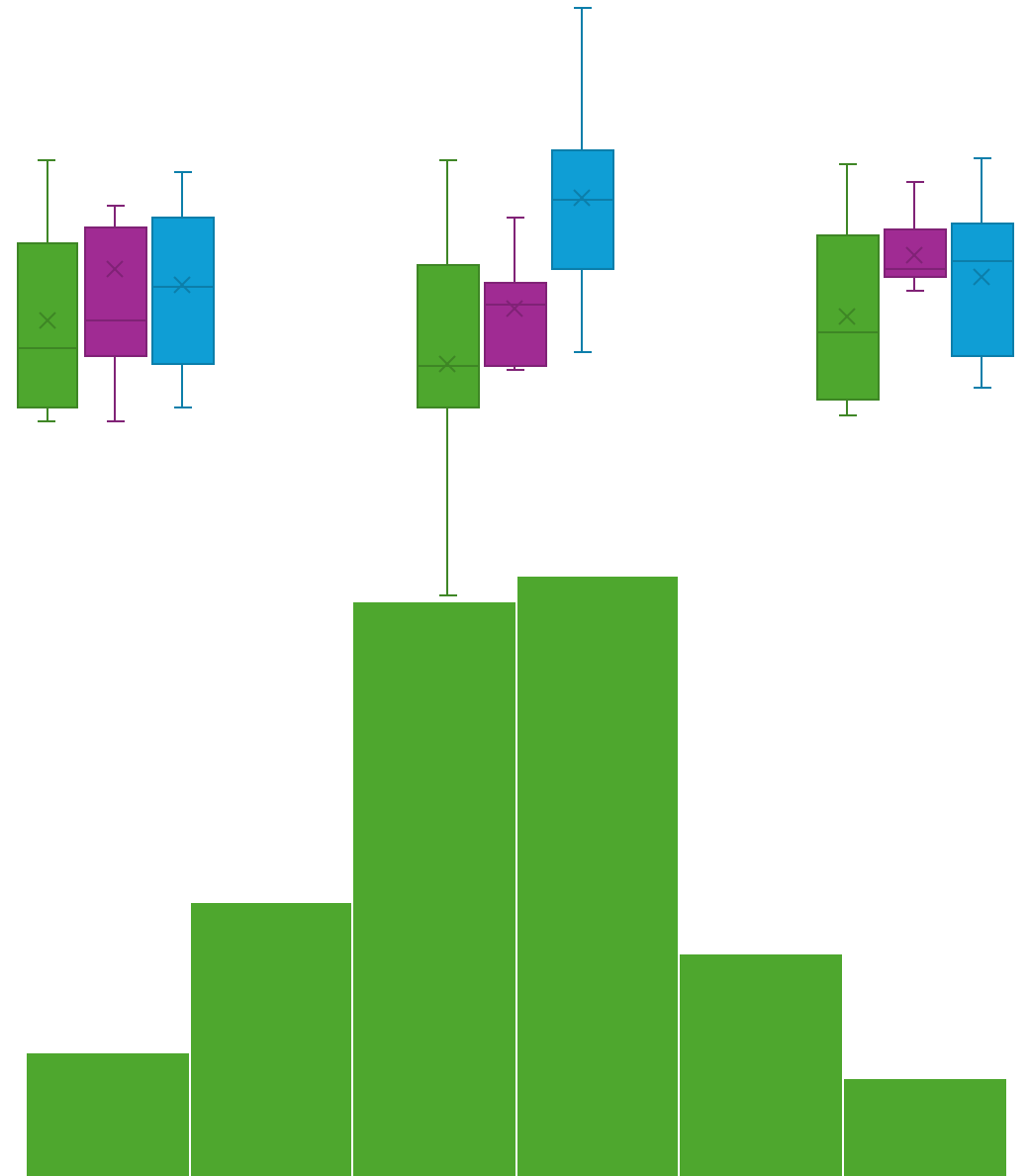
# Looking at Data



# Chart Types

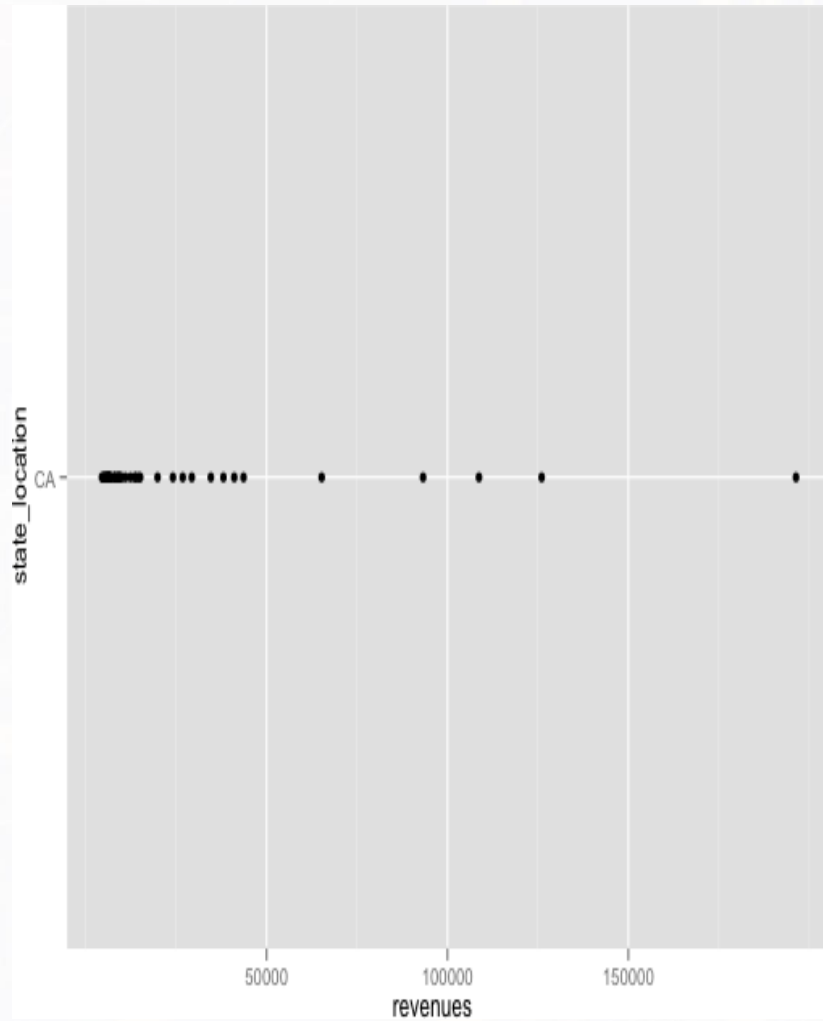
## Single variable

- 📊 Dot plot
- 📊 Jitter plot
- 📊 Histogram and bar chart
- 📊 Kernel density estimate
- 📊 Cumulative distribution function

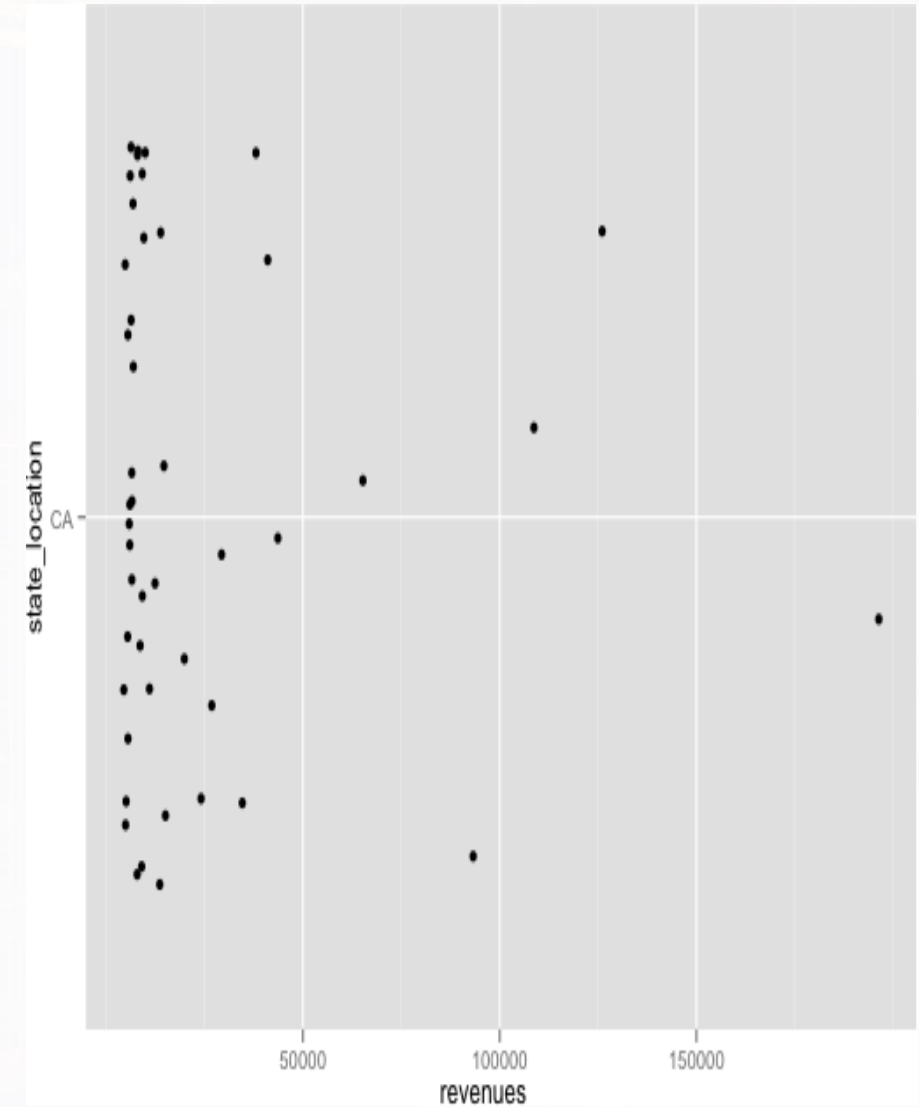


# Dot Plot and Jitter Plot

Dot



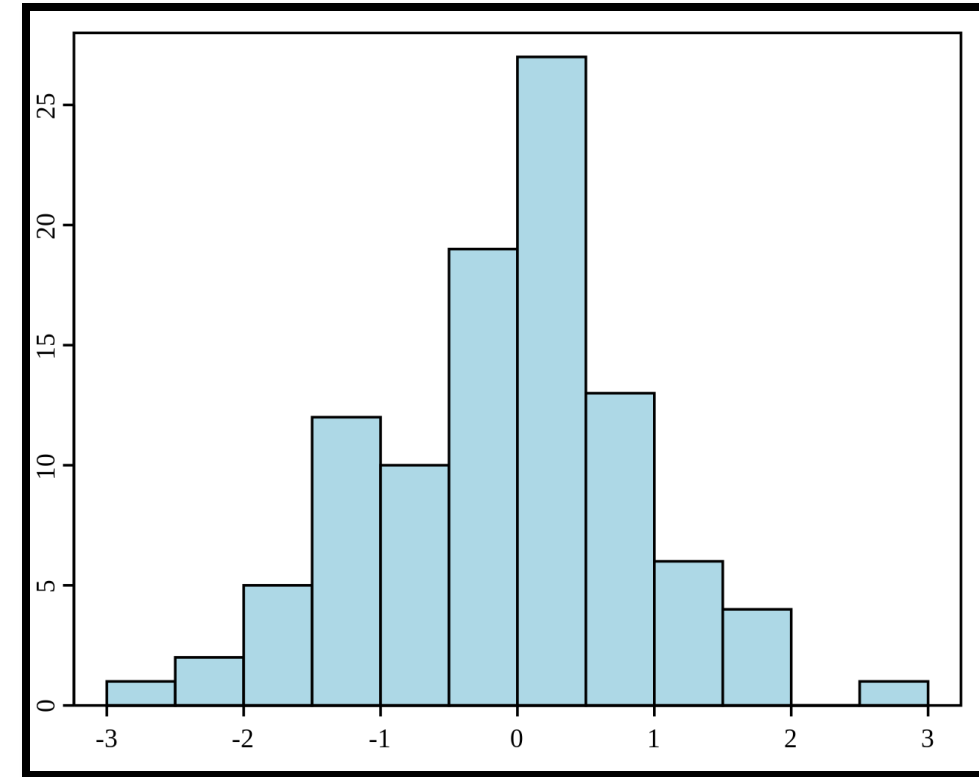
Jitter



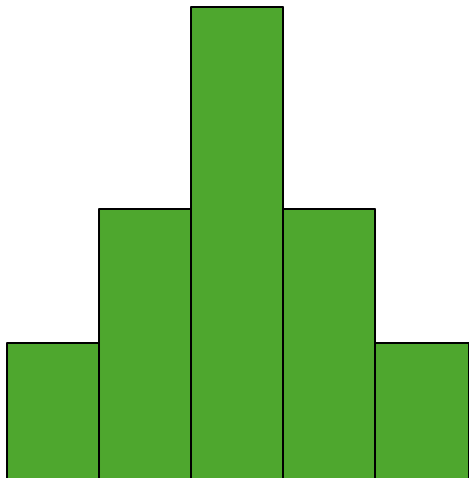
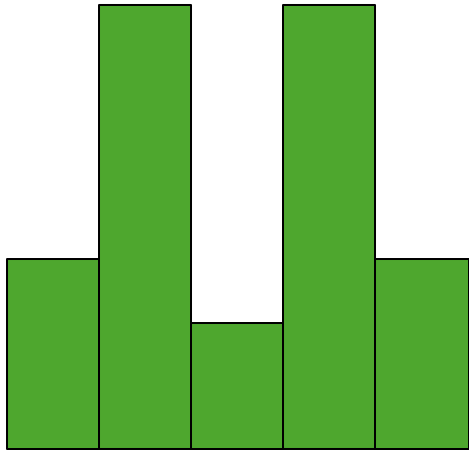


# Histogram Analysis

- ▮ Histogram: Graph display of tabulated frequencies, shown as bars
- ▮ It shows what proportion of cases fall into each of several categories
- ▮ The categories are usually specified as non-overlapping intervals of the variable. The categories (bars) must be adjacent



# Histograms



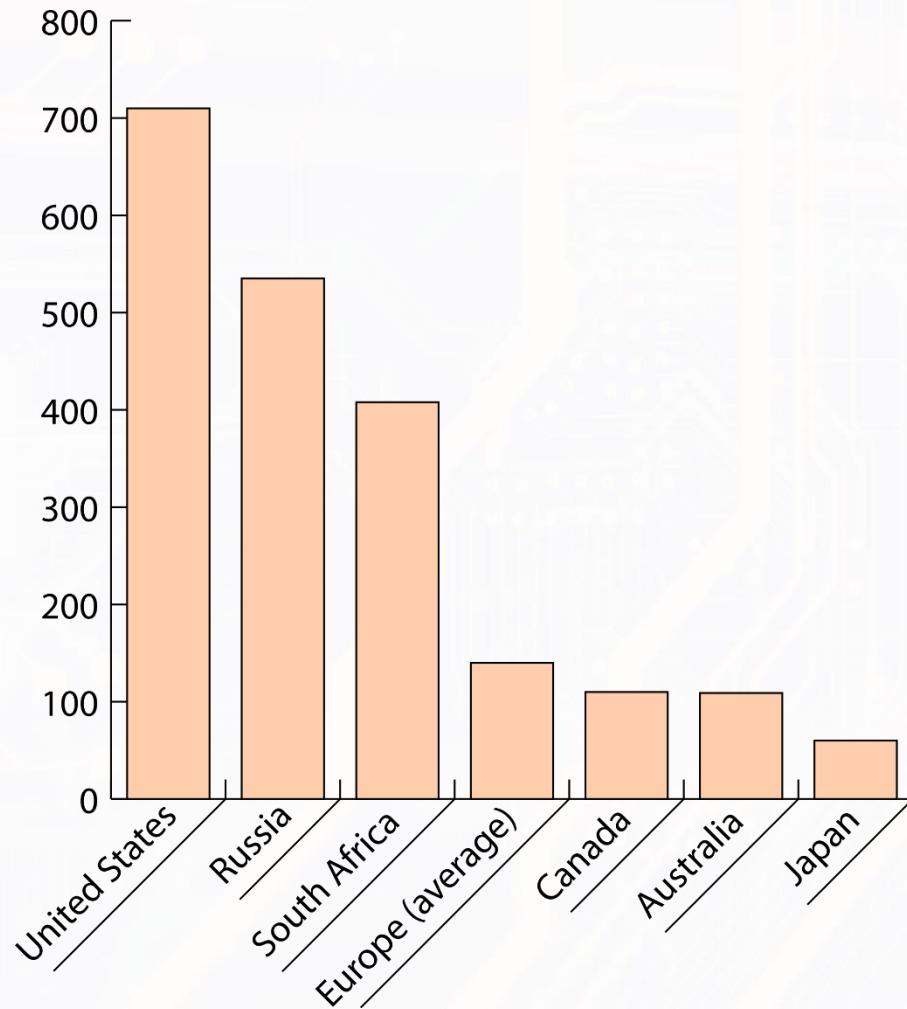
- Complement boxplots
- The two histograms shown in the left may have the same boxplot representation
- The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions

# Bar Plot and Histogram

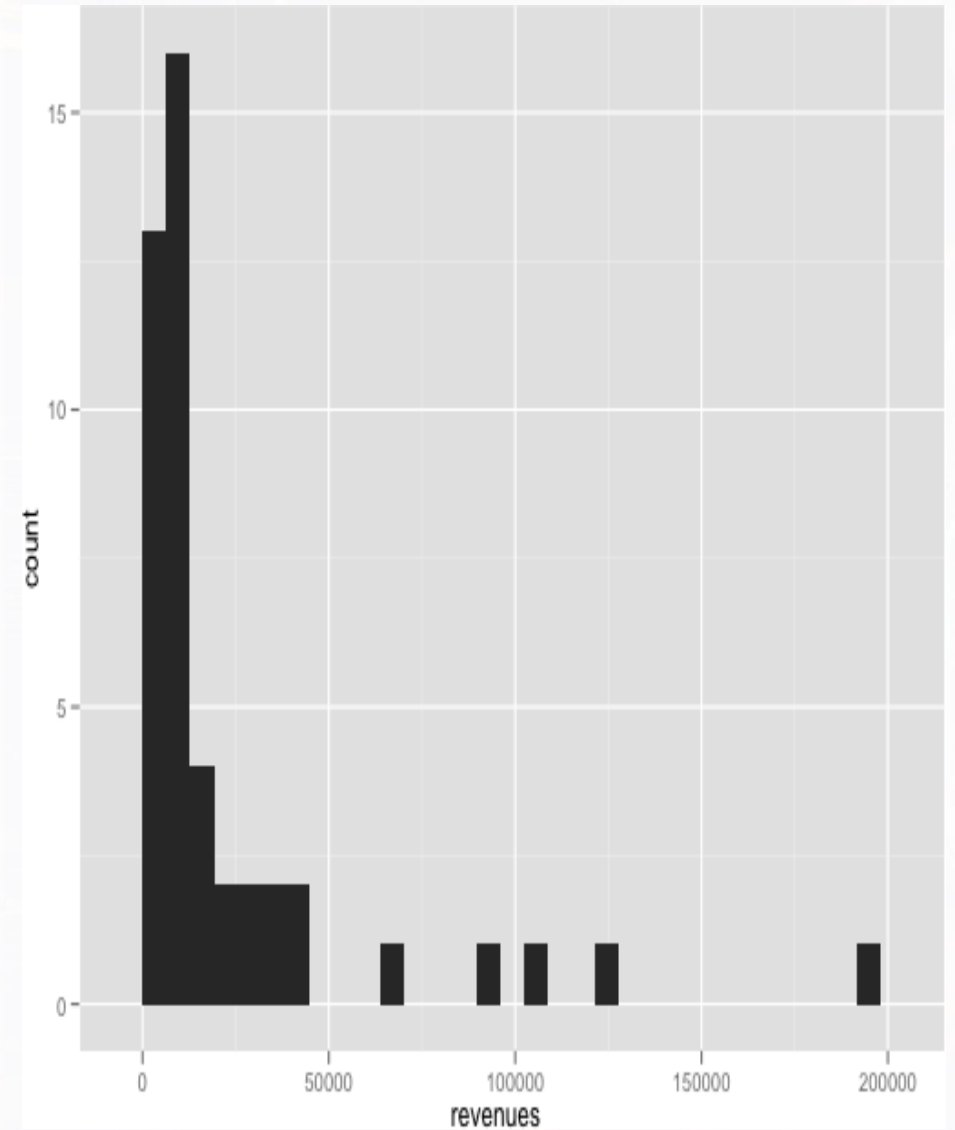


X variable is discrete

Bar

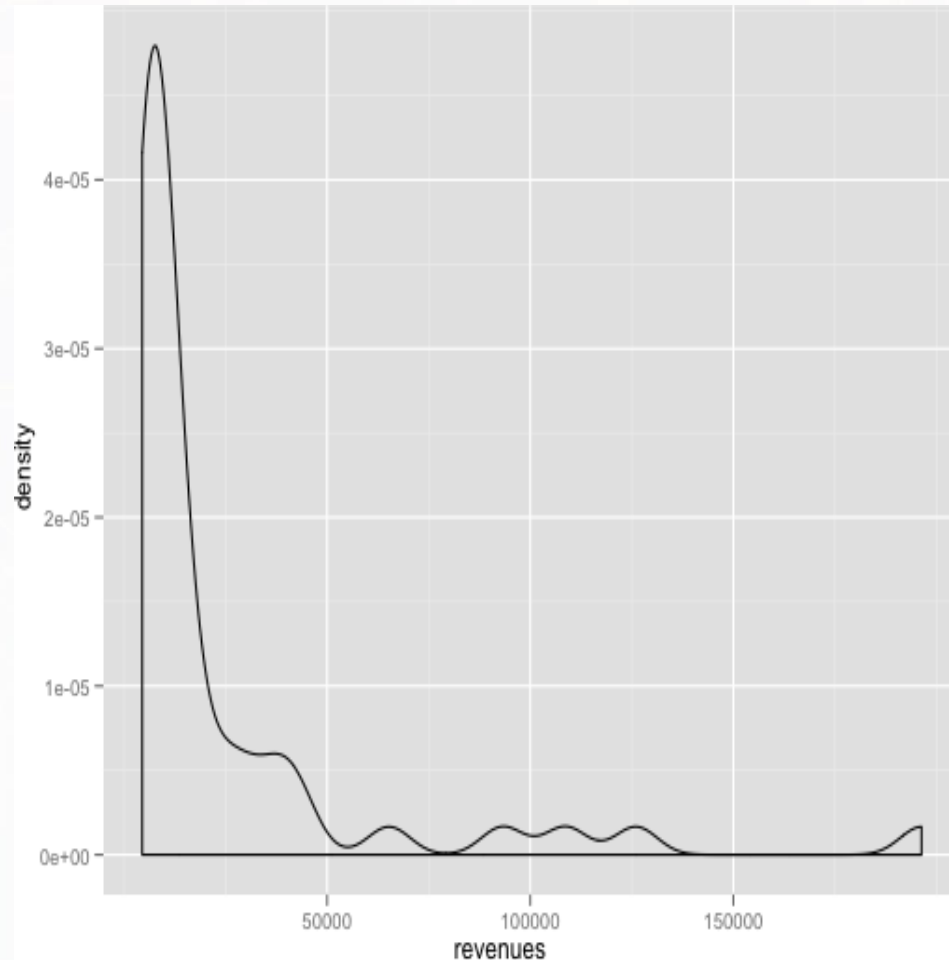


Histogram

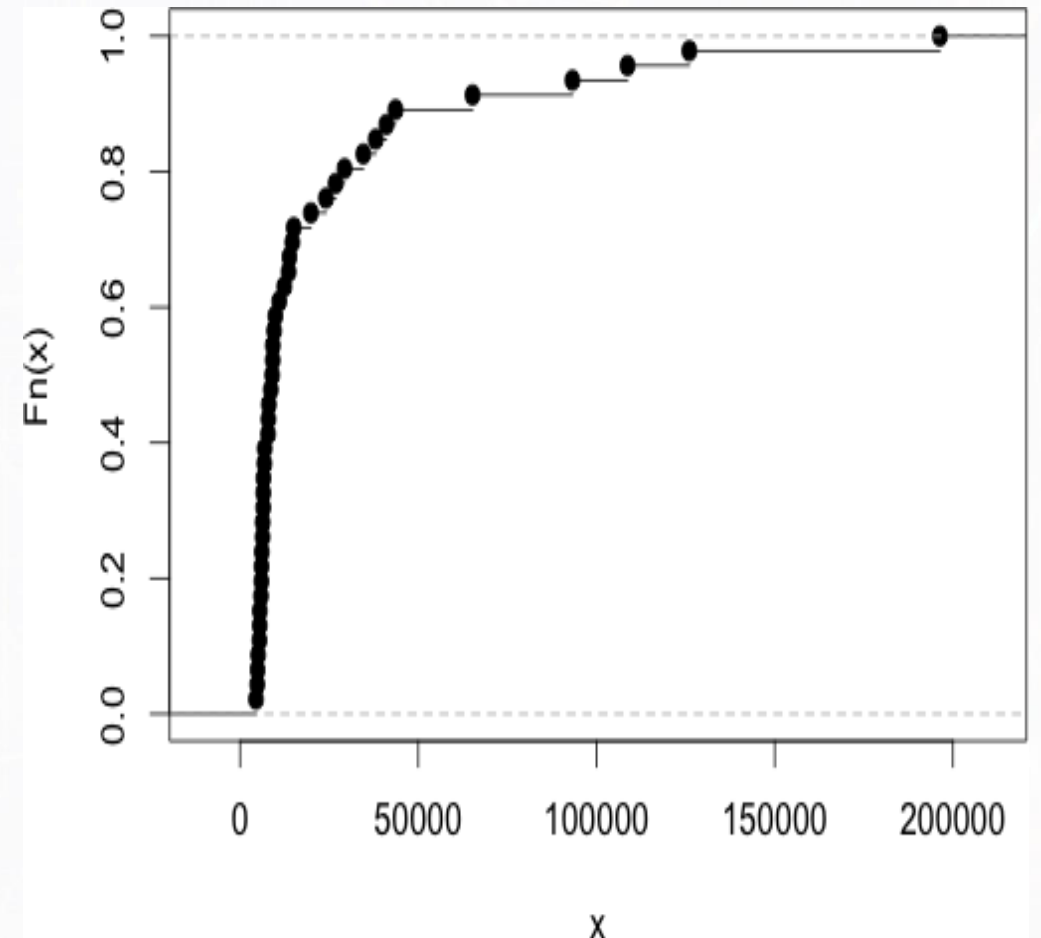


# Other Plots

## Kernel Density Estimation

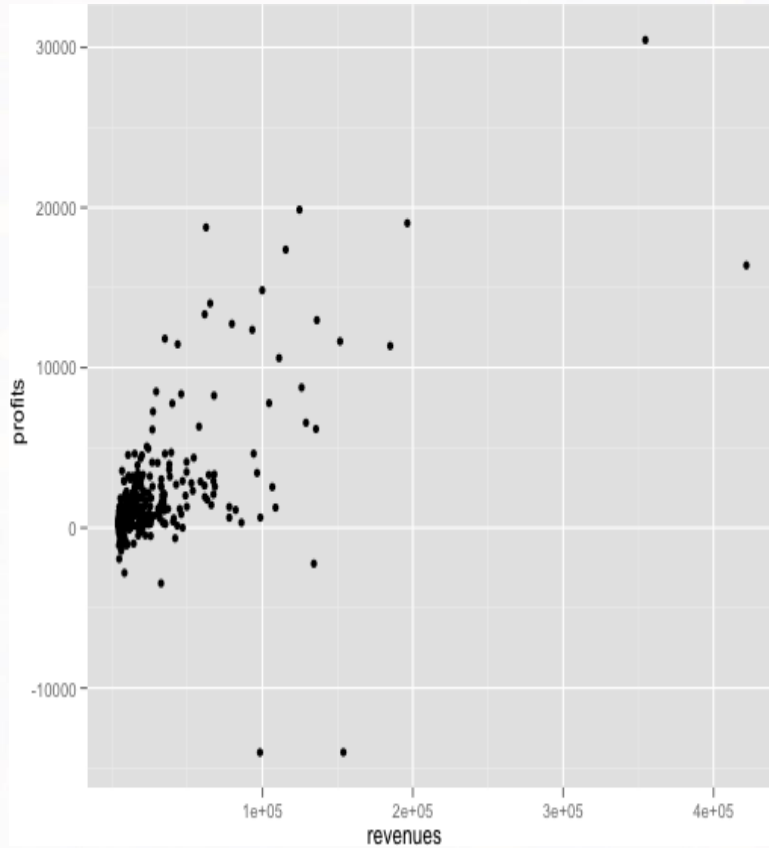


## Cumulative Distribution Function

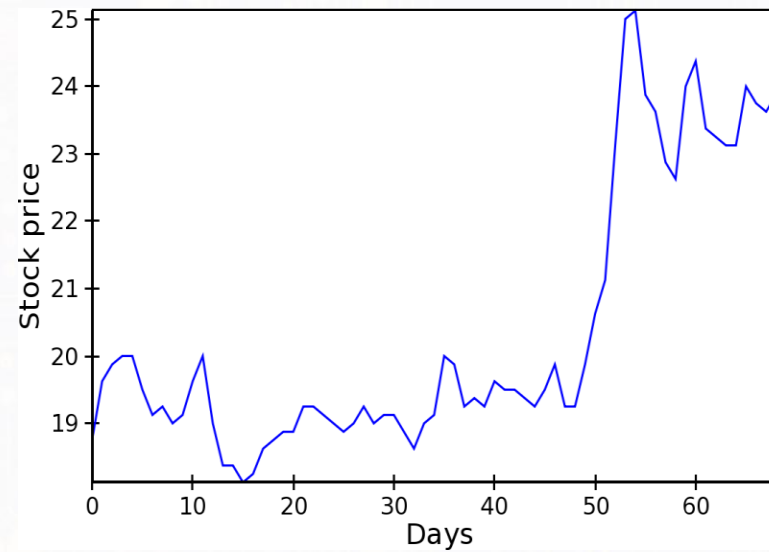


# Two-Variable Chart Types

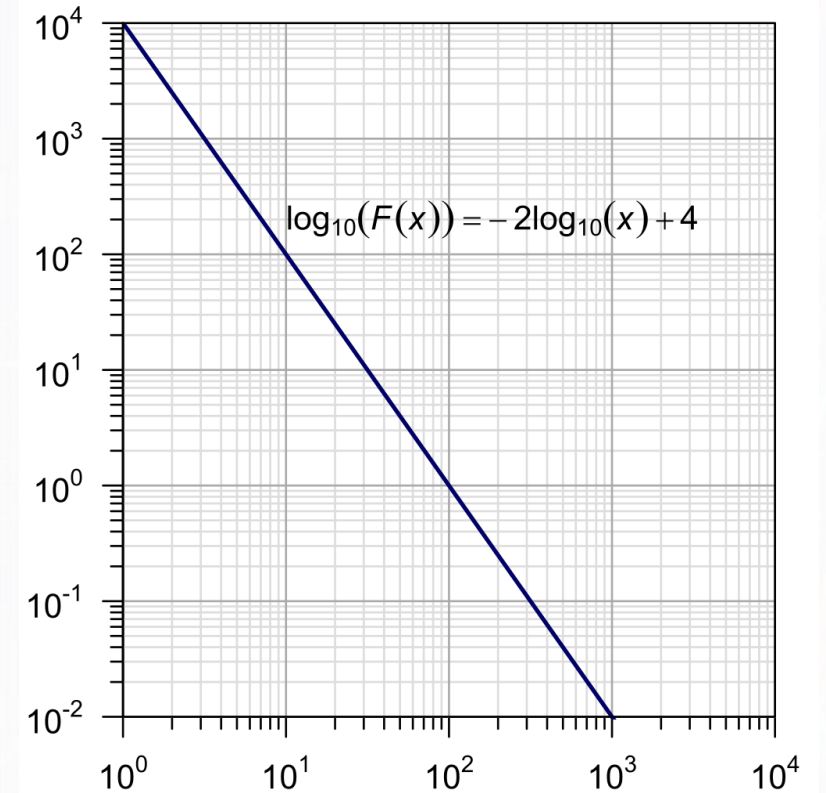
## Scatter Plot



## Line Plot

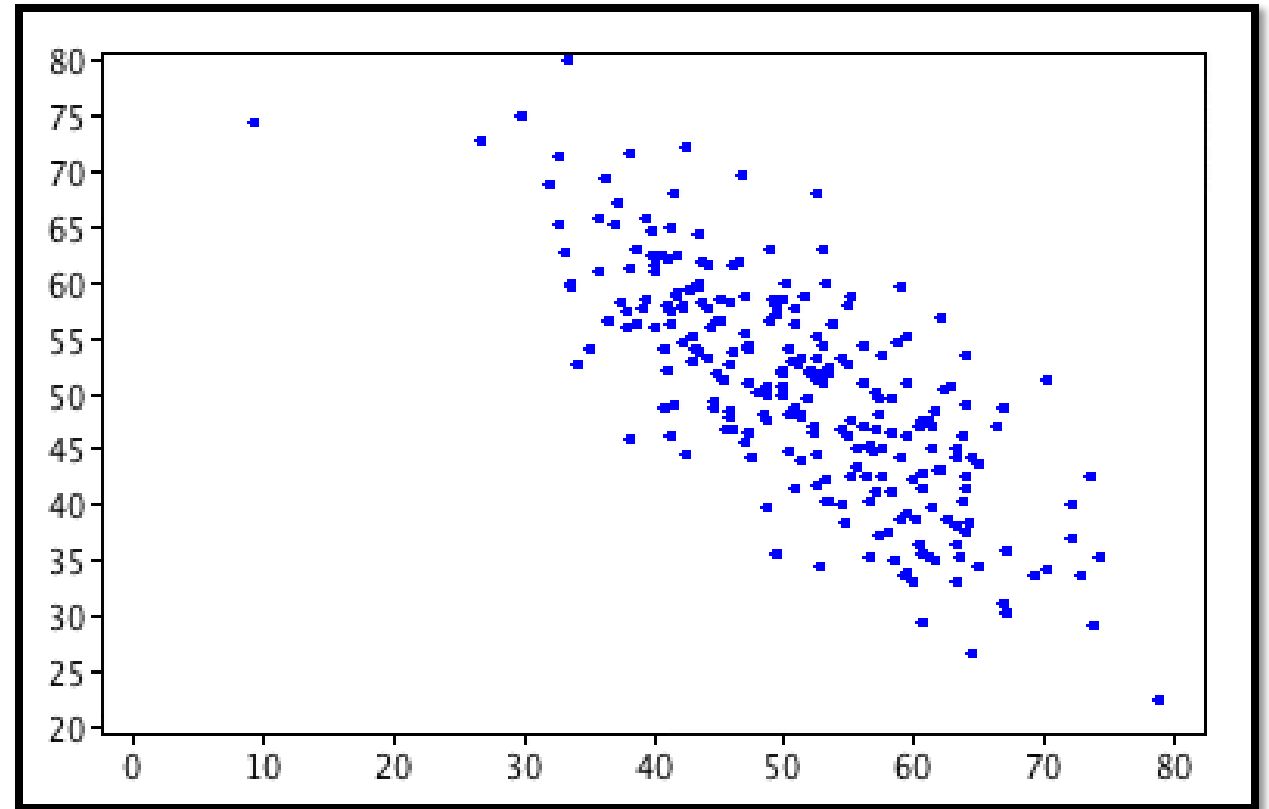


## Log-Log Plot



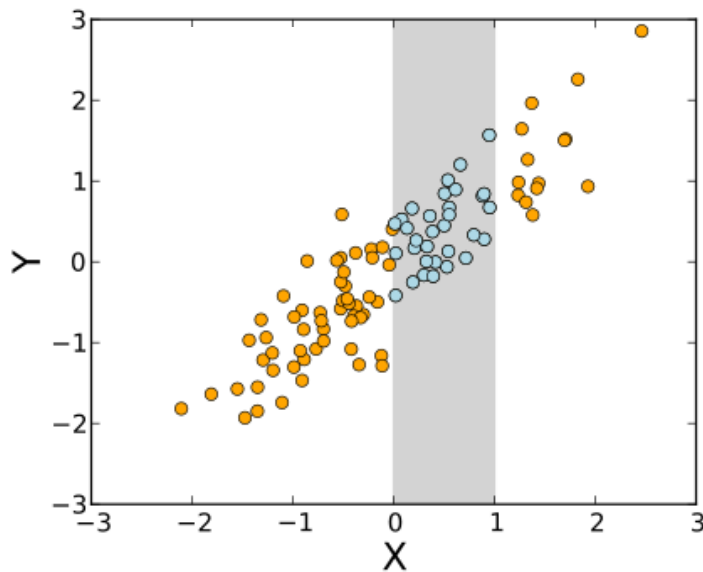
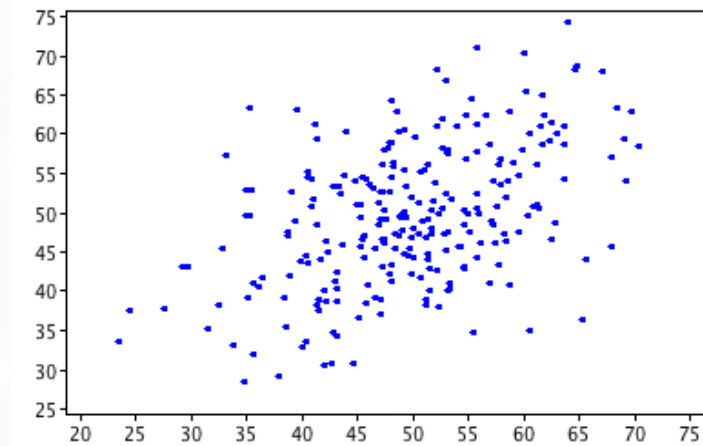
# Scatter Plot

- ↳ Provides a first look at bivariate data to see clusters of points, outliers, etc
- ↳ Each pair of values is treated as a pair of coordinates and
- ↳ Plotted as points in the plane

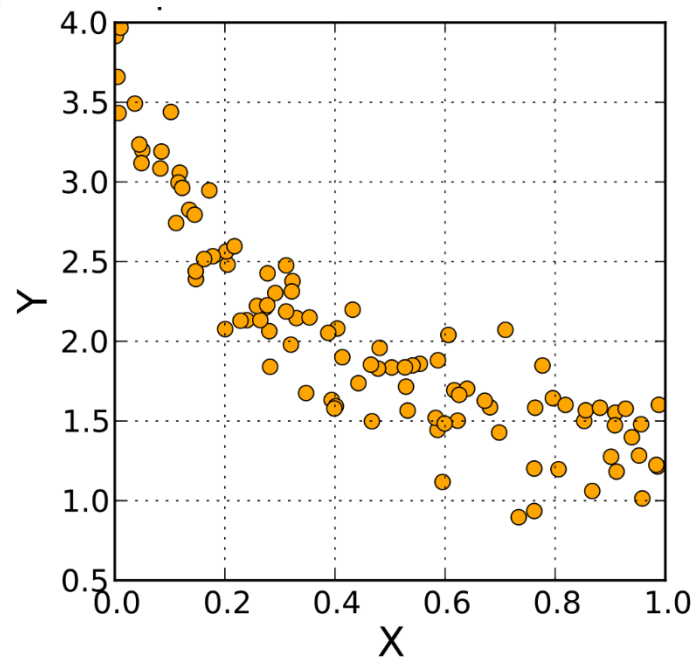
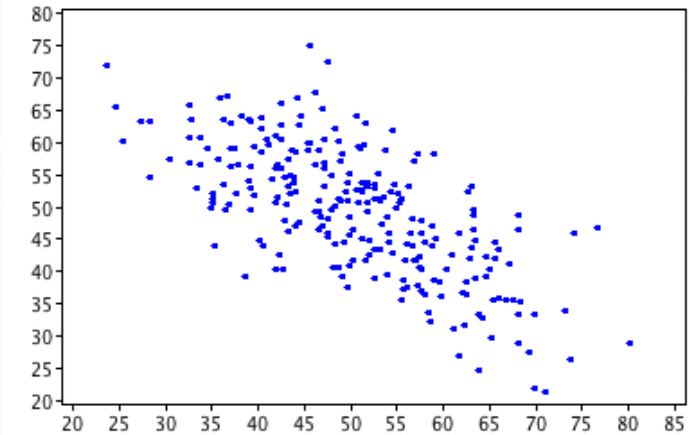


# Data Correlation

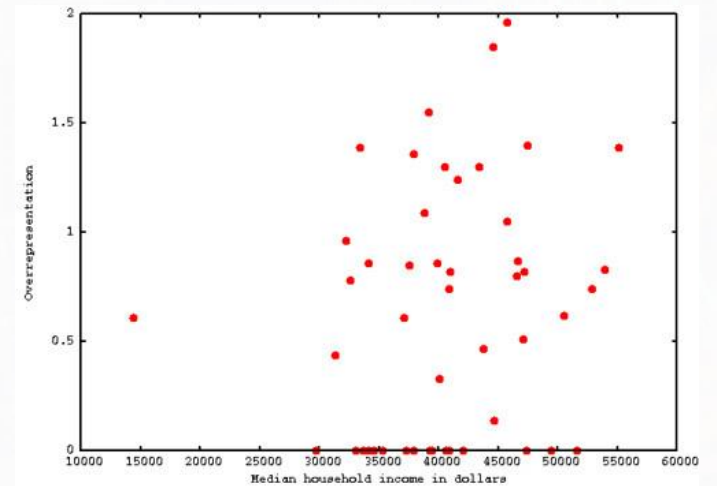
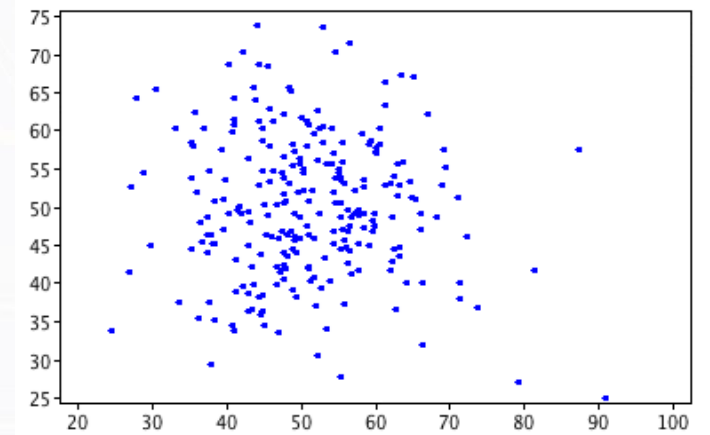
## Positively Correlated



## Negatively Correlated

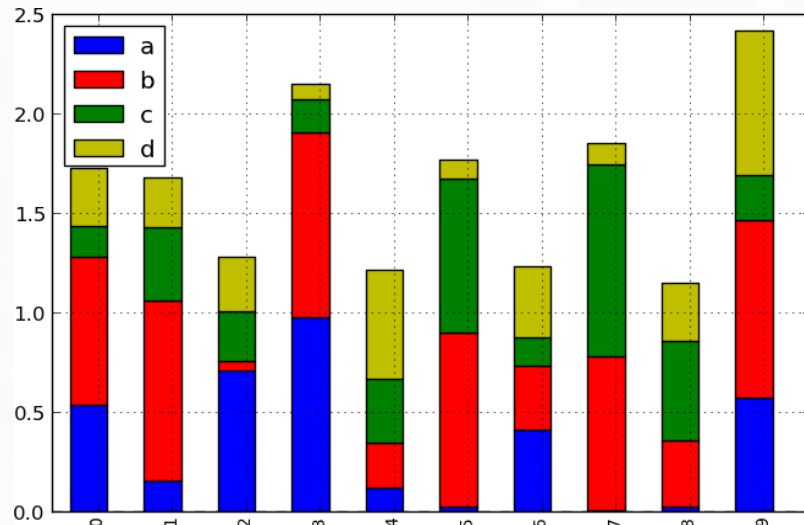


## Uncorrelated

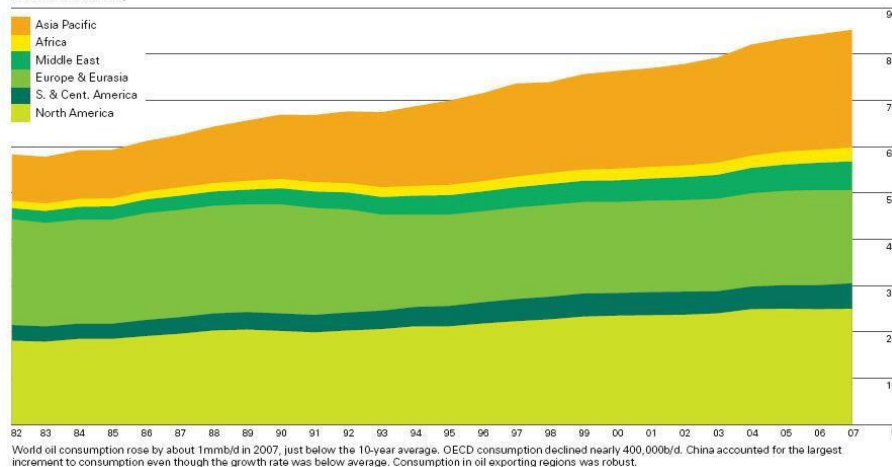


# More Than Two Variables

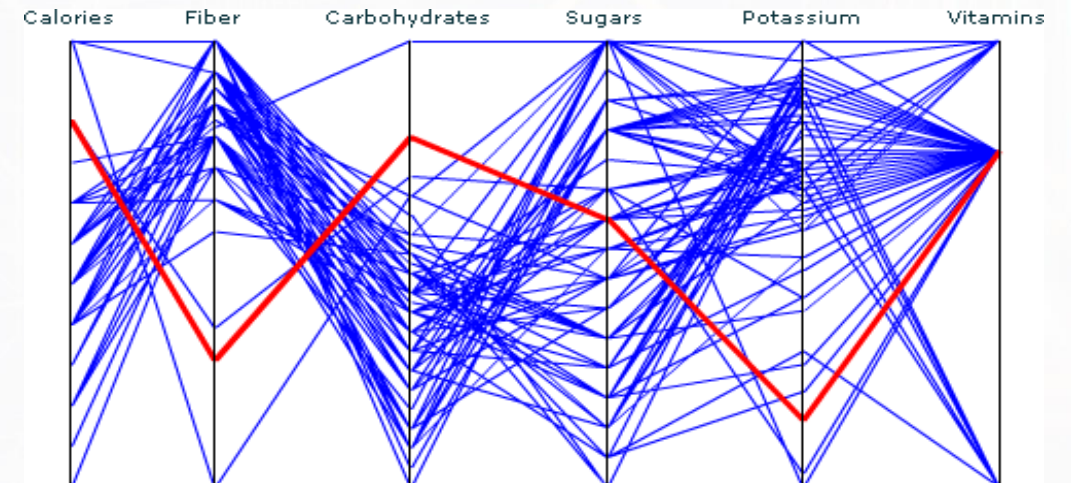
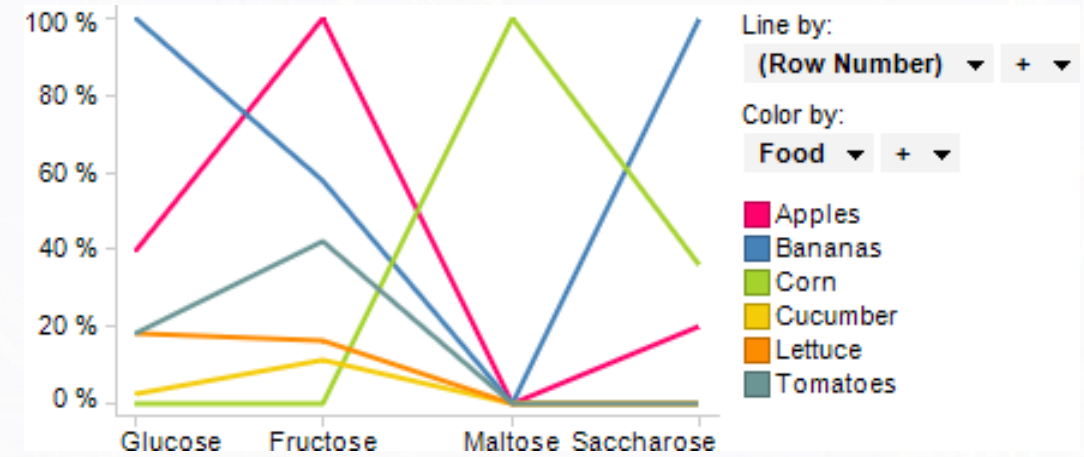
## Stacked Plots



Consumption by region  
Million barrels daily



## Parallel Coordinate Plot





# Data Presentation Can Be Art



# Summary

- Data attribute types: nominal, binary, ordinal, interval-scaled, ratio- scaled
- Many types of data sets, e.g., numerical, text, graph, Web, image.
- Gain insight into the data by:
  - Basic statistical data description: central tendency, dispersion, graphical displays
  - Data visualization: map data onto graphical primitives
  - Measure data similarity
- Above steps are the beginning of data preprocessing.
- Many methods have been developed but still an active area of research.

# Attribution

Material presented in this lecture was adapted from

Han, J., Kamber, M., & Pei, J. (2012). *Data mining concepts and techniques* (3rd ed.). Waltham, MA: Elsevier. Retrieved from [https://hanj.cs.illinois.edu/bk3/bk3\\_slidesindex.htm](https://hanj.cs.illinois.edu/bk3/bk3_slidesindex.htm)

and

Canny, J., Franklin, M., Bruckner, Sparks, E., & Venkataraman, S. (2014). CIS194 introduction to data science [PowerPoint]. Retrieved from <https://bcourses.berkeley.edu/courses/1267848/files/folder/lectures>