

# Course Introduction: GSND 5345Q

## Fundamentals of Data Science (FDS)

W. Evan Johnson, Ph.D.  
Professor, Division of Infectious Disease  
Director, Center for Data Science  
Rutgers University – New Jersey Medical School  
[w.evan.johnson@rutgers.edu](mailto:w.evan.johnson@rutgers.edu)

2026-01-05

## Section 1

### Introductions

# A Post-COVID Perspective









# Johnson Lab Research

Here is a link to the Johnson Lab Research Page

## Section 2

### Course Introduction

---

# Things you should know about this course

- Lots of diverse material
  - Not a spectator sport!
- Zoom Meeting ID for all sessions is 95146491967, passcode: 236441:
  - Click here for the direct Zoom link
  - Lectures will be recorded and posted in the “Announcements” (Canvas)
- Canvas:
  - The Canvas page will be limited confidential items: course announcements, communication, homework submissions, grades, etc.
- GitHub: [https://github.com/wevanjohnson/2026\\_Spring\\_FDS](https://github.com/wevanjohnson/2026_Spring_FDS)
  - Course information, schedule, lecture notes, homework, etc.
  - Link to Syllabus
- You need to have a basic understanding of statistics
- Learning to program in R is a requirement of this course.

# Resources for learning statistics

Here are some resources to learn basic statistics (and in some cases R simultaneously):

- Data Analysis with R Specialization (Coursera/Duke University)
- Introduction to statistics (Coursera/Stanford)

# Resources for learning R

For learning R:

- RStudio Education
- R Programming (Coursera/Johns Hopkins)
- Data Science R Basics (edx/Harvard University)
- R Training Course (LinkedIn)
- R Programming A - Z: R for Data Science (Udemy)
- Programming with R (Pluralsight)

## Section 3

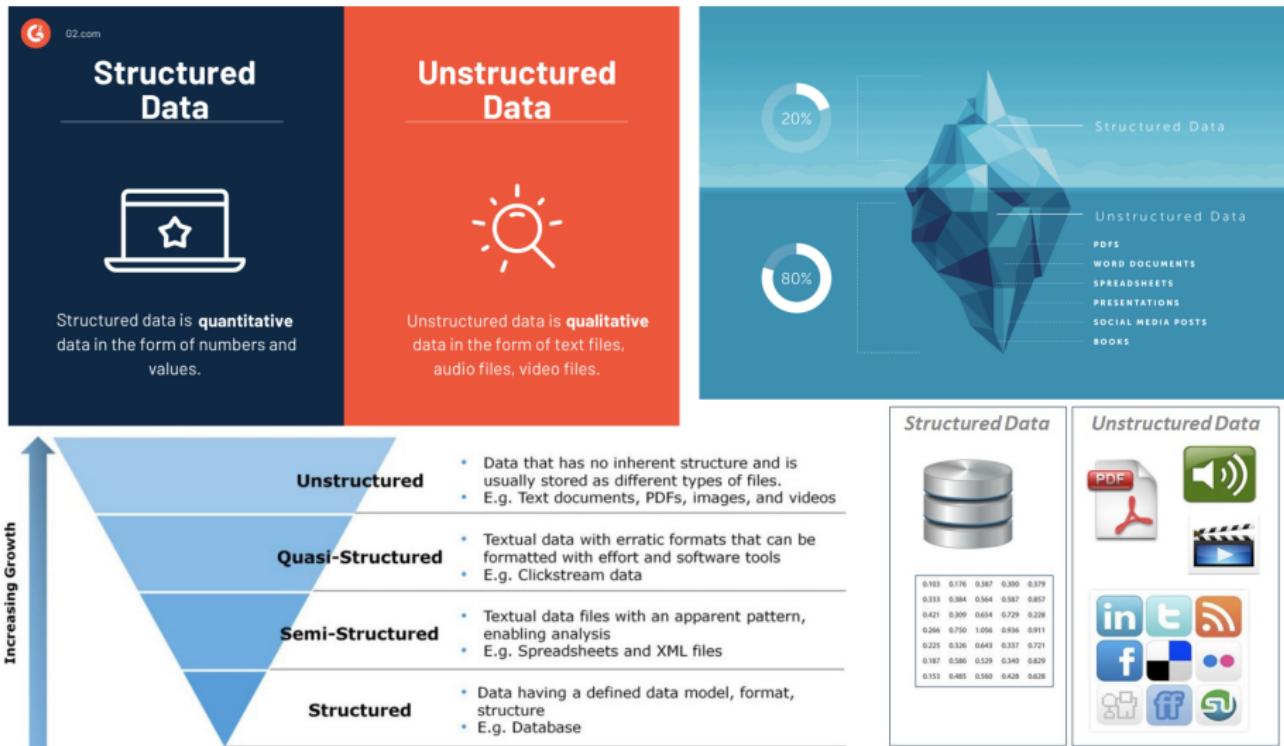
# Introduction to Data Science

# BIG DATA

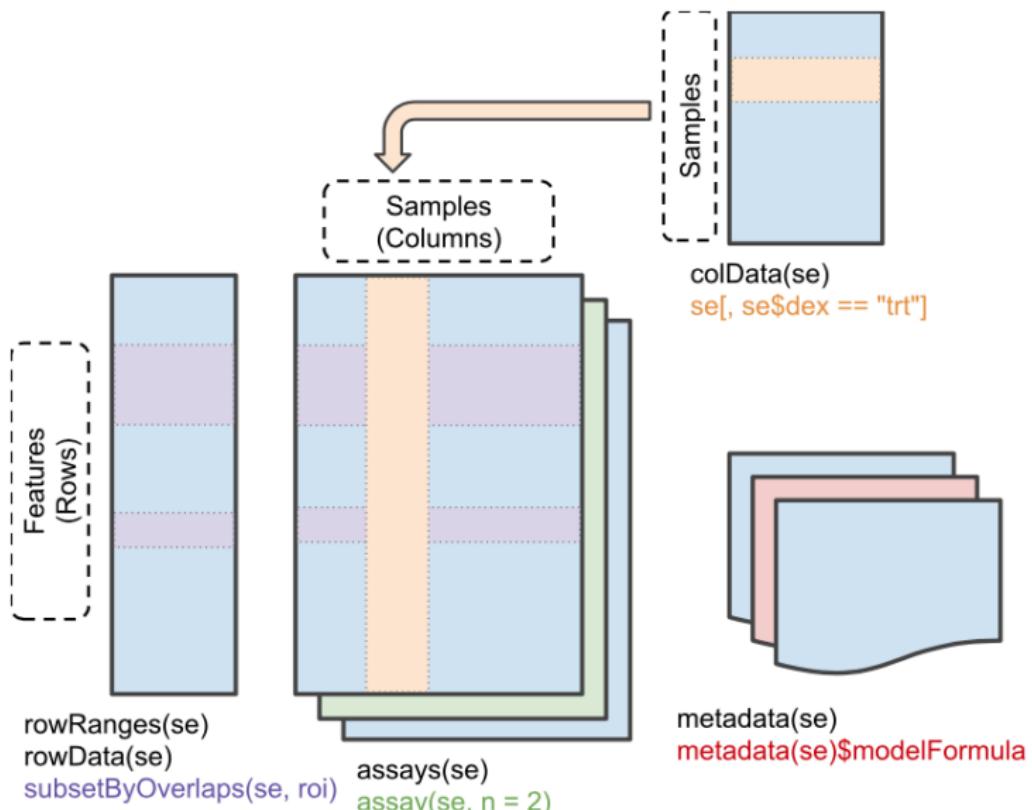


Big Data has fundamentally changed how we look at science and business. Along with advances in analytic methods, they are providing unparalleled insights into our physical world and society.

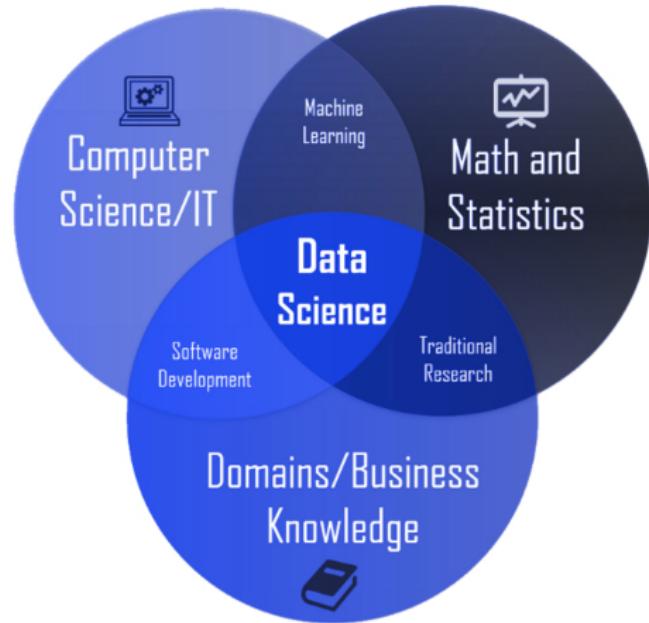
# Structured vs. Unstructured data



# Structured vs. Unstructured data



# Data Science Revolution



- Few have all the skills
- Flexibility in area (business, strategy, health care) and conditions
- Data science makes companies and data better!

# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

## MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants

## DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative



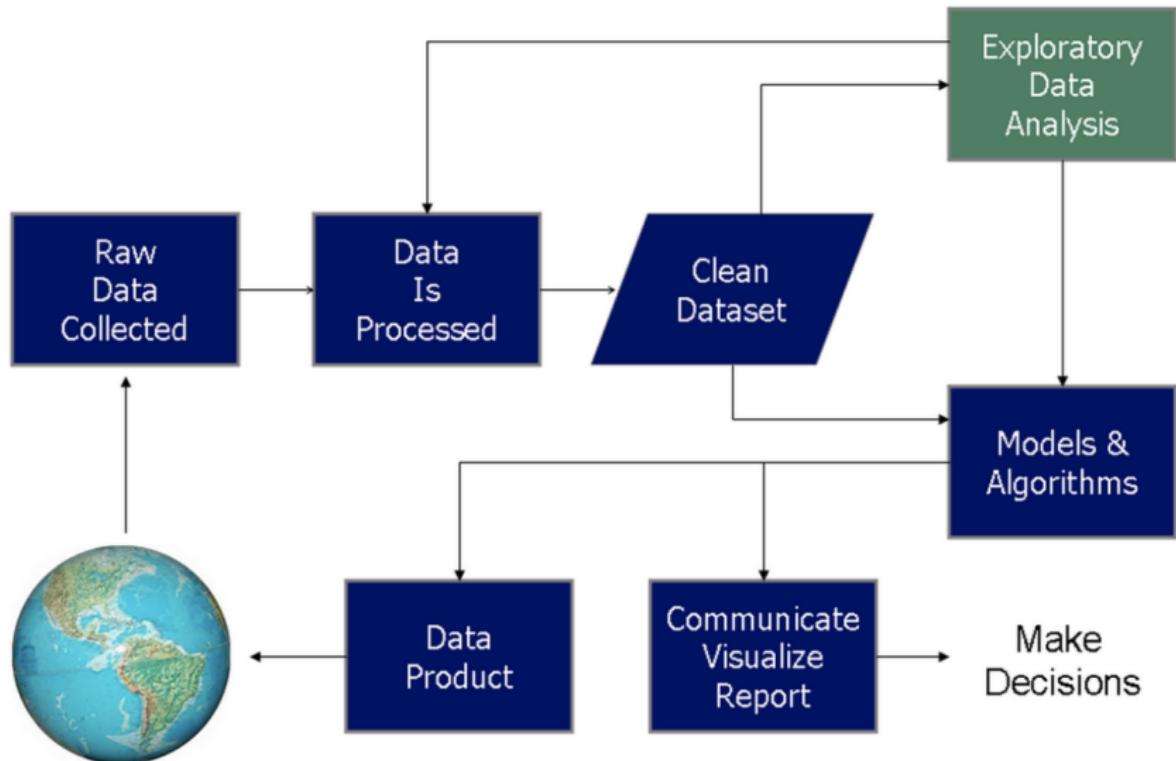
## PROGRAMMING & DATABASE

- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS

## COMMUNICATION & VISUALIZATION

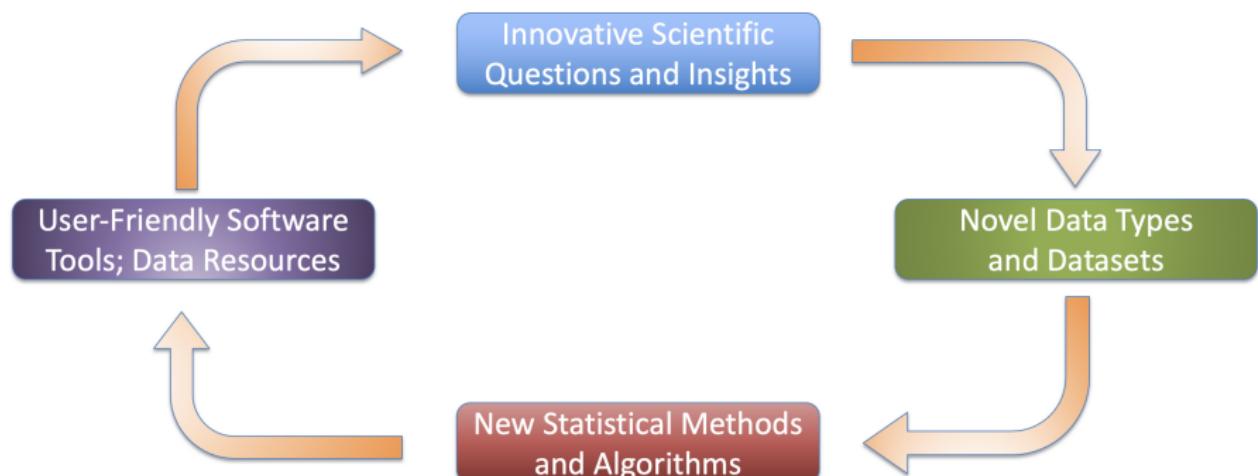
- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

# Data Science Process



# Scientific Cycle for Data Science

Johnson Lab Approach to Science:



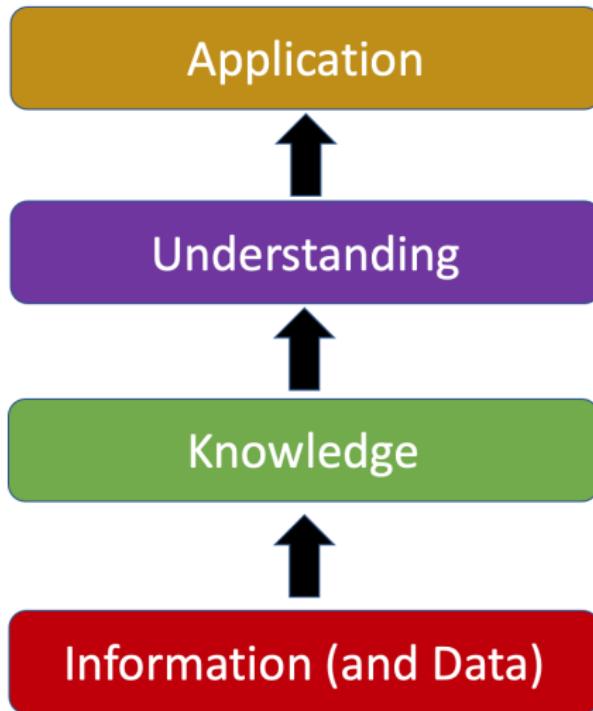
## Section 4

Keeping the “Science” in Data Science

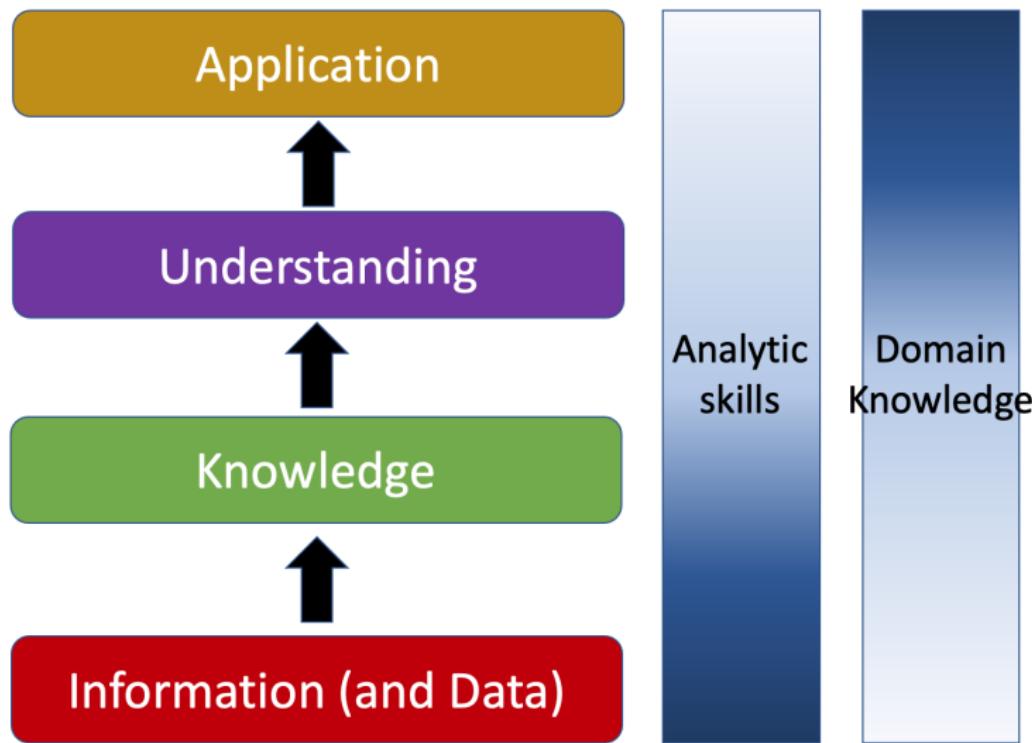
# Domain Knowledge

**Domain knowledge** is knowledge of a specific, specialized discipline or field, in contrast to general (or domain-independent) knowledge. For example, in describing a software engineer may have general knowledge of computer programming as well as domain knowledge about developing programs for a particular industry. People with domain knowledge are often regarded as specialists or experts in their field. (Wikipedia!)

# Analytics Hierarchy



# Analytics Hierarchy



# A Post-COVID Perspective







# Johnson Lab Research

Here is a link to the Johnson Lab Research Page

# Center for Data Science Updates: Courses

- ① GSND 5345Q: Fundamentals of Data Science (Jan 2025)
  - Command-line coding, literate programming, software development, version control, data wrangling and management, and visualization.
- ② GSND 5340Q: High Throughput Biomedical Data Analysis (April 2025)
  - Sequence alignment/QC, GWAS, gene expression and proteomics, epigenetics, metagenomics, and imaging data analysis.
- ③ Machine Learning for Biomedical Data (October 2025)
  - Model training and validation, regression and regularization, unsupervised learning and clustering, dimension reduction and smoothing, supervised learning and classification, neural networks, and Bayesian learning

## Section 5

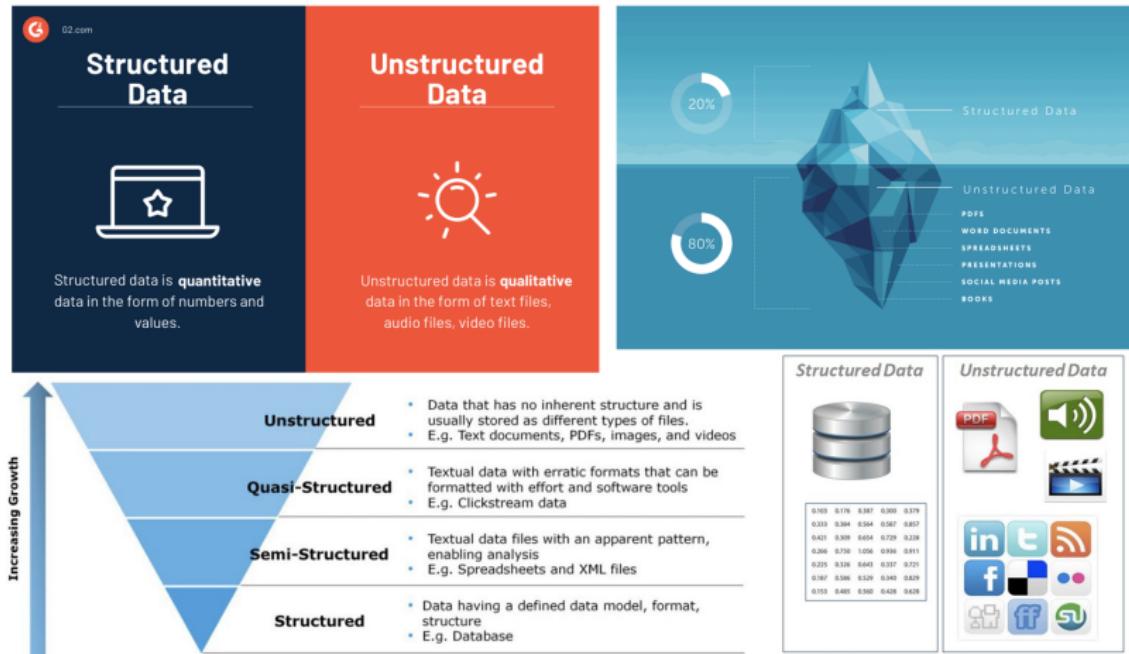
# Introduction to Data Science

# BIG DATA

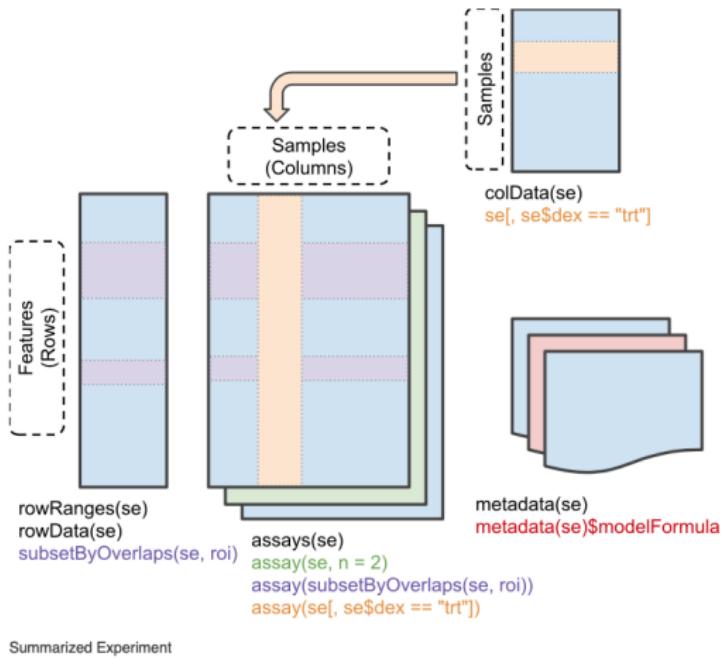


Big Data has fundamentally changed how we look at science and business. Along with advances in analytic methods, they are providing unparalleled insights into our physical world and society

# Structured vs. Unstructured data

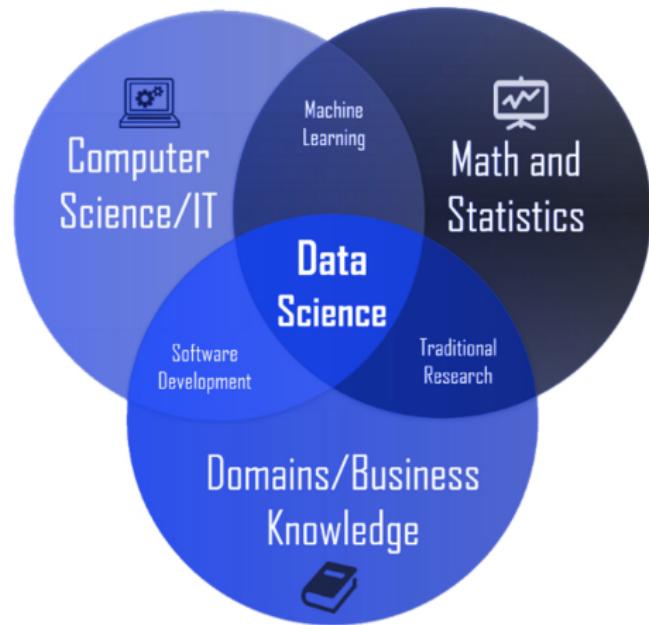


# Structured vs. Unstructured data



Summarized Experiment

# Data Science Revolution



- Few have all the skills
- Flexibility in area (business, strategy, health care) and conditions
- Data science makes companies and data better!

# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

## MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants

## DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative



## PROGRAMMING & DATABASE

- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS

## COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

# Data Science Process

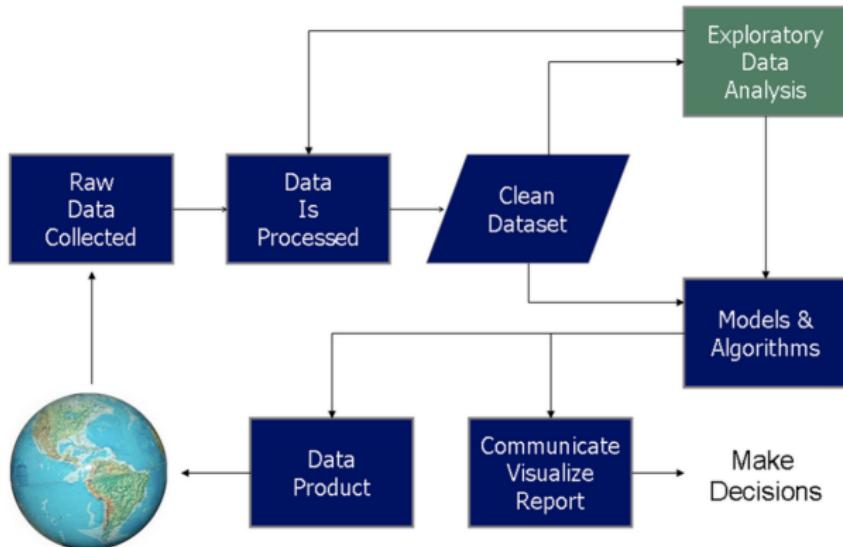
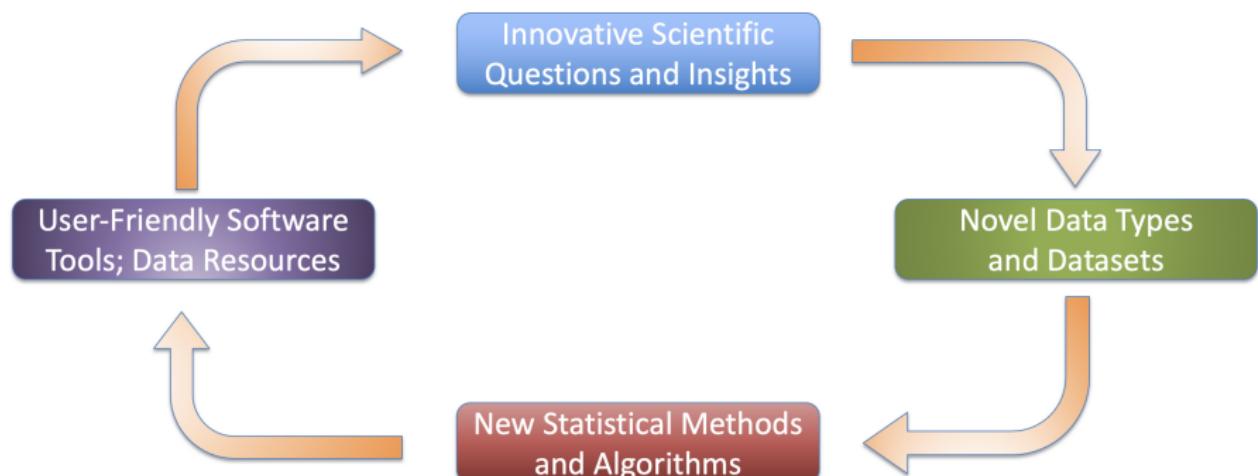


Image: [https://en.wikipedia.org/wiki/Data\\_science](https://en.wikipedia.org/wiki/Data_science)

# Scientific Cycle for Data Science

Johnson Lab Approach to Science:



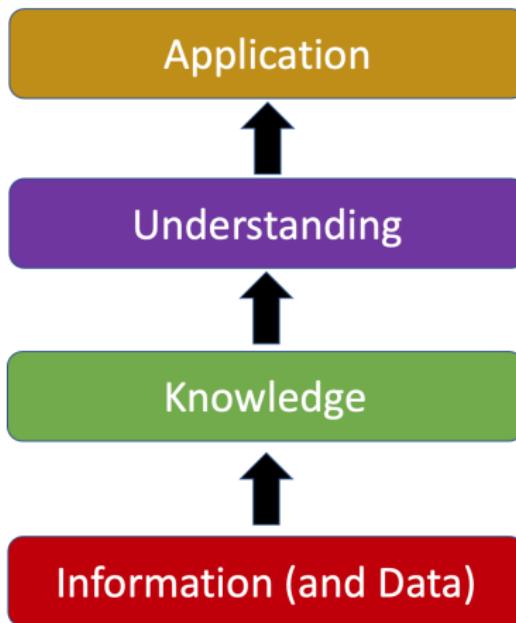
## Section 6

Keeping the “Science” in Data Science

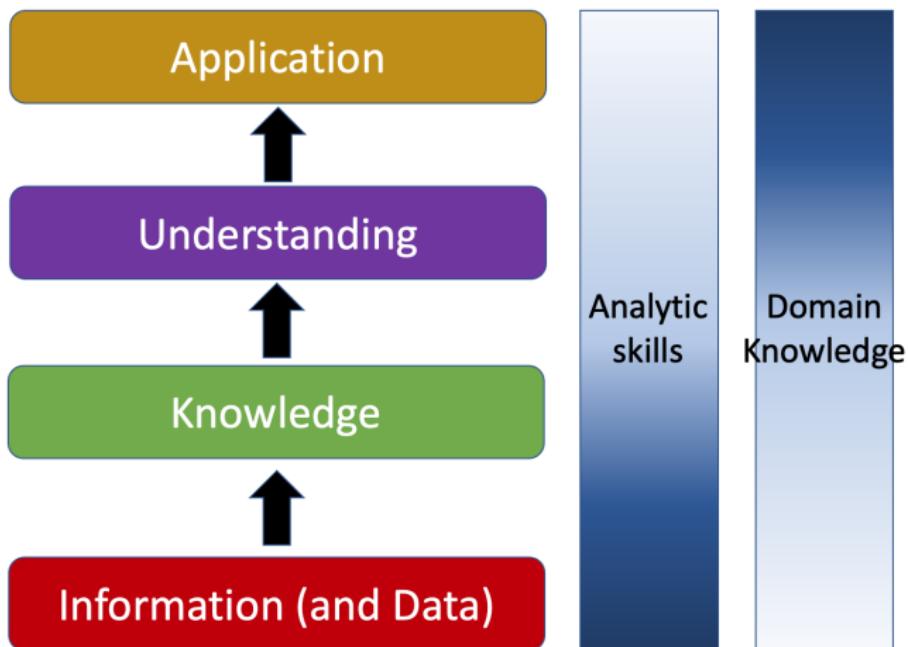
# Domain Knowledge

**Domain knowledge** is knowledge of a specific, specialized discipline or field, in contrast to general (or domain-independent) knowledge. For example, in describing a software engineer may have general knowledge of computer programming as well as domain knowledge about developing programs for a particular industry. People with domain knowledge are often regarded as specialists or experts in their field. (Wikipedia!)

# Analytics Hierarchy



# Analytics Hierarchy



# Session info

```
sessionInfo()
```

```
## R version 4.5.1 (2025-06-13)
## Platform: aarch64-apple-darwin20
## Running under: macOS Tahoe 26.1
##
## Matrix products: default
## BLAS:    /Library/Frameworks/R.framework/Versions/4.5-arm64/Resources/lib/libRblas.0.dylib
## LAPACK:  /Library/Frameworks/R.framework/Versions/4.5-arm64/Resources/lib/libRlapack.dylib;  LAPACK version 3
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: America/Denver
## tzcode source: internal
##
## attached base packages:
## [1] stats      graphics   grDevices utils      datasets   methods    base
##
## loaded via a namespace (and not attached):
## [1] compiler_4.5.1    fastmap_1.2.0    cli_3.6.5       tools_4.5.1
## [5] htmltools_0.5.8.1 rstudioapi_0.17.1 yaml_2.3.10    rmarkdown_2.29
## [9] knitr_1.50        xfun_0.52       digest_0.6.37   rlang_1.1.6
## [13] evaluate_1.0.4
```