# Introduction to Business Analytics

Ron Guymon

ILLINOIS
Gies College of Business

# Introduction to Module 2

Introduction to Business Analytics

# What Is Data?

# What Is Data?

# What Is Data?

# What Is Data?

# What Is Data?

"Factual information (such as measurements or statistics) used as a basis for reasoning, discussion, or calculation"

https://www.merriam-webster.com/dictionary/data

# Module 2 Objectives

- Appreciate some common issues associated with assembling data.

- Investigate how assembling data is related to framing a question and calculating the results.

- Execute some fundamental data assembly tasks.

# Framing a Question

Introduction to Business Analytics

# FACT Framework

- Frame the question

- Assemble the data

- Calculate the results

- Tell others the results

# Framing a Question Influences:

- Whether or not you will spend time supporting the organization's goals

- Whether you even need to use data

- If you do could benefit from data, then it will influence the data that is assembled.

# Frame a Question That

- Considers organization's goals

- Is informed by domain knowledge

- Leads to action

# Bad Examples of Framing a Question

- Why are things bad?

- Why are sales down?

# Good Examples of Framing a Question

- In what regions are sales lower than expected?

- What factors caused a decline in sales for regions in which sales are lower than expected?

# Great Examples of Framing a Question

- What factors are most influential in identifying regions that have sales that are statistically lower than expected?

- What factors cause a statistically significant decline in sales for regions in which sales are lower than expected?

# Assembling Data

Introduction to Business Analytics

# FACT Framework

- Frame the question

- Assemble the data

- Calculate the results

- Tell others the results

# Steps for Assembling Data

- Finding data

- Extract, transform, load (ETL)

- Data wrangling

Chief Data Officers "have been chartered with improving the efficiency and value-generating capacity of their organization's information ecosystem. That is, they've been asked to lead their organization in treating and leveraging information with the same discipline as its other, more traditional assets."

Laney, Doug. *Infonomics*, p. 9

# A Few Data Sources on the Web

www.data.gov

www.google.com/publicdata/directory

aws.amazon.com/opendata/public-datasets/

docs.microsoft.com/en-us/azure/sql-database/sql-database-public-data-sets
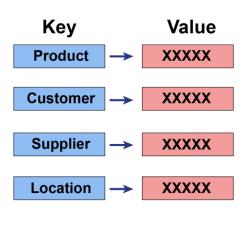
https://www.kaggle.com/datasets
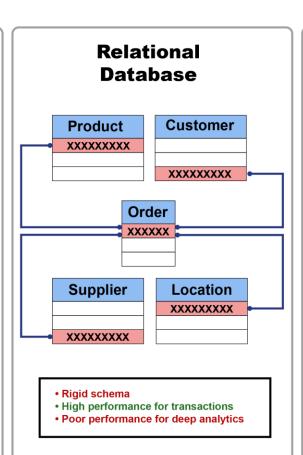
https://data.world/
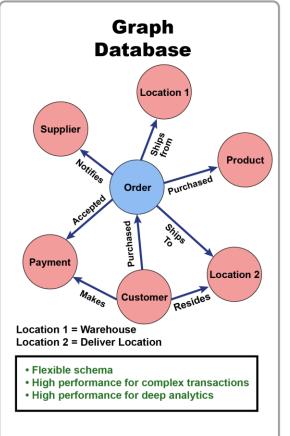
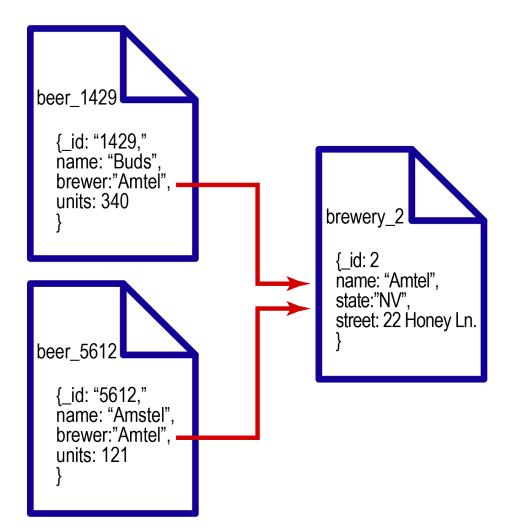https://www.gapminder.org/data/

# ETL

- Extract

- Transform

- Load

# Key-Value Database

**Key**       **Value**

| Product | → | XXXXX |
| Customer | → | XXXXX |
| Supplier | → | XXXXX |
| Location | → | XXXXX |

- Highly fluid schema/no schema
- High performance for simple transactions
- Poor performance deep analytics

# Relational Database

**Product**
XXXXXXXX

**Customer**
XXXXXXXXX

**Order**
XXXXXX

**Supplier**
XXXXXXXXX

**Location**
XXXXXXXXX

- Rigid schema
- High performance for transactions
- Poor performance for deep analytics

# Graph Database

Supplier — Notifies / Accepted — Order
Order — Ships from — Location 1
Order — Purchased — Product
Order — Ships To — Location 2
Order — Purchased — Customer
Payment — Accepted — Order
Customer — Makes — Payment
Customer — Resides — Location 2

Location 1 = Warehouse
Location 2 = Deliver Location

- Flexible schema
- High performance for complex transactions
- High performance for deep analytics

"timestamp":"2017-06-03T18:42:18.018", "deltaStartMillis"
"class":"com.orgmanager.handlers.RequestHandler", "method":
"sizeChars":"5022", "message":"Duration Log", "durationMillis"
"webURL":"/app/page/analyze", "webParams":"null", "durationMillis"
"requestID":"8249868e-afd8-46ac-9745-839146a20f09", "class":"com
"durationMillis":"36"}{"timestamp":"2017-06-03T18:43:335.030", "sessionID":"144o2n620
"webParams":"file=chartdata_new.json", "class":"com.orgmanager.handlers.
"sessionID":"144o2n620jm9trnd3s3n7wg0k", "sizeChars":"48455", "message":"Duration Log
"deltaStartMillis":"0", "level":"INFO", "webURL":"/app/page/report", "webParams":"null
"requestID":"789d89cb-bfa8-4e7d-8047-498454af885d", "sessionID":"144o2n620jm9trnd3s3n
"durationMillis":"7"}{"timestamp":"2017-06-03T18:46:921.000", "deltaStartMillis":"0",
"class":"com.orgmanager.handlers.RequestHandler", "method":"handle", "requestID":"7ac6c
"sizeChars":"10190", "message":"Duration Log", "durationMillis":"10"}{"timestamp":"2017
"webURL":"/app/rest/json/file", "webParams":"file=chartdata_new.json", "class":"com.
"requestID":"7ac6ce95-19e2-4a60-88d7-6ead86e273d1", "sessionID":"144o2n620jm9trnd3s3n
"durationMillis":"23"}{"timestamp":"2017-06-03T18:42:18.018", "deltaStartMillis":"0",
"class":"com.orgmanager.handlers.RequestHandler", "method":"handle", "requestID":"b886
"sizeChars":"5022", "message":"Duration Log", "durationMillis":"508"}{"timestamp":"2017
"webURL":"/app/page/analyze", "webParams":"null", "class":"com.orgmanager.handlers.
"requestID":"8249868e-afd8-46ac-9745-839146a20f09", "class":"com.orgmanager.handlers.
"durationMillis":"36"}{"timestamp":"2017-06-03T18:43:335.030", "sessionID":"144o2n620
"webParams":"file=chartdata_new.json", "class":"com.orgmanager.handlers.
"sessionID":"144o2n620jm9trnd3s3n7wg0k", "sizeChars":"48455", "webURL":"/app/page/report
"deltaStartMillis":"0", "level":"INFO", "webURL":"/app/page/report", "sessionID":"144o2n620
"requestID":"789d89cb-bfa8-4e7d-8047-498454af885d", "timestamp":"2017-06-03T18:46:921.000
"durationMillis":"7"}{"timestamp":"2017-06-03T18:46:921.000", "method":"handle", "method":
"orgmanager.handlers.RequestHandler", "durationMillis":

# Data Wrangling

- Cleaning data

- Combining data with other data

- Cleaning data again

- Combining data with other data again

- Cleaning again

- Then changing its shape

# Calculate the Results

Introduction to Business Analytics

# FACT Framework

- Frame the question

- Assemble the data

- Calculate the results

- Tell others the results

# Summary Statistics = Descriptive Statistics

Examples of Summary Statistics:

- Minimum

- Mean

- Median

- Mode

- Maximum

- Variance

- Standard Deviation

- Range

# Sprinkler Sales

| | Abnormal Sales | Abnormal Temp. | Abnormal Precip. |
|---|---|---|---|
| Northwest | 15 | 5 | -7.5 |
| Midwest | 5 | 2 | -2.5 |
| Northeast | -50 | -17 | 25 |
| East Coast | -60 | -20 | 30 |
| Southeast | -80 | -27 | 40 |
| South | -40 | -13 | 20 |
| Southwest | 10 | 3 | -5 |
| Central Plains | 0 | 0 | 0 |
| Rocky Mountains | 20 | 7 | -10 |
| West Coast | 50 | 17 | -25 |

# Sprinkler Sales Model

Weekly Sales = 250 – 500*Precip + 250*Temp – 46*Wind

# Categories of Machine Learning Algorithms

# Considerations for Selecting an Algorithm

- Tradeoff between accuracy and the time to create the model

- How well does the algorithm deal with:

  - Nonlinearity

  - High dimensionality

  - Multicollinearity

  - Constant stream of new data

- Ability to explain the model to others

# How Is Model Accuracy Evaluated

- Using a "loss" function that measures the difference between the right answers and the predicted answer

- Creating a model on a "training" dataset and then evaluating the performance on a separate "testing" dataset

# Data Types

Introduction to Business Analytics

# Data Types in This Lesson

- Character strings

- Numeric

- Datetimes

- Factors

~/ ⤳

```
> as.numeric('$5,678.34')
[1] NA
Warning message:
NAs introduced by coercion
> |
```

~/ ⤳

```
> as.numeric('5678.34')
[1] 5678.34
>
```

# Dates and Times

- Stored in a special integer format that displays them as dates

- Number of days that have passed since the start of a specific epoch

- March 3, 2005 is stored as the integer 38,415 in Excel.

- Timestamps represent the number of seconds that have passed since the start of an epoch.

# What does 03/04/05 represent?

March 4, 2005

April 3, 2005

April 5, 2003

# Additional Considerations for Dates and Times

- Time zones

- Daylight savings time

- Leap year

- Number of work days in a quarter

- Weekends

- Holidays

~/

```
> read.csv('ities.csv', stringsAsFactors = F)
```

# Conclusion to Module 2

Introduction to Business Analytics