



Introduction to Applied Machine Learning in R

Dr. Brian J. Spiering

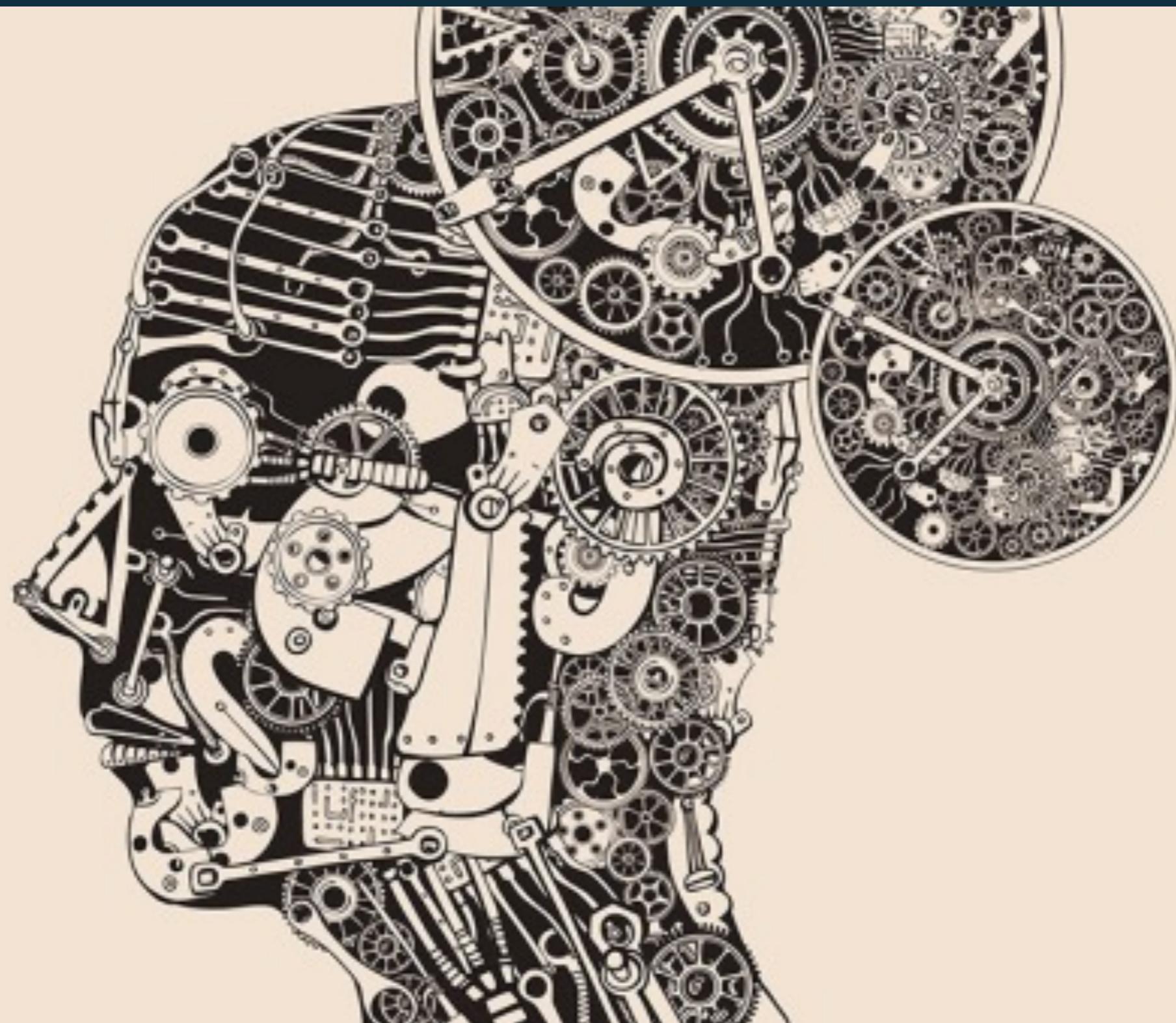


Introduction to Applied Machine Learning in R

Introduction

$$\begin{aligned} & P(Z_{(m,n)} = k | \mathcal{Z}, \mathbf{W}; \alpha, \beta) \\ & \propto P(Z_{(m,n)} = k | \mathcal{Z}, \mathbf{W}; \alpha, \beta) \\ & = \left(\frac{\Gamma(\sum_{i=1}^M \alpha_i)}{\prod_{i=1}^M \Gamma(\alpha_i)} \right)^M \prod_{j \neq m} \frac{\prod_{i=1}^K \Gamma(n_{(i),r}^j + \alpha_i)}{\prod_{i=1}^K \Gamma(n_{(i),r}^j + \alpha_i)} \\ & \quad \left(\frac{\left(\sum_{r=1}^V \beta_r \right)_+}{\prod_{r=1}^V \Gamma(\beta_r)} \prod_{r \neq v} \Gamma(n_{(\cdot),r}^i + \beta_r) \right) \\ & \quad \times \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i + \alpha_m)} \prod_{i=1}^K \frac{\Gamma(n_{(i),v}^i + \beta_v)}{\Gamma(n_{(i),v}^i + \alpha_i + \alpha_m)} \\ & \propto \frac{\prod_{i=1}^K \Gamma(\sum_{r=1}^V n_{(i),r}^i + \alpha_i)}{\Gamma\left(\sum_{i=1}^K n_{m,(\cdot)}^i + \sum_{r=1}^V \sum_{i=1}^K n_{(i),r}^i + \sum_{r=1}^V \beta_r\right)}. \end{aligned}$$

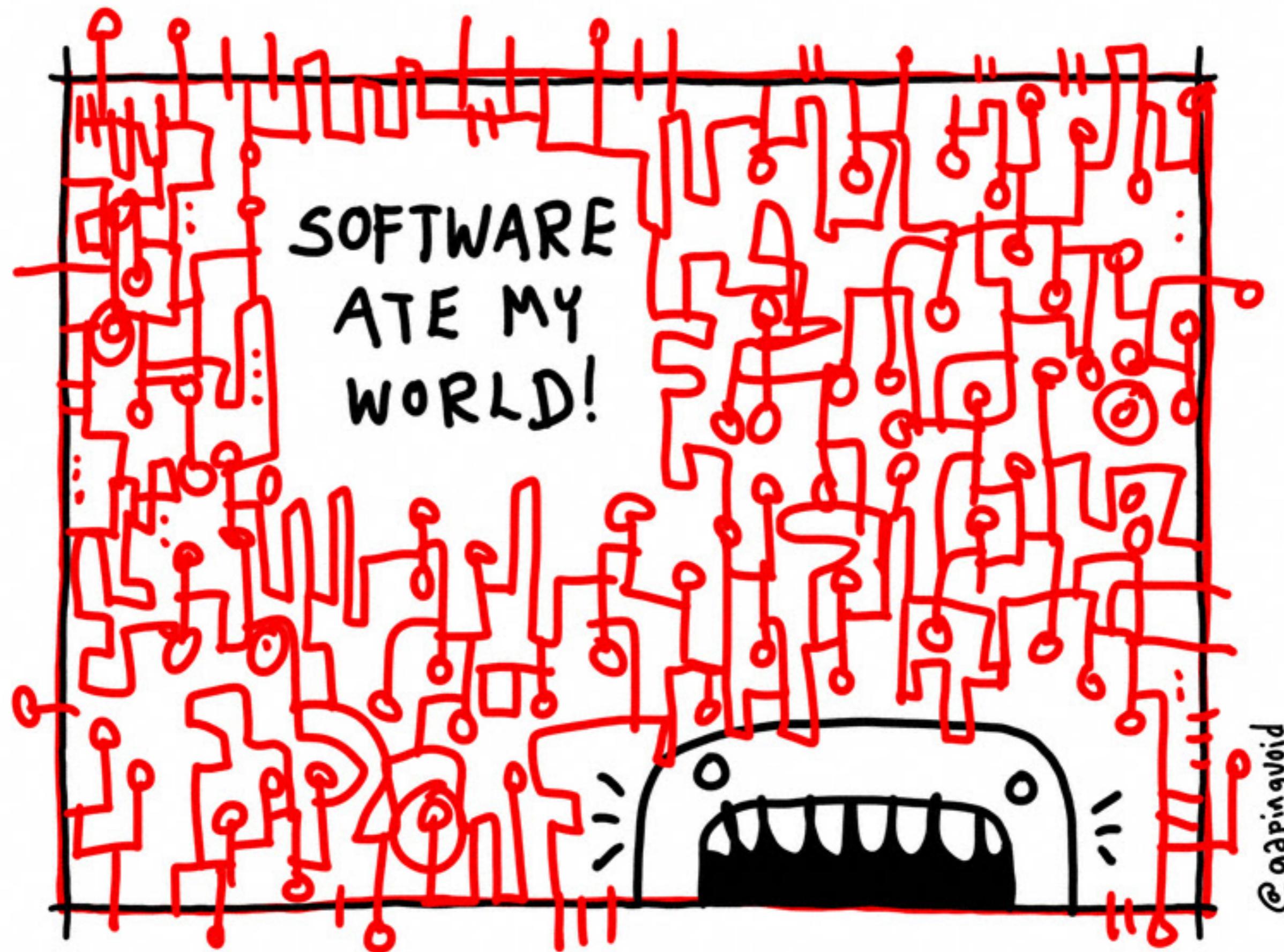
Applied (vs. Theoretical)



What I talk about
when I talk about
machine learning
(apologies to Raymond Carver & Birdman)

machine learning (ml):

Programming computers
to automatically learn and
generalize from examples



Gregoriusoid

Software might eat the world,
but algorithms will digest it

-Bruno Aziza

Machine Intelligence LANDSCAPE

CORE TECHNOLOGIES

ARTIFICIAL INTELLIGENCE



DEEP LEARNING



MACHINE LEARNING



NLP PLATFORMS



PREDICTIVE APIs



IMAGE RECOGNITION



SPEECH RECOGNITION



RETHINKING ENTERPRISE

SALES



SECURITY / AUTHENTICATION



FRAUD DETECTION



HR / RECRUITING



MARKETING



PERSONAL ASSISTANT



INTELLIGENCE TOOLS



RETHINKING INDUSTRIES

ADTECH



AGRICULTURE



EDUCATION



FINANCE



LEGAL



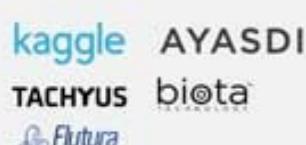
MANUFACTURING



MEDICAL



OIL AND GAS



MEDIA / CONTENT



CONSUMER FINANCE



PHILANTHROPIES



AUTOMOTIVE



DIAGNOSTICS



RETAIL



RETHINKING HUMANS / HCI

AUGMENTED REALITY



GESTURAL COMPUTING



ROBOTICS



EMOTIONALrecognition



HARDWARE



DATA PREP

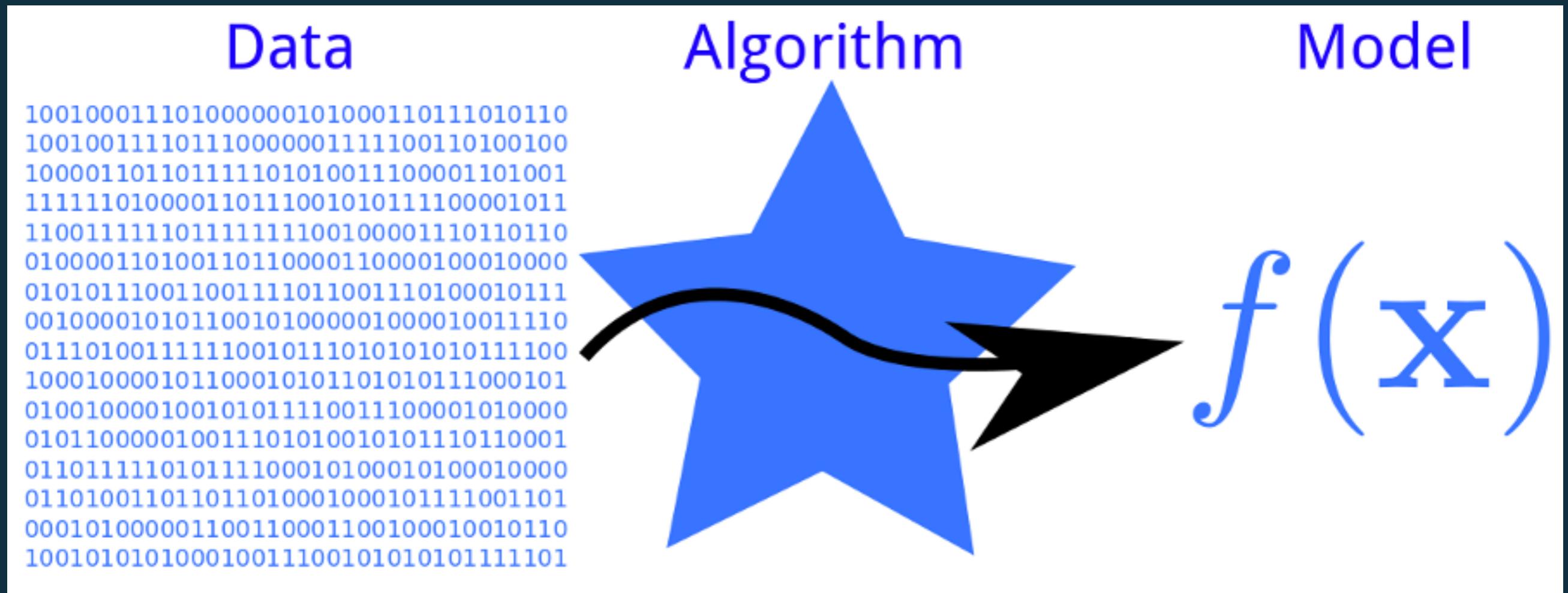


DATA COLLECTION



PP

ml process





Take a step back.
Look at the bigger picture.

- Frank Underwood

Analytics Workflow

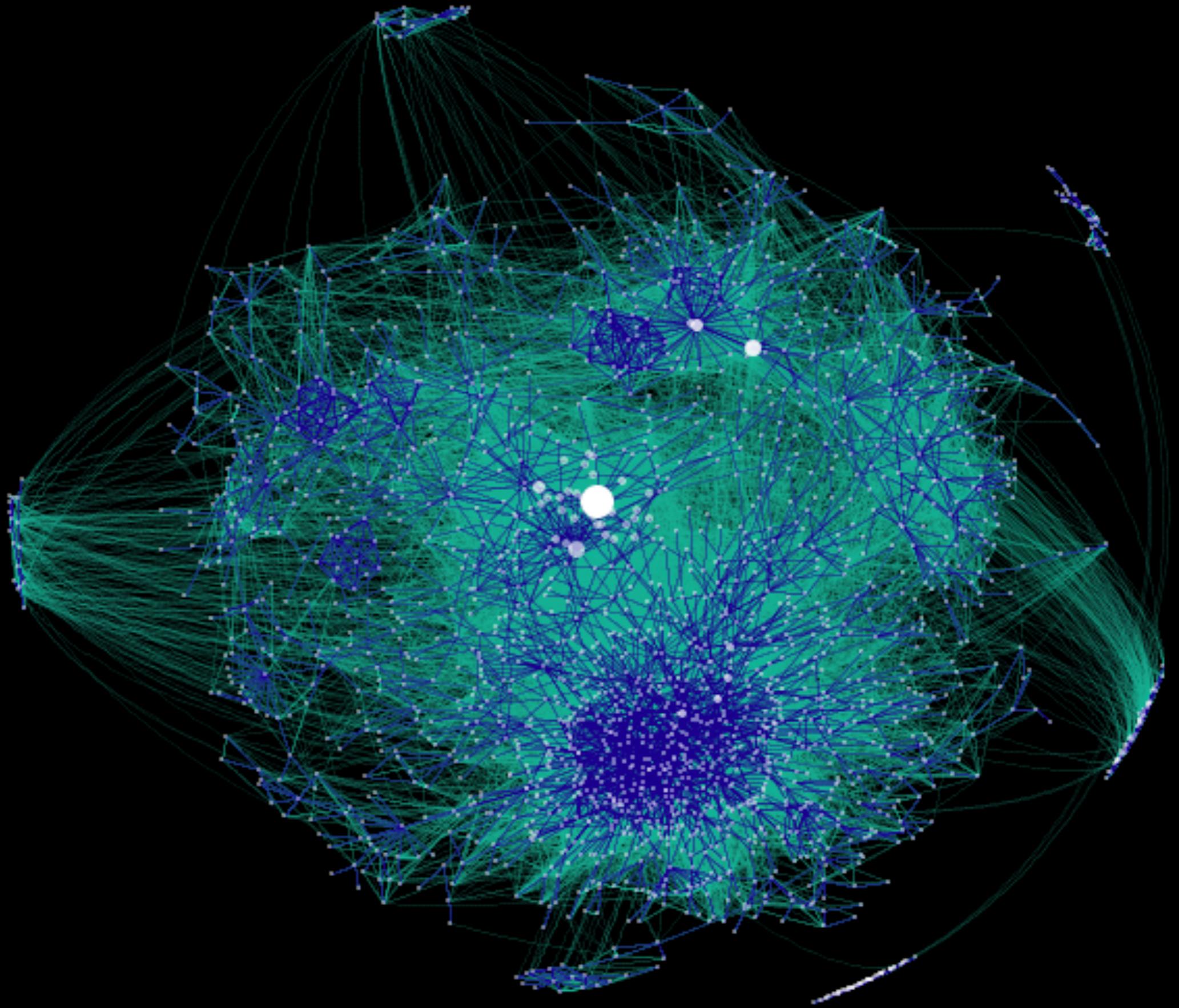
A black and white portrait of W. Edwards Deming, an elderly man with glasses, wearing a suit and tie.

If you do not know how to ask the right question,
you discover nothing.

(W. Edwards Deming)

1. Ask the right question

2. Find the right data





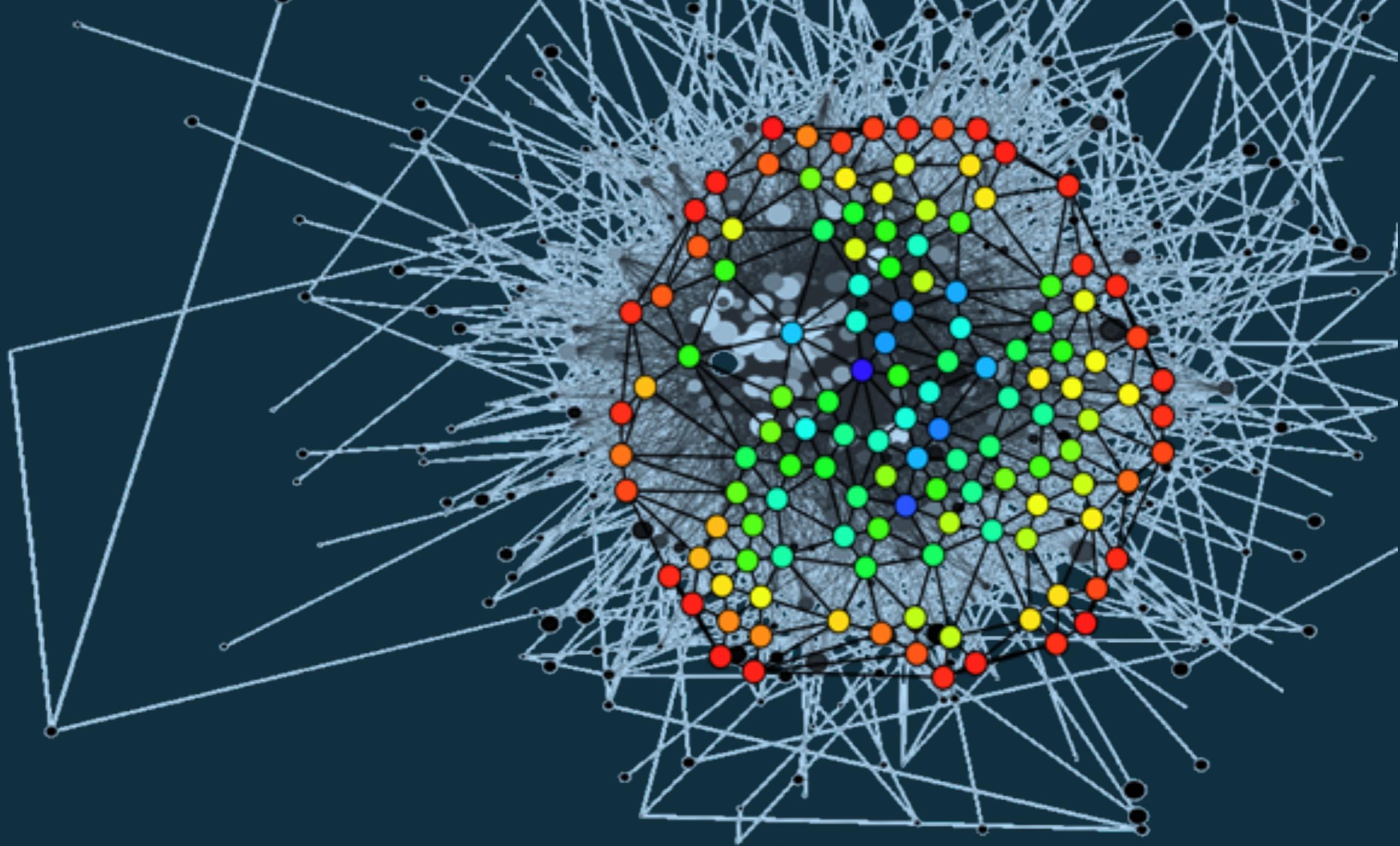
Data Manipulation with R

Perform group-wise data manipulation and deal with large datasets using R efficiently and effectively

Jaynal Abedin

[PACKT] open source[®]
PUBLISHING

3. Process the data



4. Model the data

5. Share the results



Analytics Workflow

1. Ask the right question
2. Find the right data
3. Process the data
4. Model the data
5. Share the results

ml =

representation + evaluation + optimization

Representation

We need to represent the data and model
in a way computers can understand.

Evaluation

Separate good models from bad models

Optimization

A method to search among models.

Efficiently search for effective models.

Paradigms in ml

- Supervised Learning
- Unsupervised Learning

Supervised Learning

Thick



Thin



Infer a function from labeled data

Unsupervised Learning



Find hidden structure in unlabeled data



Linear Regression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + \varepsilon$$

Fit a line to the data

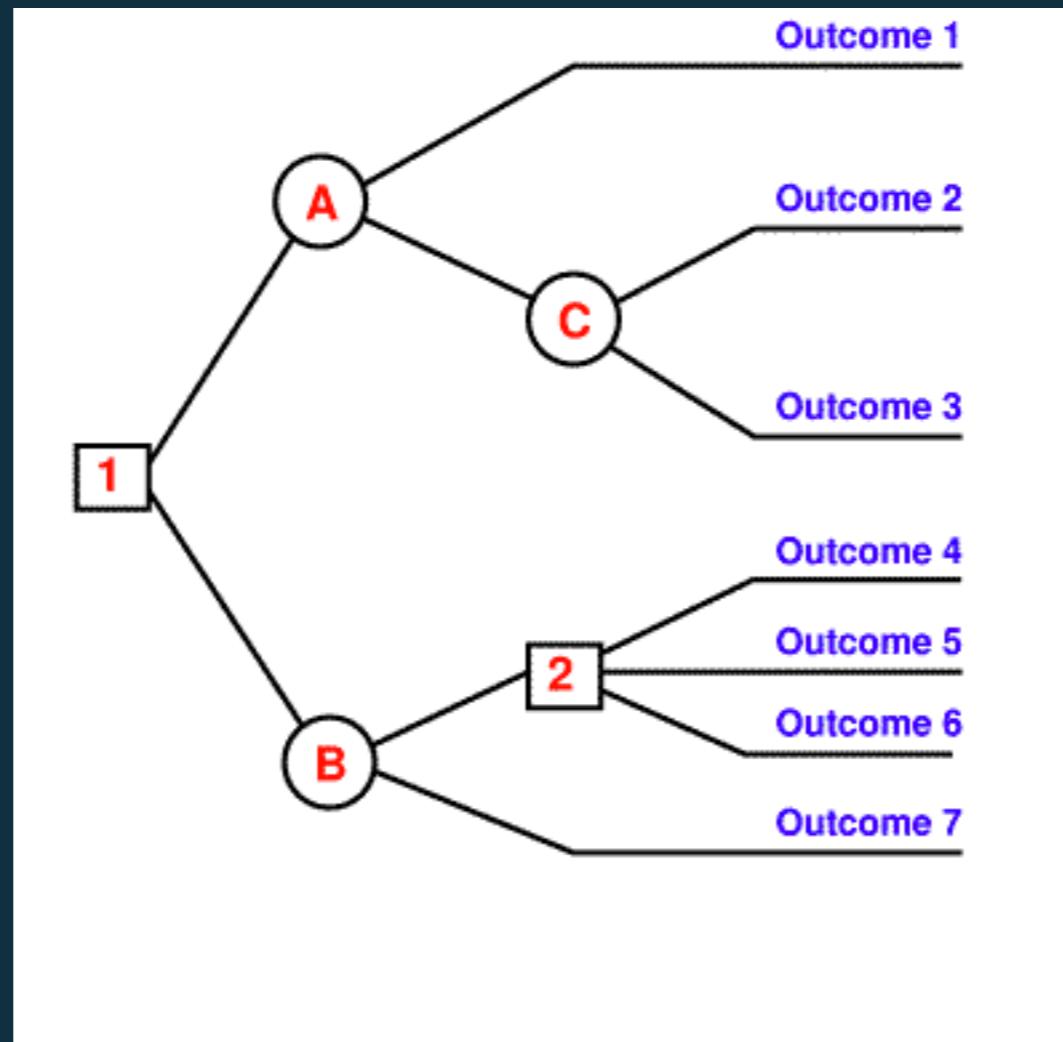
R code interlude

Pros & Cons

- + Easy to fit
- + Easy to interpret
- Life is nonlinear
- Limited to predicting numeric output



Decision Trees



A series of branches that divide the data

R code interlude

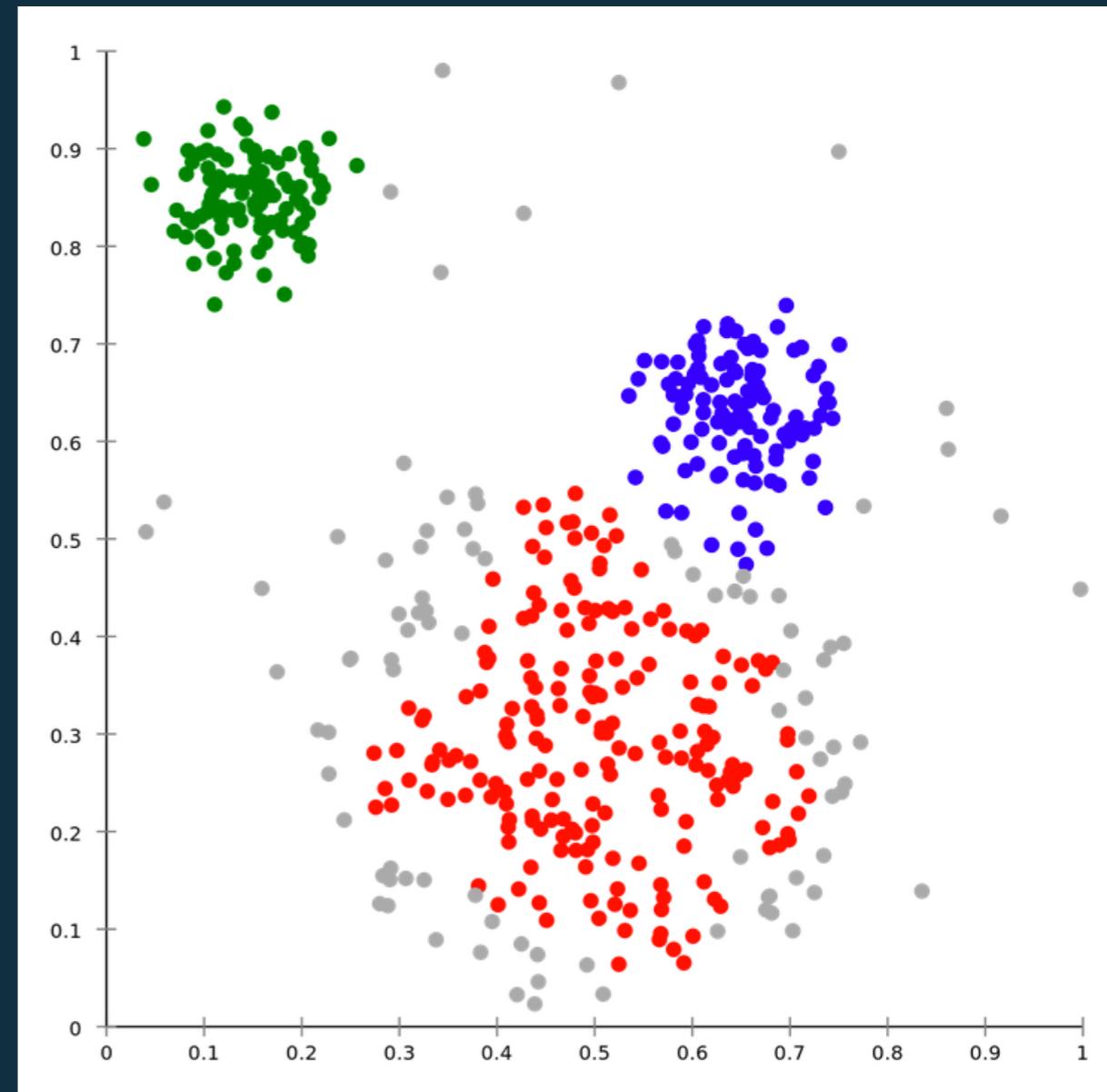
Pros & Cons

- + White box model
- + Can fit most types of data
- Fragile
- Overfitting
- Out-of-sample prediction



Clustering

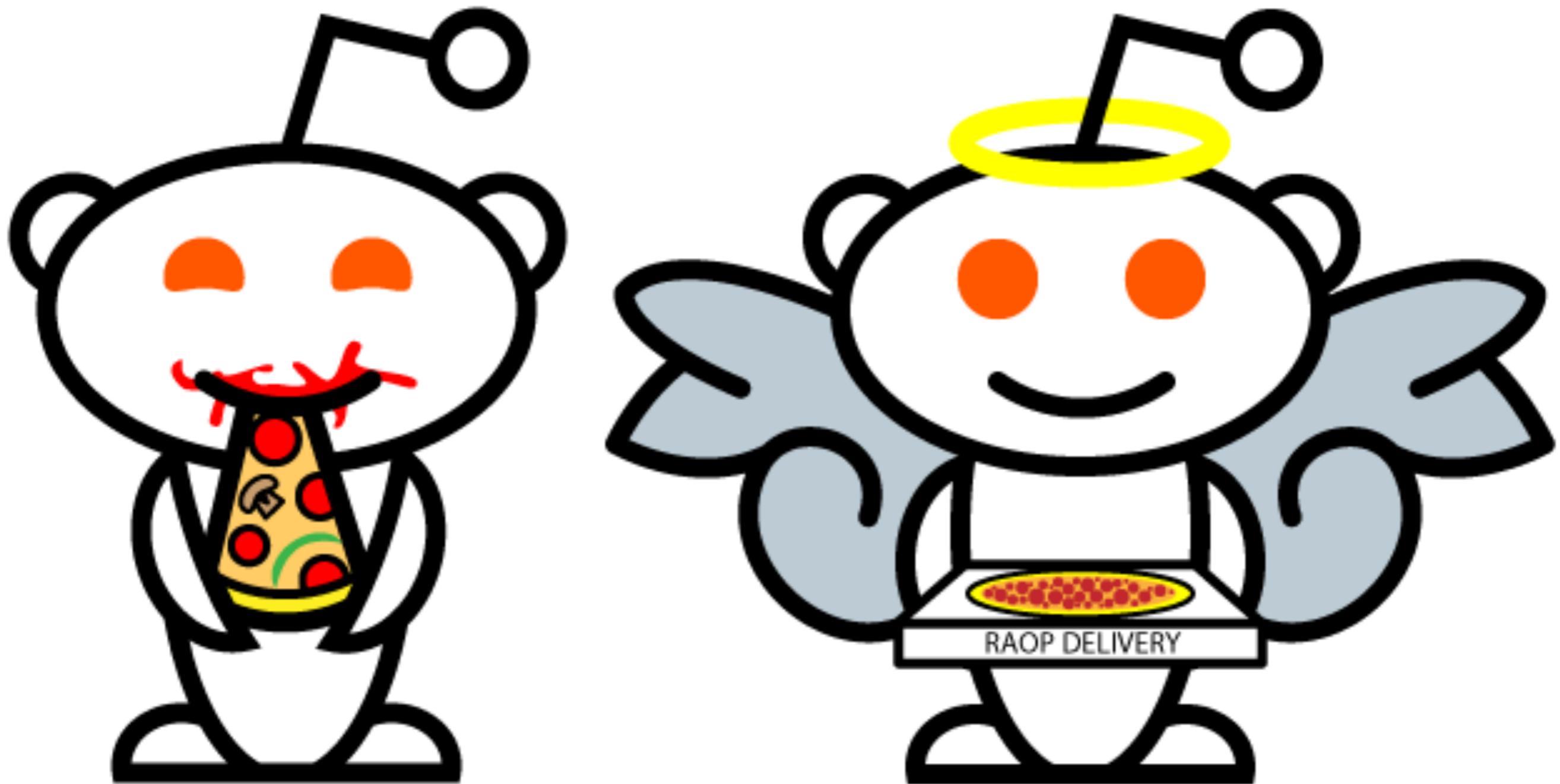
Grouping data so that data in the same group is more similar to each other than to data in the other groups



R code interlude

Pros & Cons

- + Powerful and common unsupervised model
- + Good for a quick insights into data
- Difficult to identify or define clusters
- Slow to fit



ROAP
(Random Acts of Pizza)

Suggested Resources:

general ml

1. Data Mining and Analysis
2. Machine Learning: A Probabilistic Perspective

Suggested Resources: R

1. Try R from Code School
2. R in a Nutshell
3. The Art of R Programming

Suggested Resources: ml in R

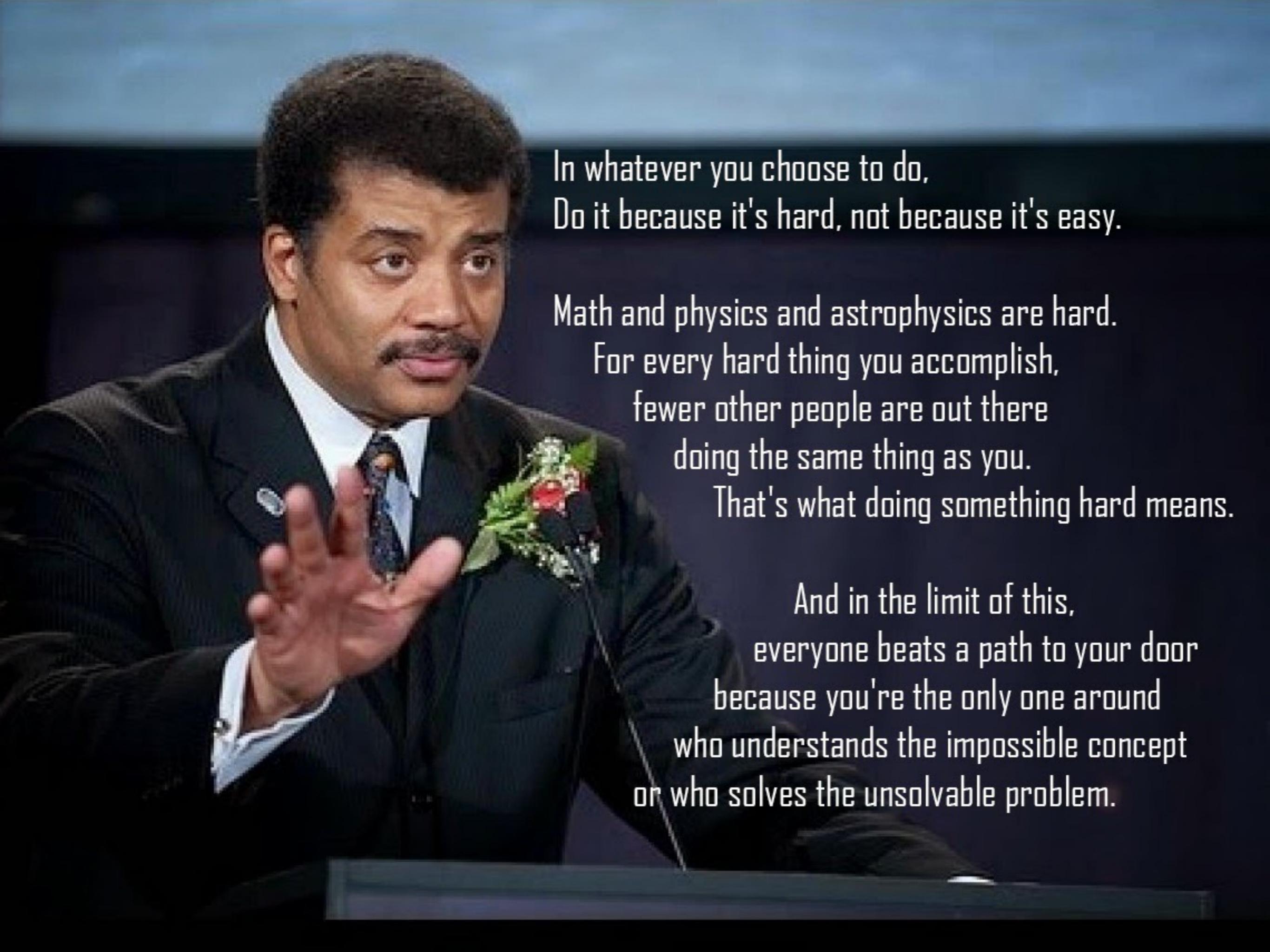
- An Introduction to Statistical Learning:
with Applications in R
 - Book
 - Video & Slides
 - Applied Predictive Modeling



LiveCareer
is Hiring



drop me a line.



In whatever you choose to do,
Do it because it's hard, not because it's easy.

Math and physics and astrophysics are hard.
For every hard thing you accomplish,
fewer other people are out there
doing the same thing as you.
That's what doing something hard means.

And in the limit of this,
everyone beats a path to your door
because you're the only one around
who understands the impossible concept
or who solves the unsolvable problem.