# INTRO TO DATA SCIENCE

# ANDREW WORSLEY

**National Data Science Lead - Velrada**

**Lead Instructor (Part-time), Data Science - GA**

Expertise in production-grade machine learning, high-performance cloud computing and statistical analysis.

generalassemb.ly/instructors/andrew-worsley/16428
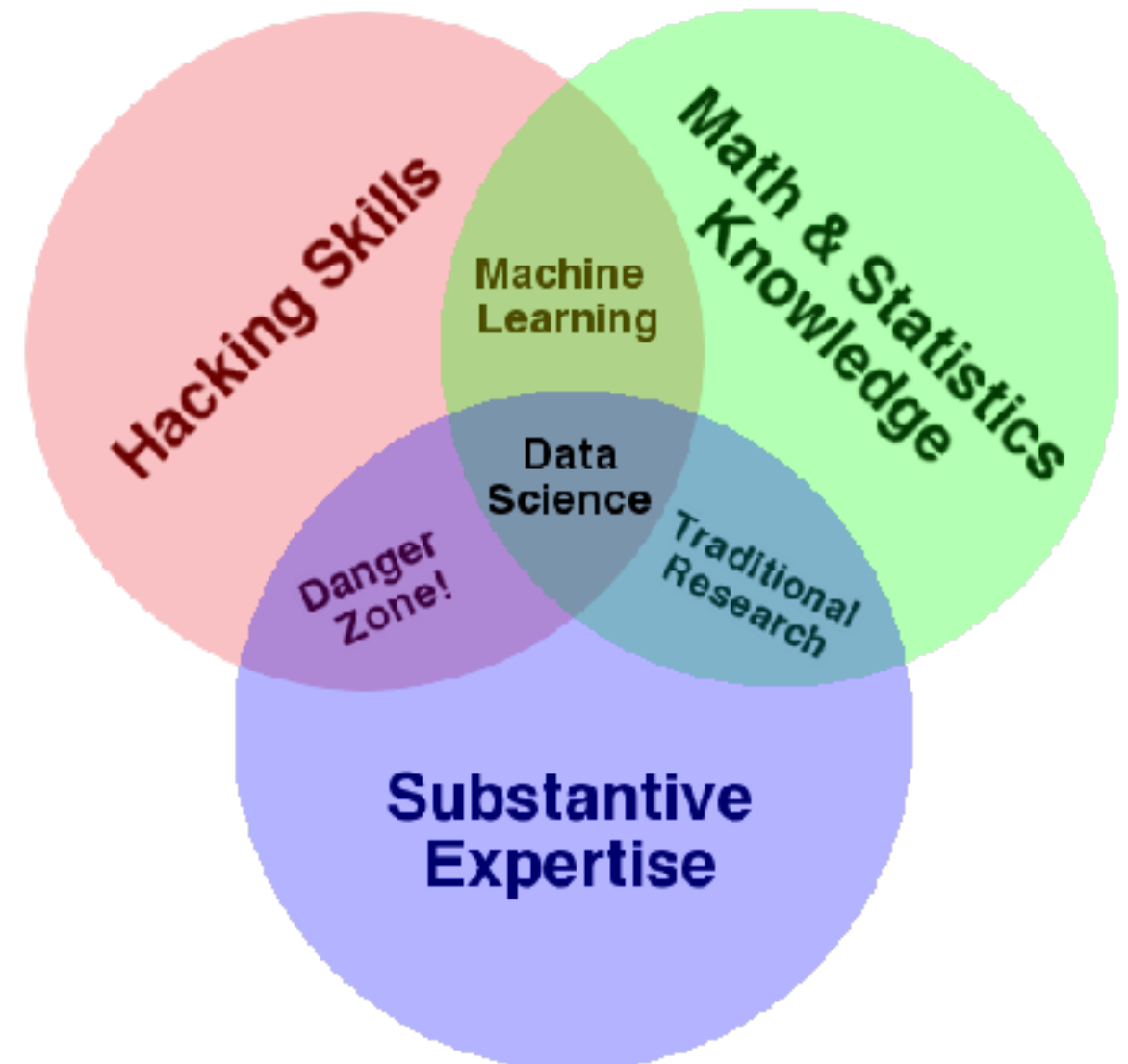
www.linkedin.com/in/andrew-worsley/

# WHAT IS DATA SCIENCE?

# WHAT IS A DATA SCIENTIST?

‣ "Data Scientist' is a Data Analyst who lives in California"

‣ "A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician."

‣ Someone who can collect, statistically explore and analyse data in an efficient and reproducible manner… but who can also translate from Dataese to Peoplese. Oh, and something something machine learning.

# WHAT IS DATA SCIENCE?

‣A set of tools and techniques for data

‣Interdisciplinary problem-solving

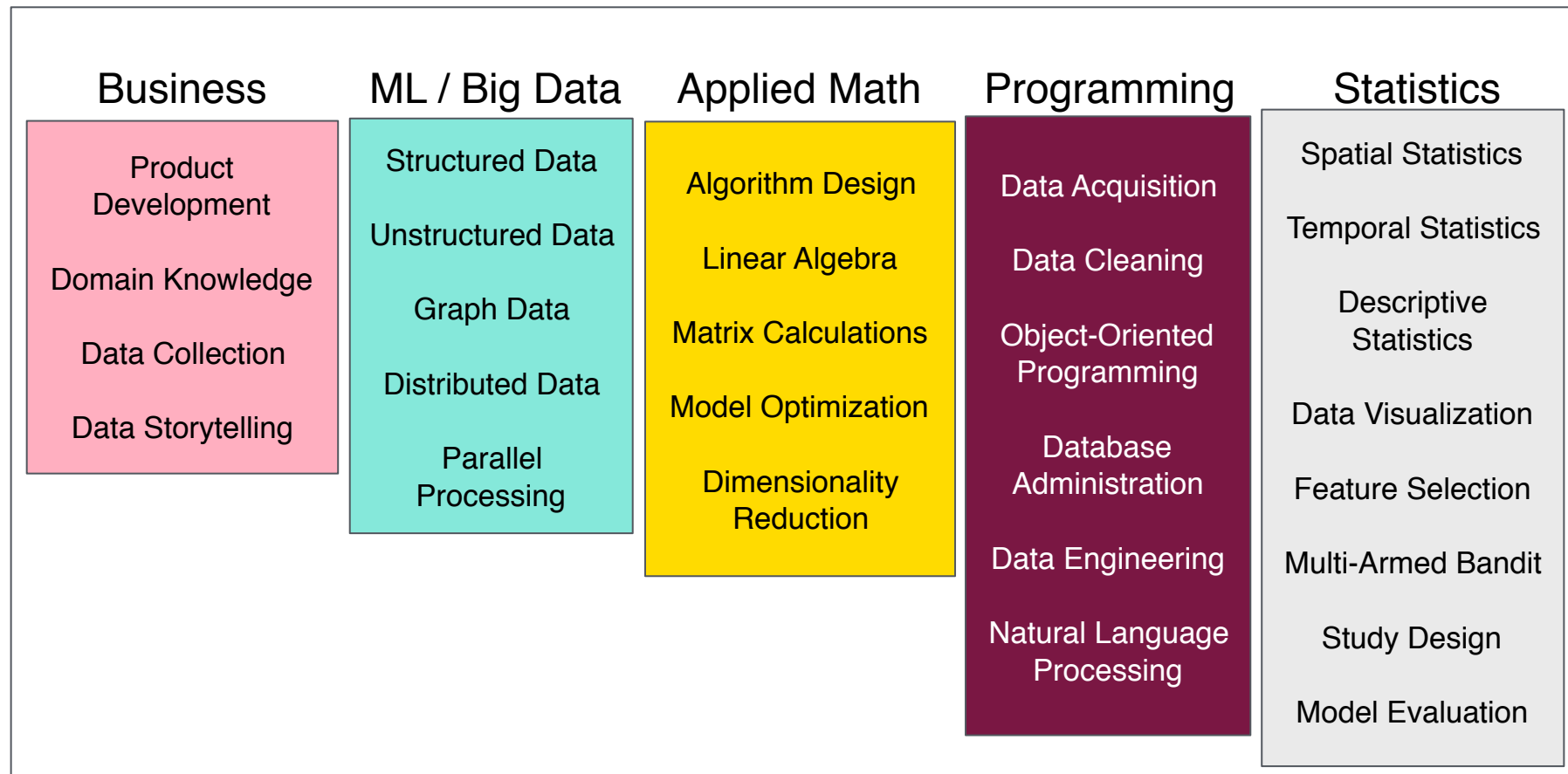‣Application of scientific techniques to practical problems

# WHAT ARE THE ROLES IN DATA SCIENCE?

‣ Data Science involves a variety of roles, not just one.

| | | | |
|---|---|---|---|
| Data Developer | Developer | Engineer | |
| Data Researcher | Researcher | Scientist | Statistician |
| Data Creative | Jack of All Trades | Artist | Hacker |
| Data Businessperson | Leader | Businessperson | Entrepeneur |

# WHAT ARE THE ROLES IN DATA SCIENCE?

▸ Data Science involves a variety of skill sets, not just one.

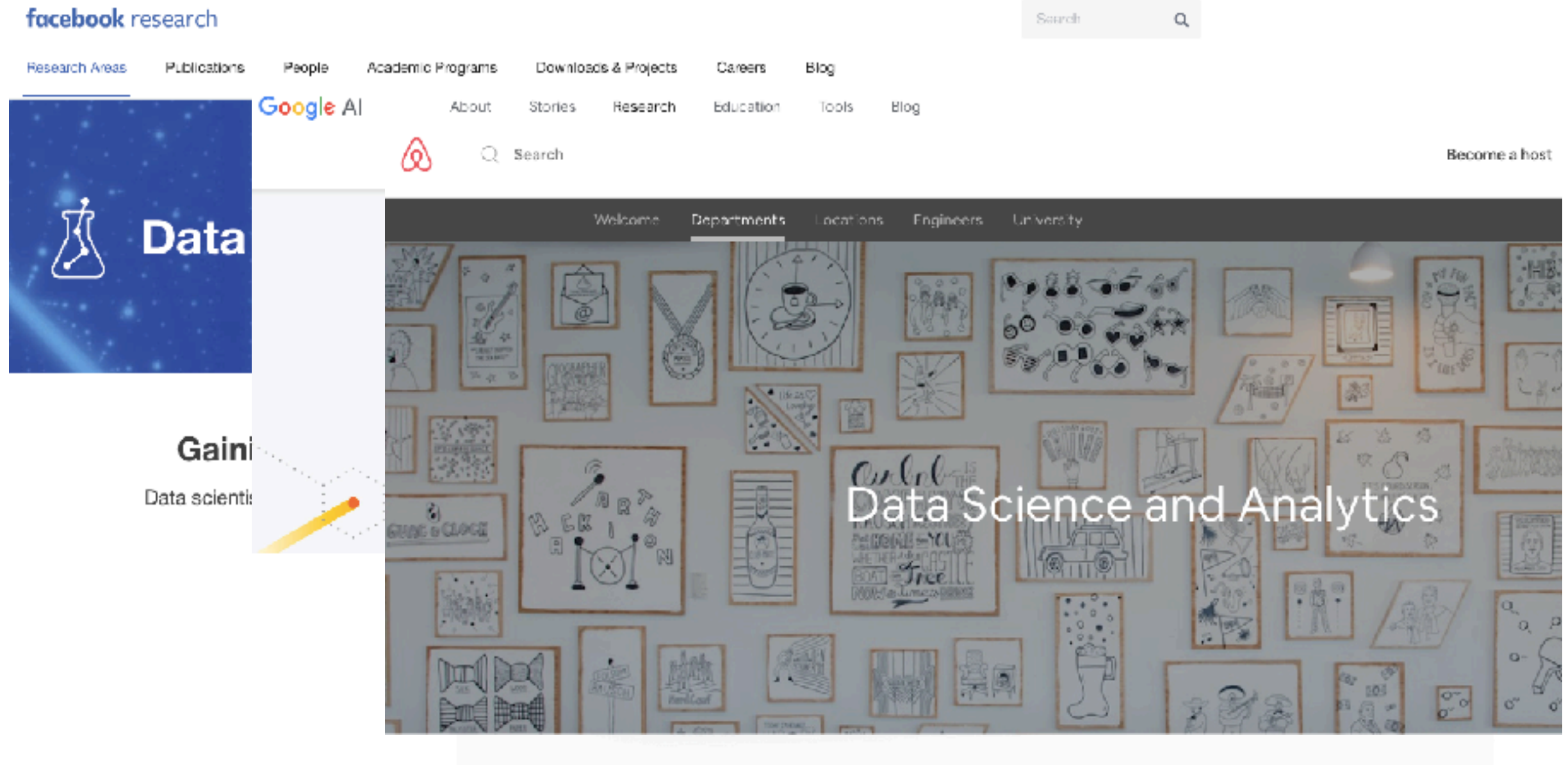| Business | ML / Big Data | Applied Math | Programming | Statistics |
|---|---|---|---|---|
| Product Development | Structured Data | Algorithm Design | Data Acquisition | Spatial Statistics |
| Domain Knowledge | Unstructured Data | Linear Algebra | Data Cleaning | Temporal Statistics |
| Data Collection | Graph Data | Matrix Calculations | Object-Oriented Programming | Descriptive Statistics |
| Data Storytelling | Distributed Data | Model Optimization | Database Administration | Data Visualization |
| | Parallel Processing | Dimensionality Reduction | Data Engineering | Feature Selection |
| | | | Natural Language Processing | Multi-Armed Bandit |
| | | | | Study Design |
| | | | | Model Evaluation |

# WHAT ARE THE ROLES IN DATA SCIENCE?

‣ These roles prioritize different skill sets.

‣ However, all roles involve some part of each skillset.
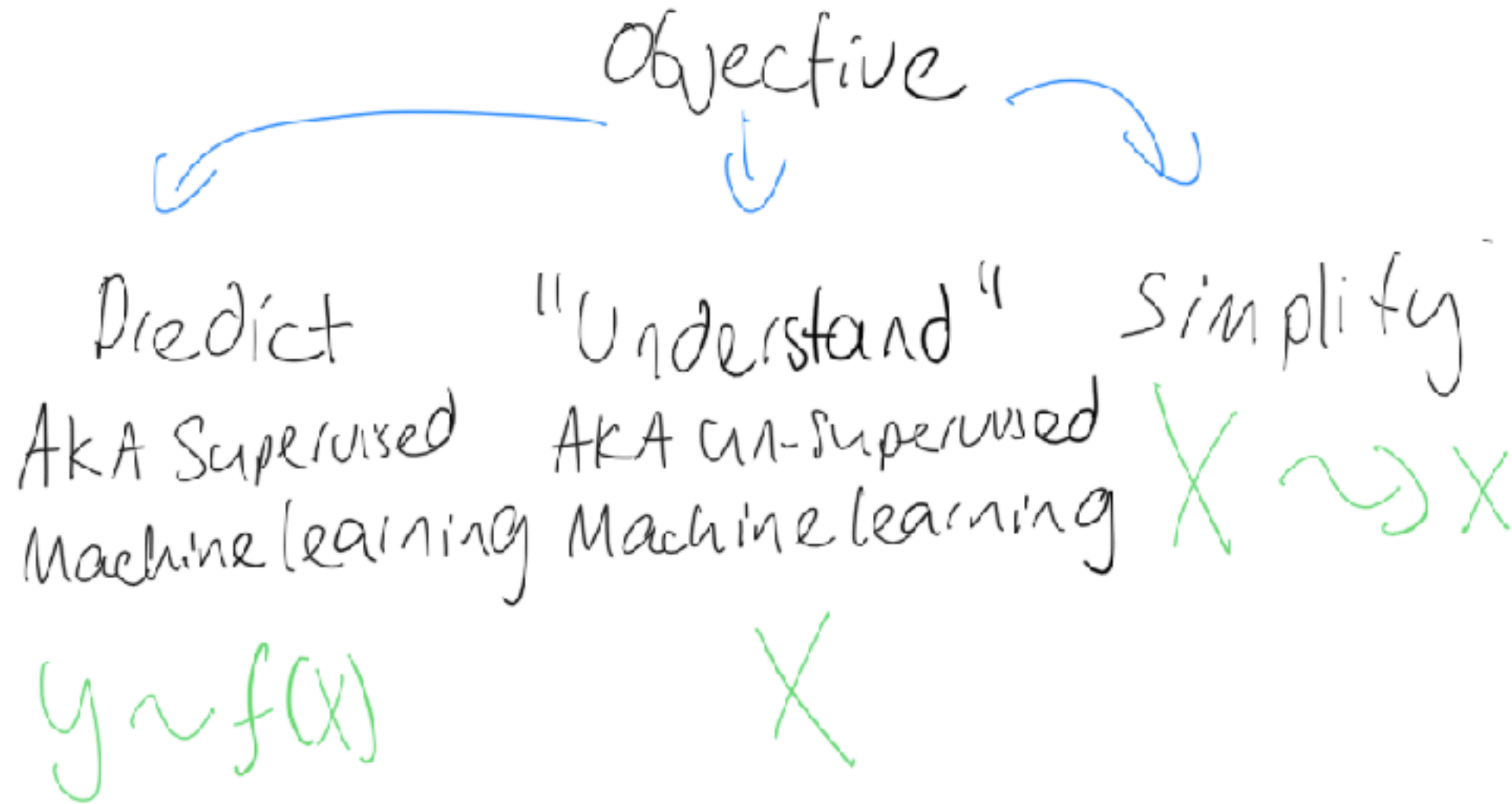
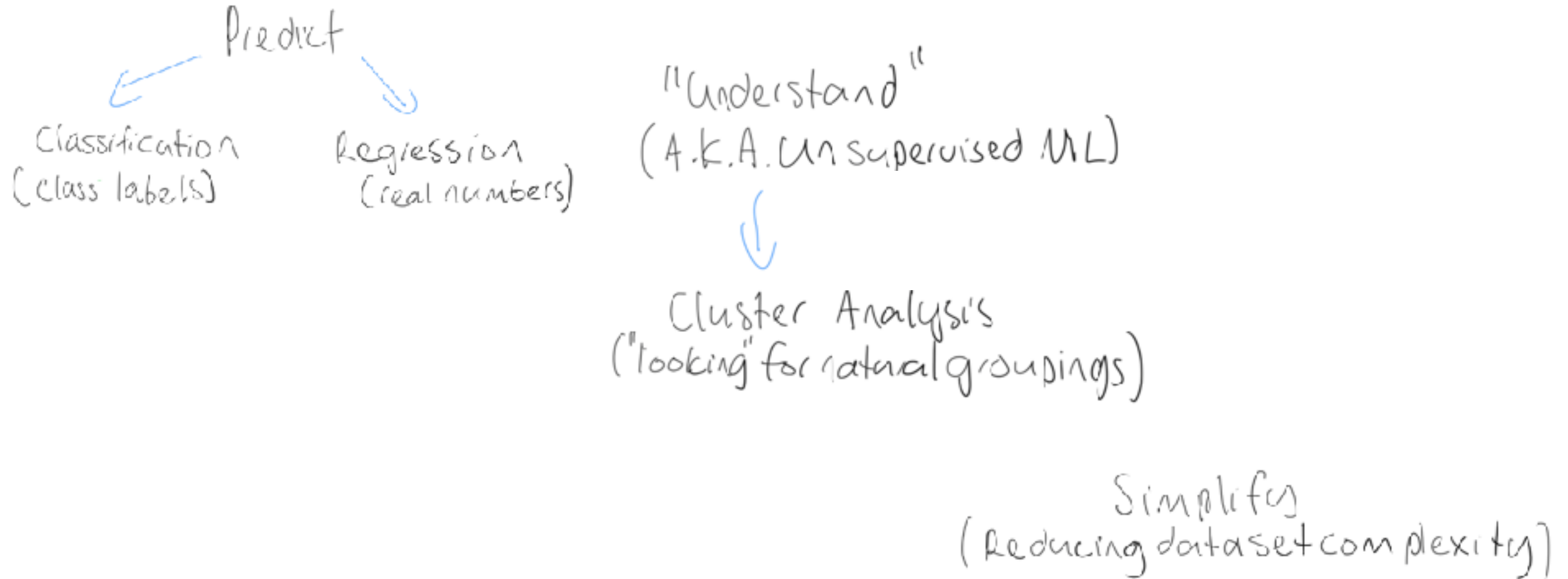‣ Where are your strengths and weaknesses?

# DATA SCIENCE A PRIORTIY

## WHY NOW?

'Big Data' + Algorithms + **Computing Power** = Data Science

# PROBLEMS WE SOLVE

Objective

Predict     "Understand"     Simplify

AKA Supervised    AKA Un-Supervised

Machine Learning   Machine Learning    $X \rightsquigarrow X$

$y \sim f(x)$        $X$

# PROBLEMS WE SOLVE

Predict

Classification
(class labels)

Regression
(real numbers)

"Understand"
(A.K.A. Unsupervised ML)

Cluster Analysis
("looking" for natural groupings)

Simplify
(Reducing dataset complexity)

# THE DATA SCIENCE WORKFLOW

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

▸ A methodology for doing Data Science

▸ Similar to the scientific method

▸ Helps produce *reliable* and *reproducible* results

　▸ *Reliable*:  Accurate findings

　▸ *Reproducible*:  Others can follow your steps and get the same results

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

The steps:

1. Identify the problem
2. Acquire the data
3. Parse the data
4. Mine the data
5. Refine the data
6. Build a data model
7. Present the results

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

**Identify**

## IDENTIFY THE PROBLEM

- ☐ Identify business/product objectives
- ☐ Identify and hypothesize goals and criteria for success
- ☐ Create a set of questions for identifying correct data set

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

## ACQUIRE THE DATA

- ☐ Identify the "right" data set(s)
- ☐ Import data and set up local or remote data structure
- ☐ Determine most appropriate tools to work with data

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

**Parse**

## PARSE THE DATA

- ☐ Read any documentation provided with the data
- ☐ Perform exploratory data analysis
- ☐ Verify the quality of the data

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

**Mine**

## MINE THE DATA

- ☐ Determine sampling methodology and sample data
- ☐ Format, clean, slice, and combine data in Python
- ☐ Create necessary derived columns from the data (new data)

# OVERVIEW OF THE DATA SCIENCE WORKFLOW



**Build**

## BUILD A DATA MODEL

- ☐ Select appropriate model
- ☐ Build model
- ☐ Evaluate and refine model

# OVERVIEW OF THE DATA SCIENCE WORKFLOW

**Present**

**PRESENT THE RESULTS**

☐ Summarize findings with narrative, storytelling techniques

☐ Present limitations and assumptions of your analysis

☐ Identify follow up problems and questions for future analysis

# EXAMPLE DATA SCIENCE TOOLBOX

# LETS CODE!