



Applied Supervised Learning for Cyber Security Module 0: Introduction

Presented by Brian Genz, CISSP, GREM, GCIH, GNFA, GCFA
brian.genz@gtkcyber.com

Course Agenda

Day 1: Introduction & Data Engineering

- Intro: What is Data Science?
- What is Machine Learning?
 - Overview of Machine Learning & Cyber Applications
- Overview of Data Preparation and Feature Engineering
- Data Visualization

Day 2: Supervised Learning

- Overview of Supervised Learning techniques

Day 3: Machine Learning

- Feature Engineering
- Supervised Machine Learning
- Unsupervised Machine Learning
- Hacking Machine Learning Models

Day 4: Advanced Topics

- Introduction to Big Data Tools
- Anomaly Detection
 - PySpark, ELK & Kafka
- Hunting with Data Science

Expectations

- Please participate and **ask questions**.
- Please follow along and **TRY OUT** the examples yourself during the class
- All the answers are in the slide decks or GitHub repository, but please try to complete the exercises **without looking at the answers**.
- Join the conversation in slack!
- Have fun!

Introduction

Our Lawyers Make Us Say This



All materials presented in this training and those provided as an adjunct to the program are copyrighted 2020 by GTK Cyber LLC.

They are intended solely for the use of registered program participants and may not be reproduced or redistributed in any manner for any other reason.

Brian Genz

Brian leads the Threat & Vulnerability Management team at Splunk and serves as a Senior Instructor with GTK Cyber.

Prior to leading the red team at Splunk, he led threat hunting and security orchestration, automation and response (SOAR) efforts with Splunk Phantom at a Fortune 100 financial company. He has experience in the defense intelligence, manufacturing, and financial sectors and has worked in both offensive and defensive security. Degrees and certifications include MBA, M.S. in Information Technology Management, CISSP, GREM, GNFA, GCFA, and GCIH.

Brian has taught with GTK Cyber at Black Hat USA and in Europe. He has presented at multiple industry conferences, including the DEFCON AI Village, the Conference on Applied Machine Learning for Information Security (CAMLIS), Derby Con, Circle City Con, Cyphercon, the ISSA International Conference, ISACA Milwaukee, Infragard, and others. He has served as adjunct faculty at Wayne State University and Milwaukee Area Technical College. He is also featured in the book, "Tribe of Hackers: Red Team Edition."

Who are you?

- Your name (or what you want us to call you)
- Your job role
- What you hope to get out of this class
- Your level of experience with coding

What is Data Science?

**Data Science is the
automated extraction of
information from raw data.**

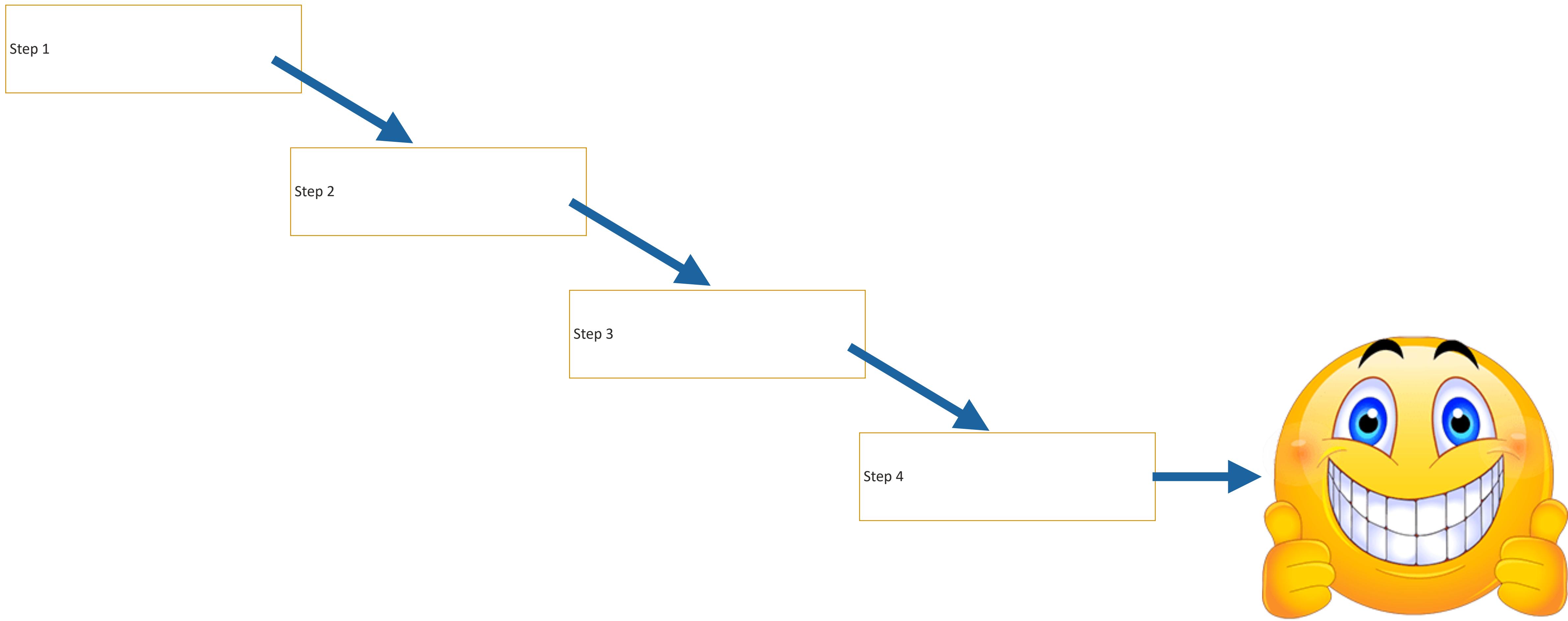
Data Science is the art of **turning data into actions. This is accomplished through the creation of data products, which provide actionable information without exposing decision makers to the underlying data or analytics**

Booz Allen Hamilton, Field Guide to Data Science, Pg. 17

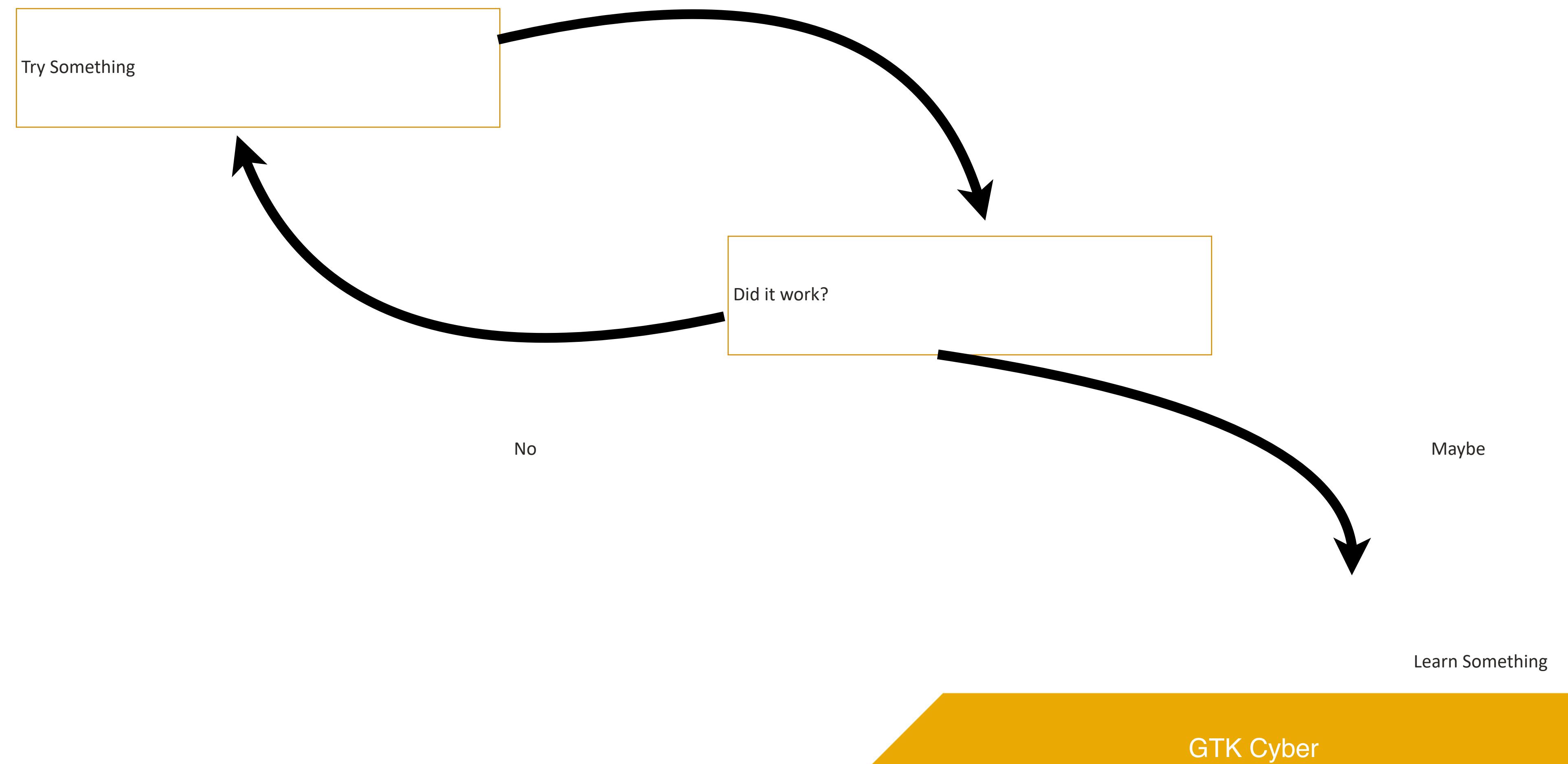
Analyst ← → Developer

Analyst + Developer

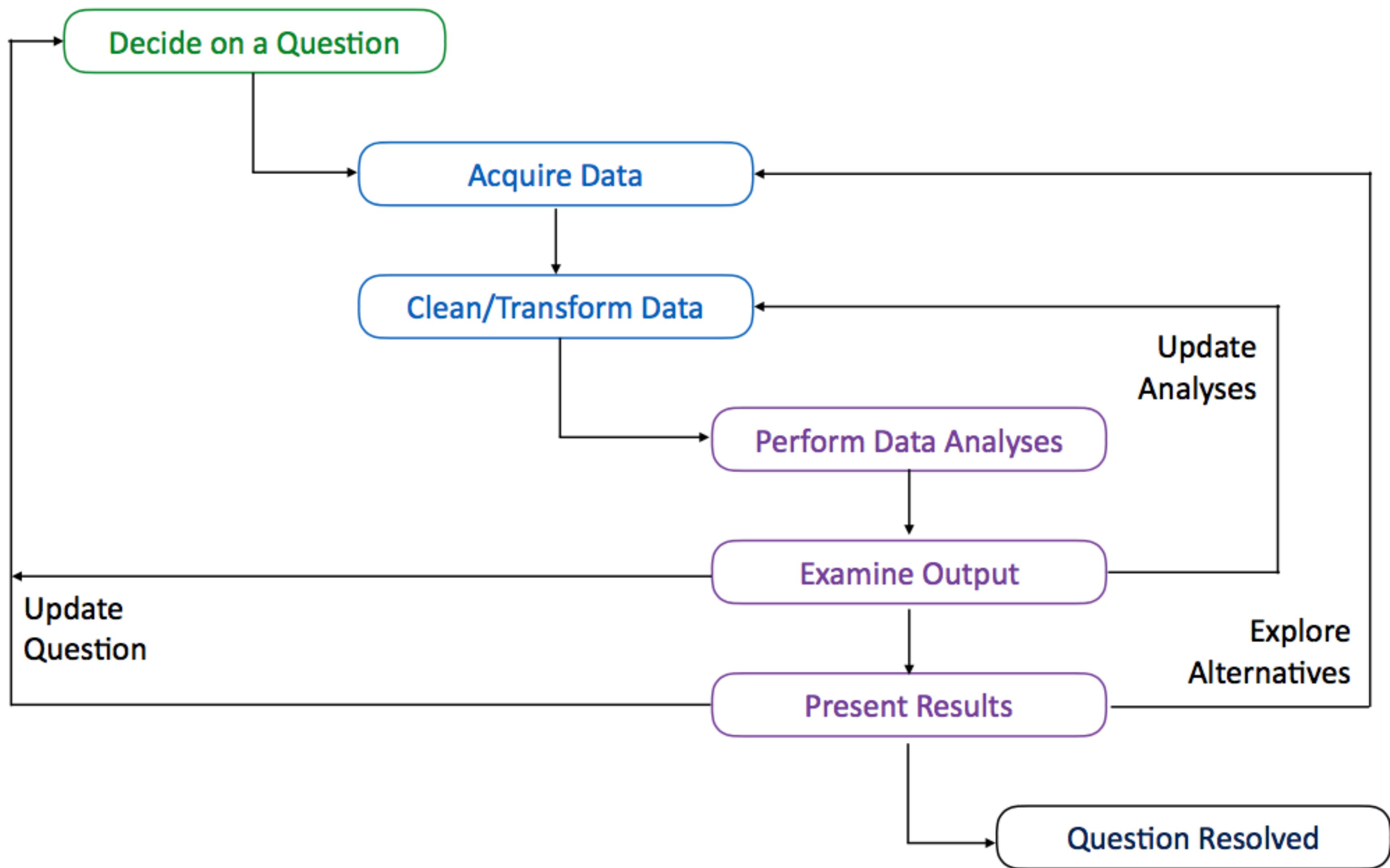
What Data Science is Not



What Data Science Is

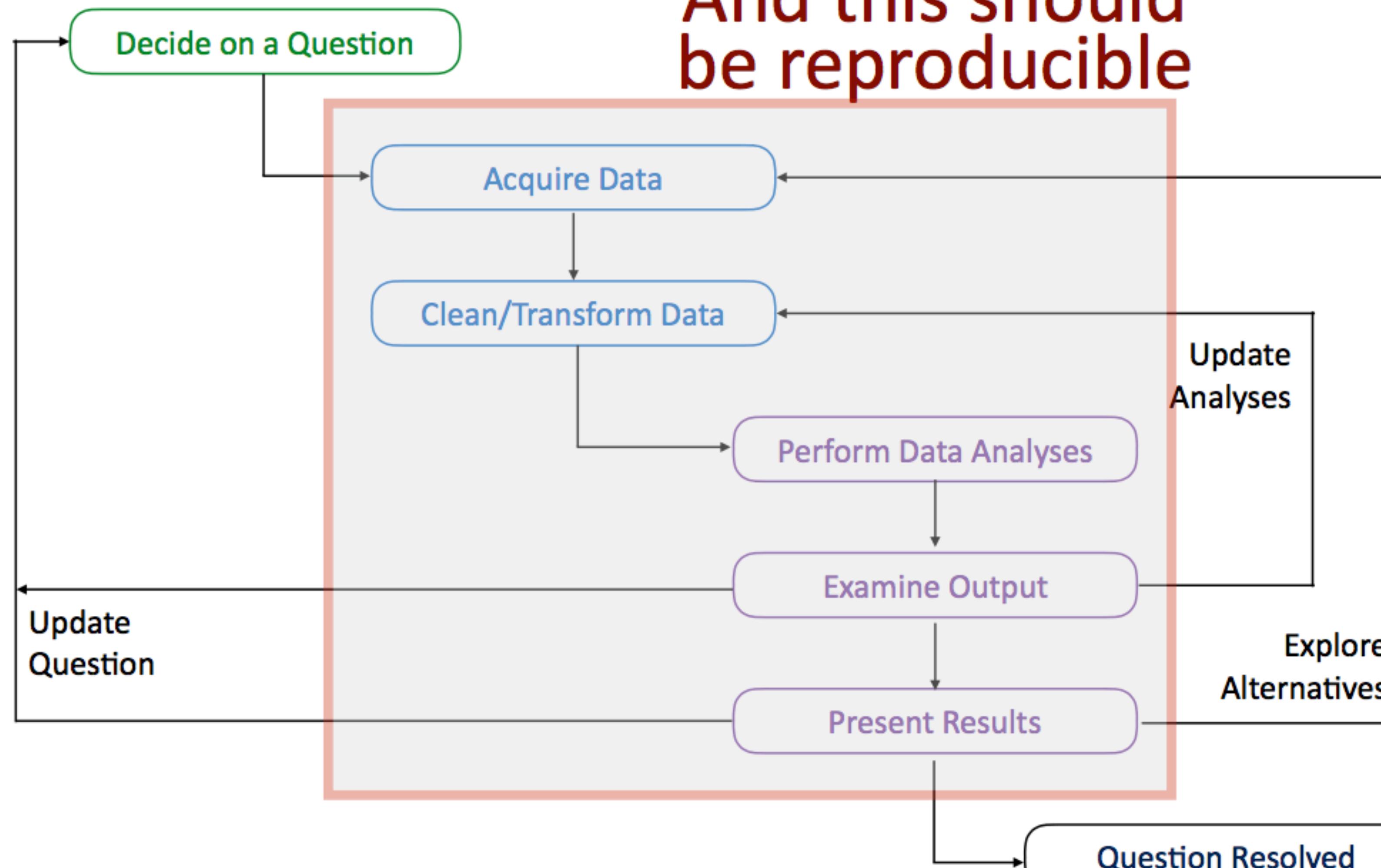


Research Process



Research Process

And this should
be reproducible



"The term "data scientist" will subside and may well sound dated five years from now. **The skills will become more commonplace and commoditized. When that happens, the real boom will begin**, because the technology will become widely adopted and thus more useful. **Instead, we need self-service tools that empower smart and tenacious business people to perform Big Data analysis themselves.**

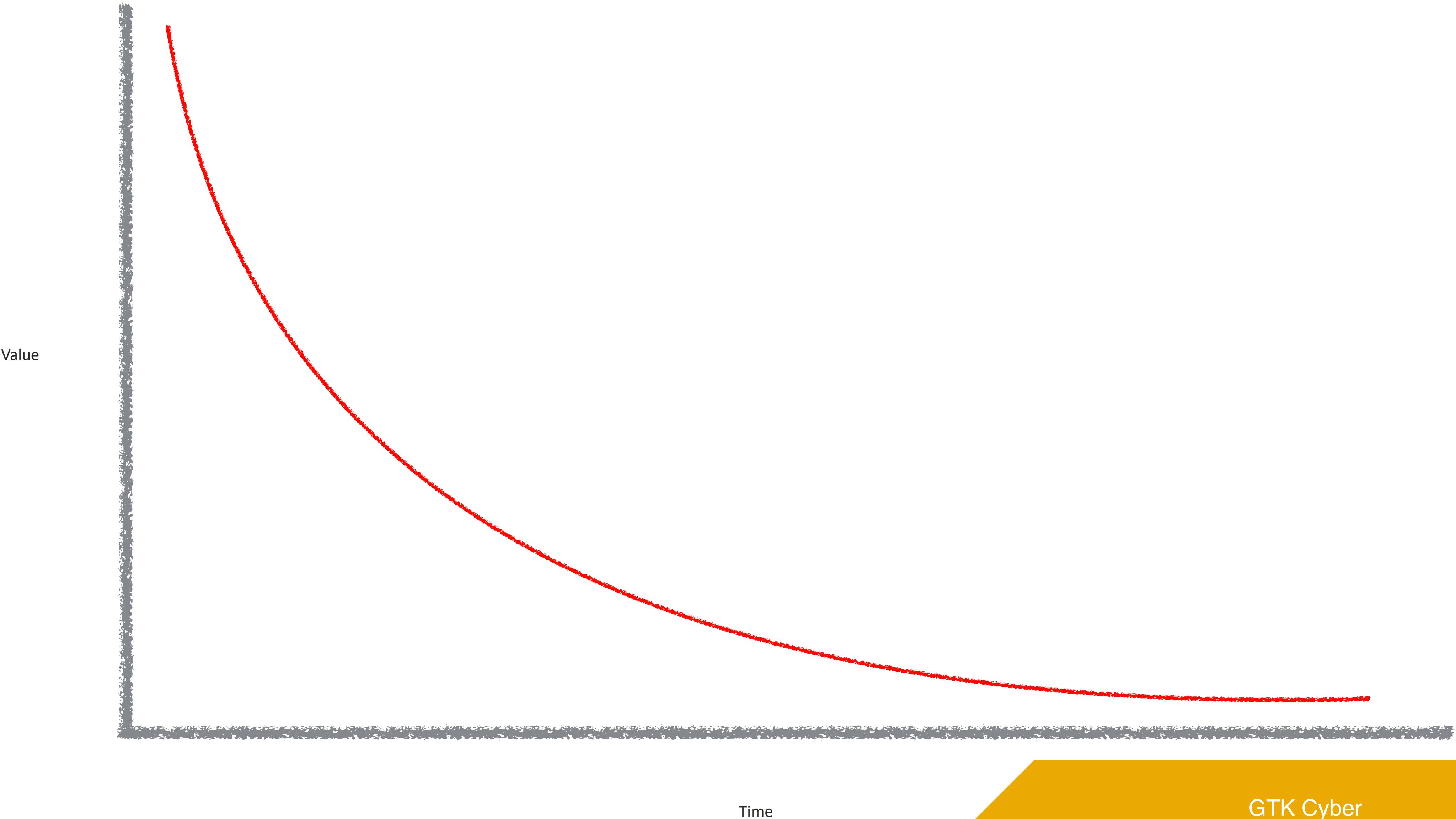
—Andrew Brust, "Data scientists don't scale", <http://www.zdnet.com/article/data-scientists-dont-scale/>

Time to insight

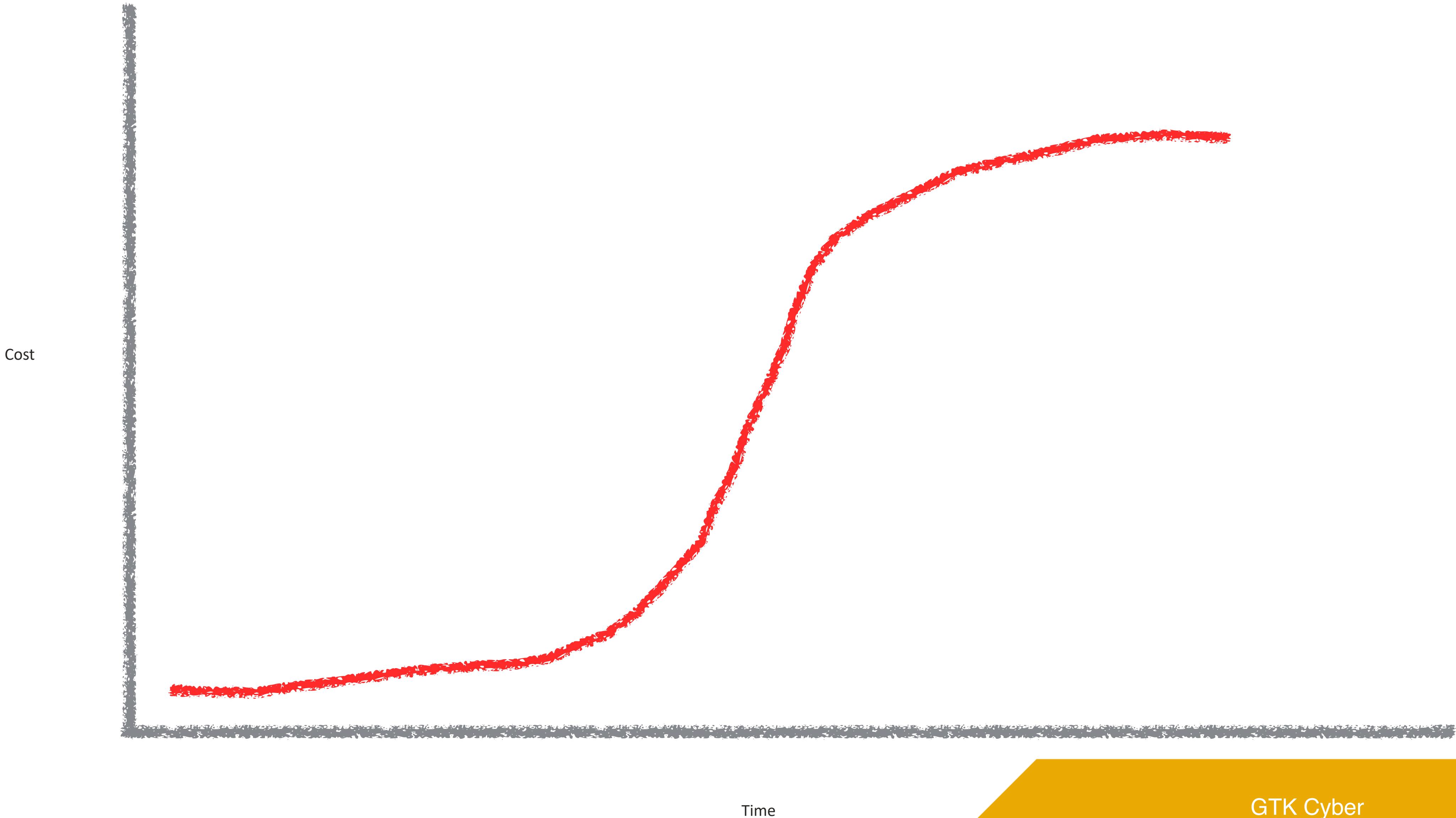
Time to Insight

Time = \$\$

Value of Insights Over Time

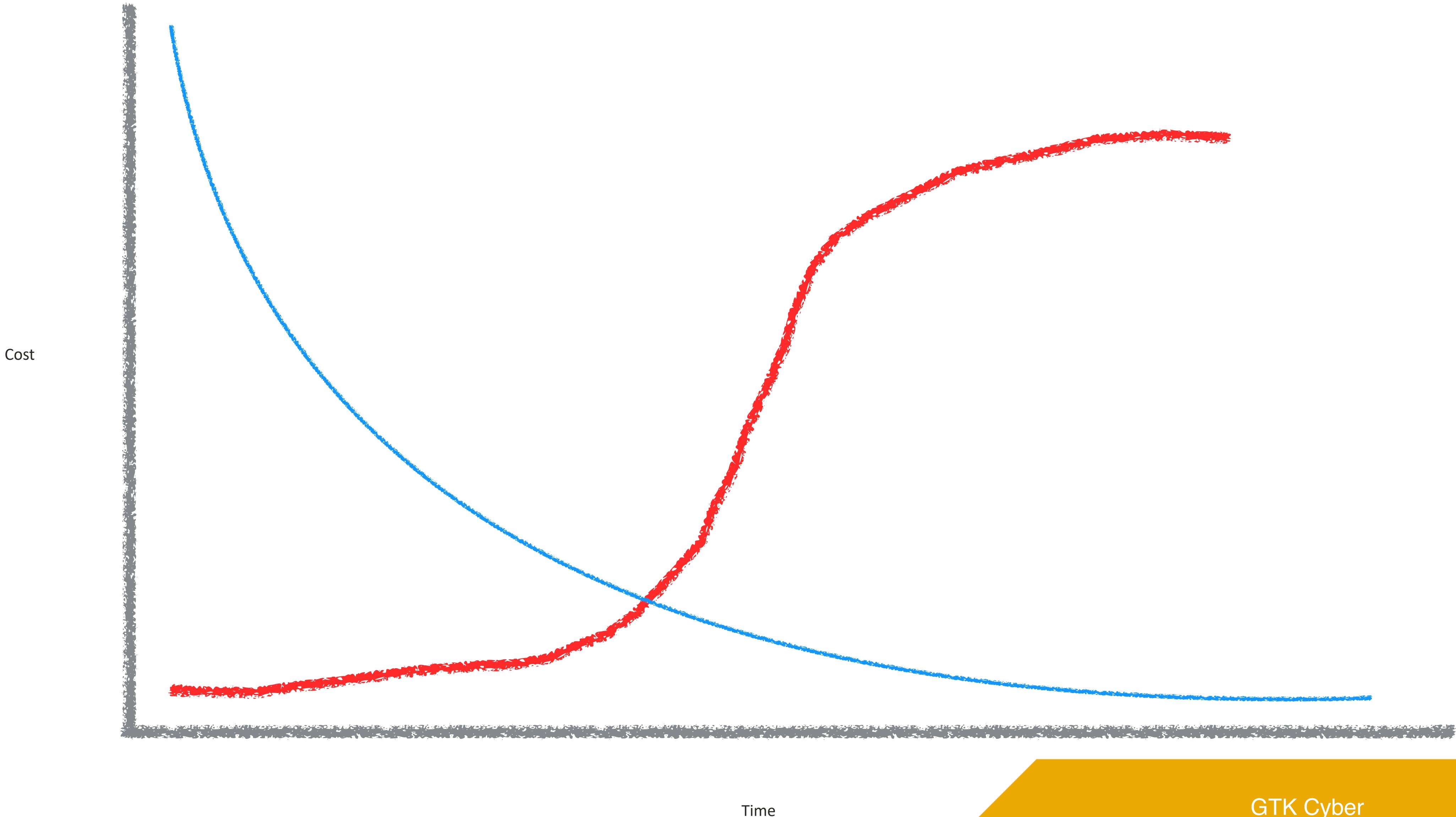


Costs of Insights Over Time

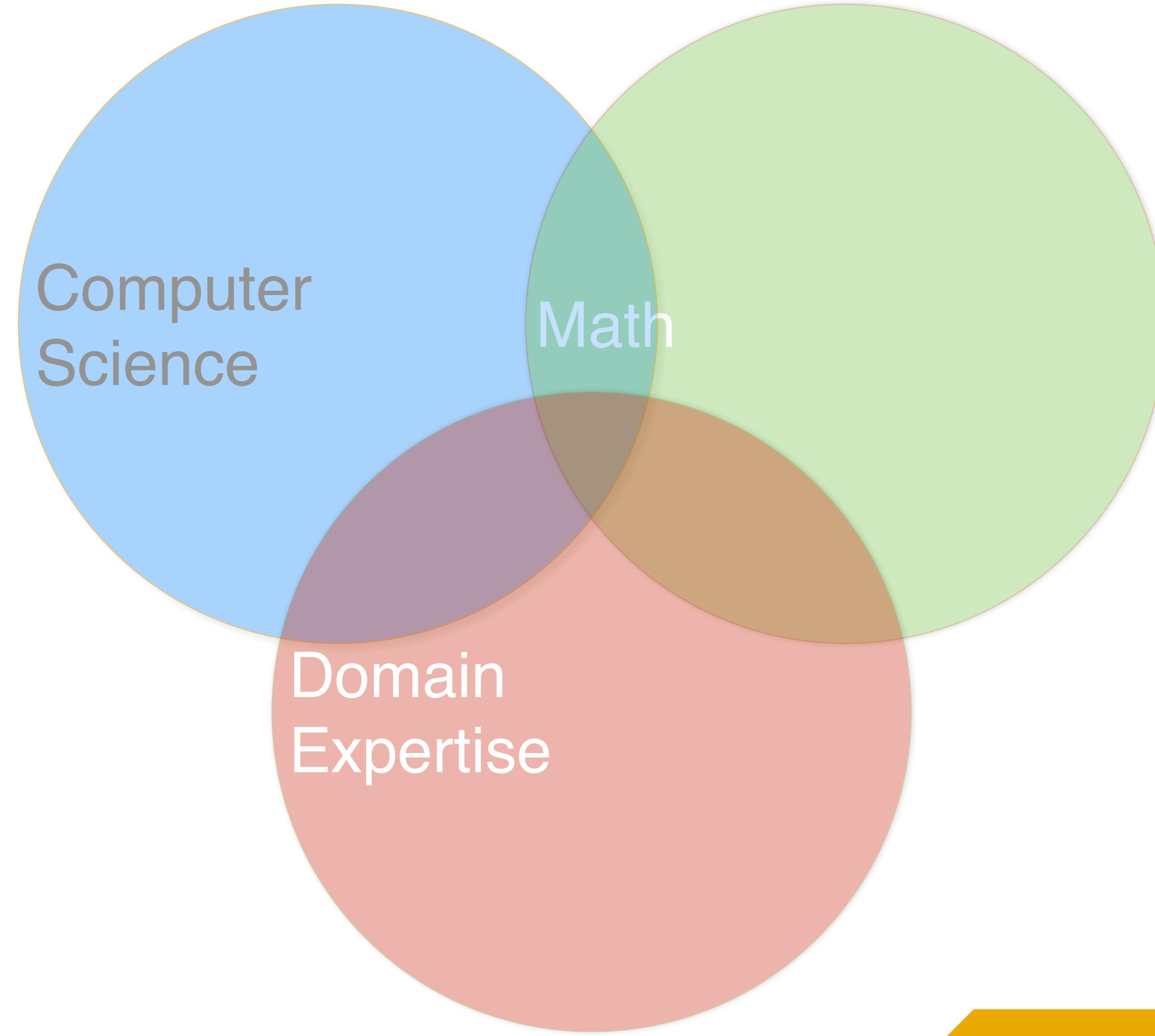


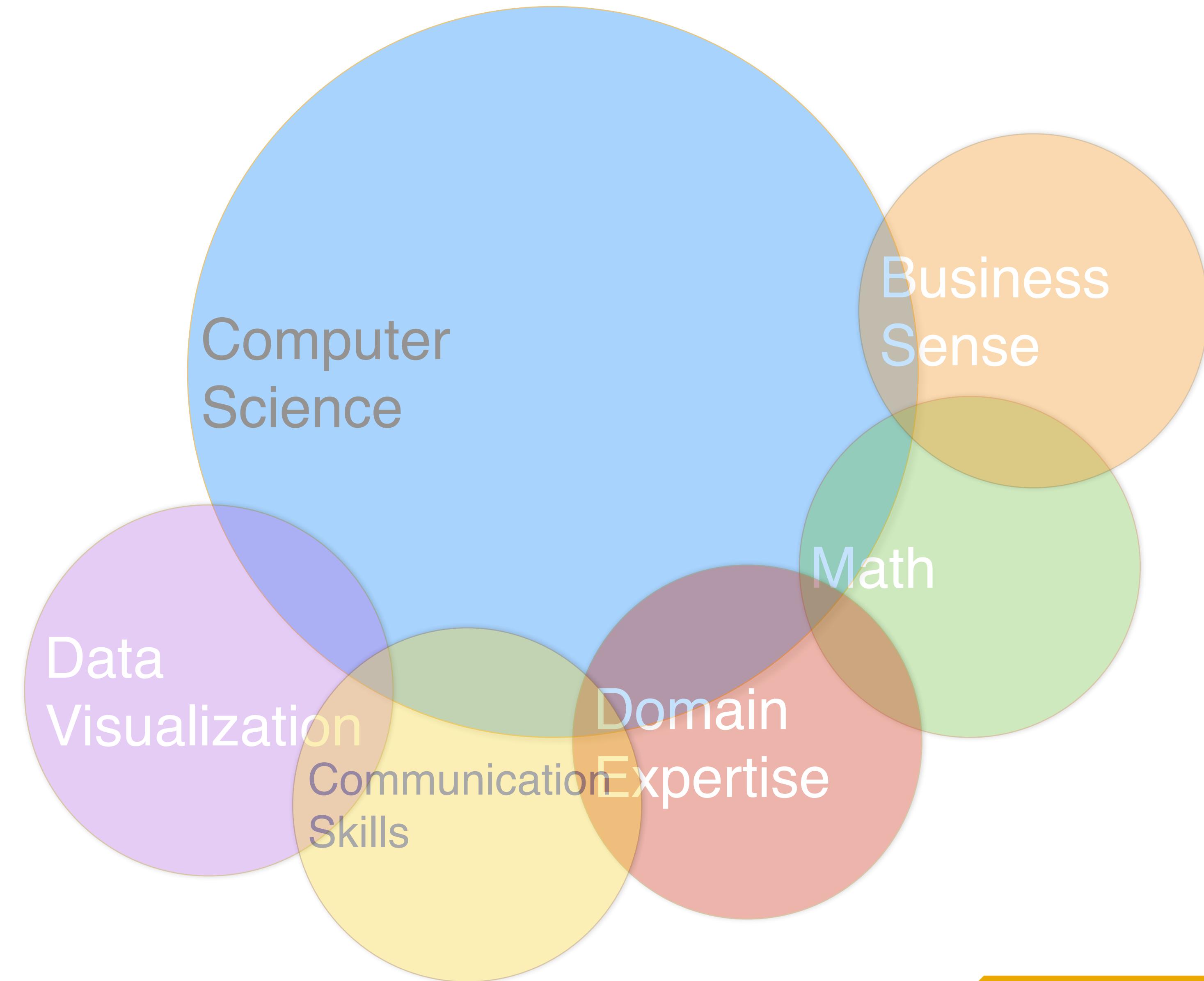
GTK Cyber

Cost of Insights Over Time

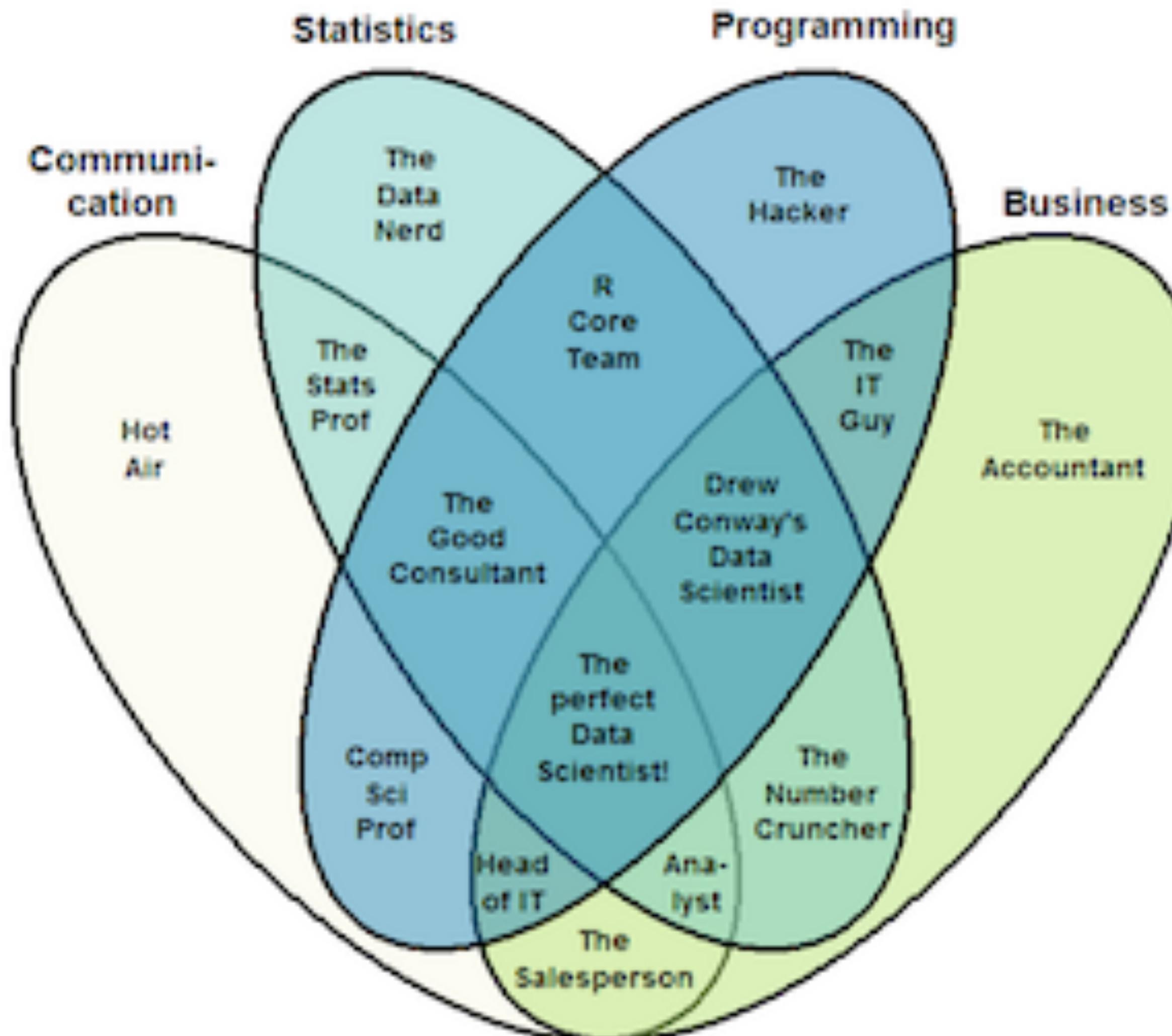


What Skills Does a Data Scientist Need?





The Data Scientist Venn Diagram

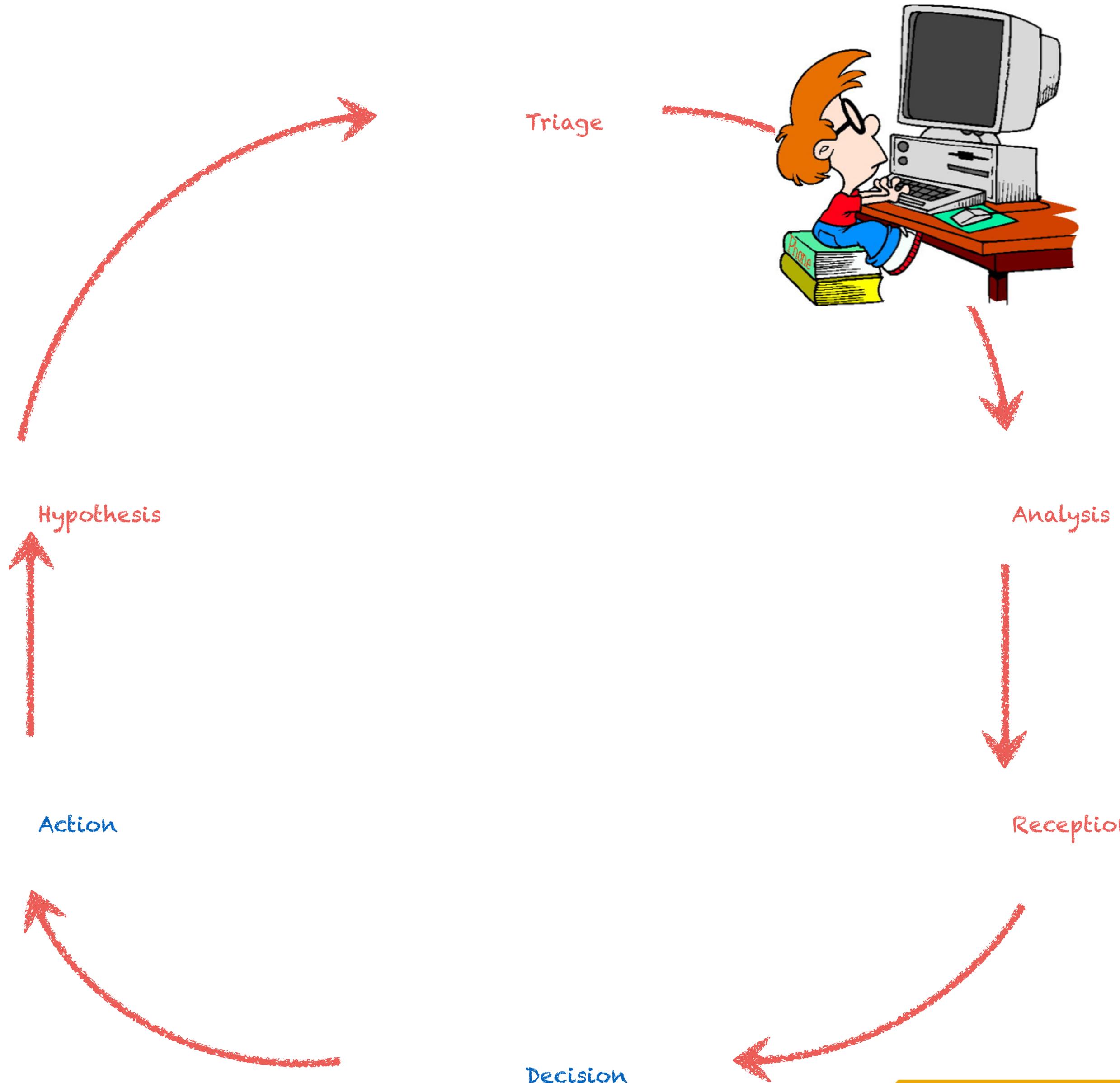


**Data Scientists spend
50-90% of their time
being...**

Data Janitors



GTK Cyber



GTK Cyber

Thoughts for Data Science Success

Data is a Strategic Asset... not a cost



Align Projects to Corporate Strategy

Align Projects to Corporate Strategy



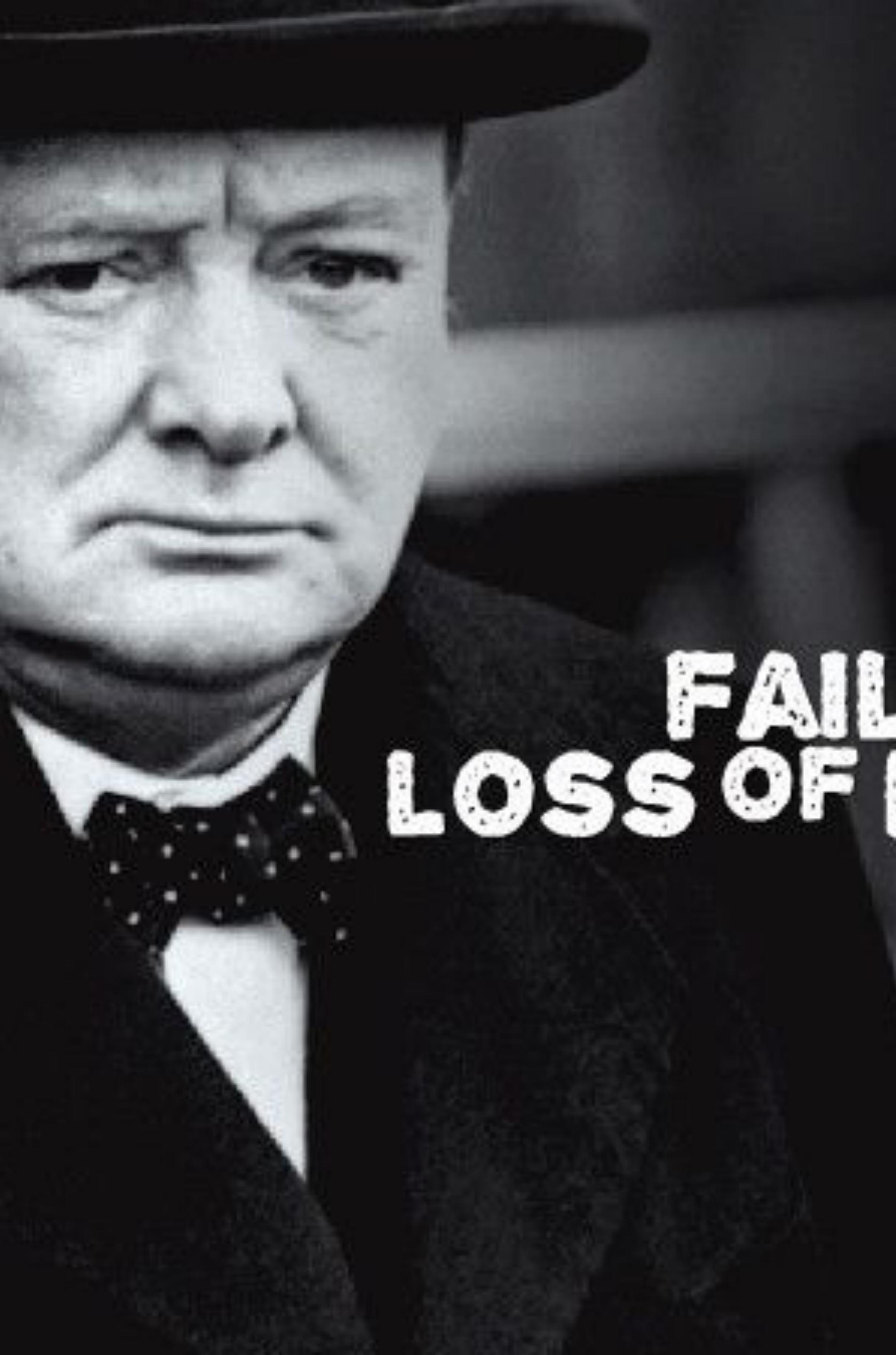
Your:
Time
Money
Job?

Build the right team for Data Initiatives



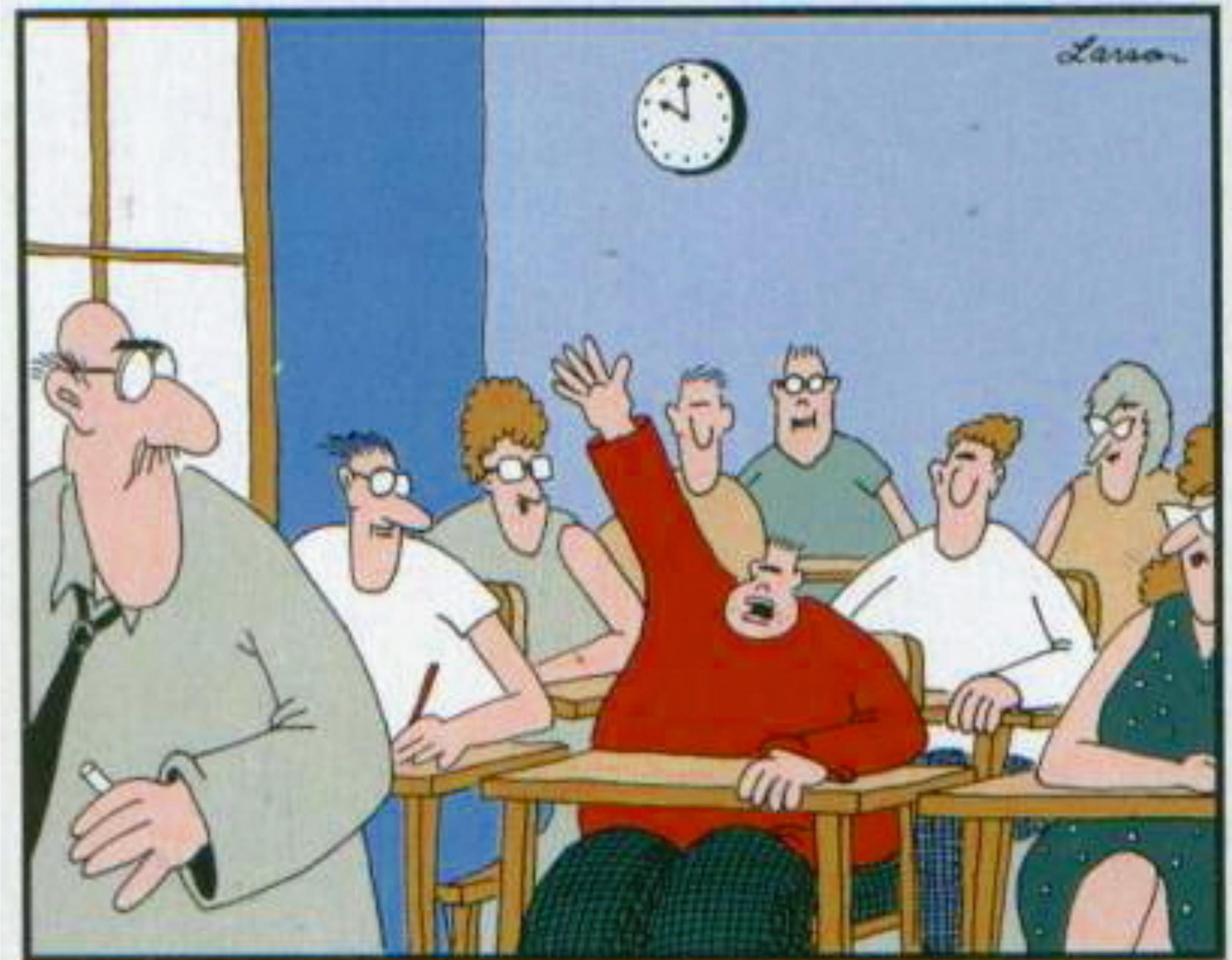
Prioritize building appropriate data platform





**"SUCCESS
CONSISTS OF
GOING FROM
FAILURE TO
FAILURE WITHOUT
LOSS OF ENTHUSIASM."**

Winston Churchill



**"Mr. Osborne, may I be excused?
My brain is full."**

Building a Data Science Team

Building a Data Science Team

Programmer + Statistician + SME = DS Team?

Building a Data Science Team

Programmer + Statistician + SME = DS Team?

Sort of...

Building a Data Science Team

- Data Ambassador
- Data Scientist
- Data Storyteller
- Data Engineer

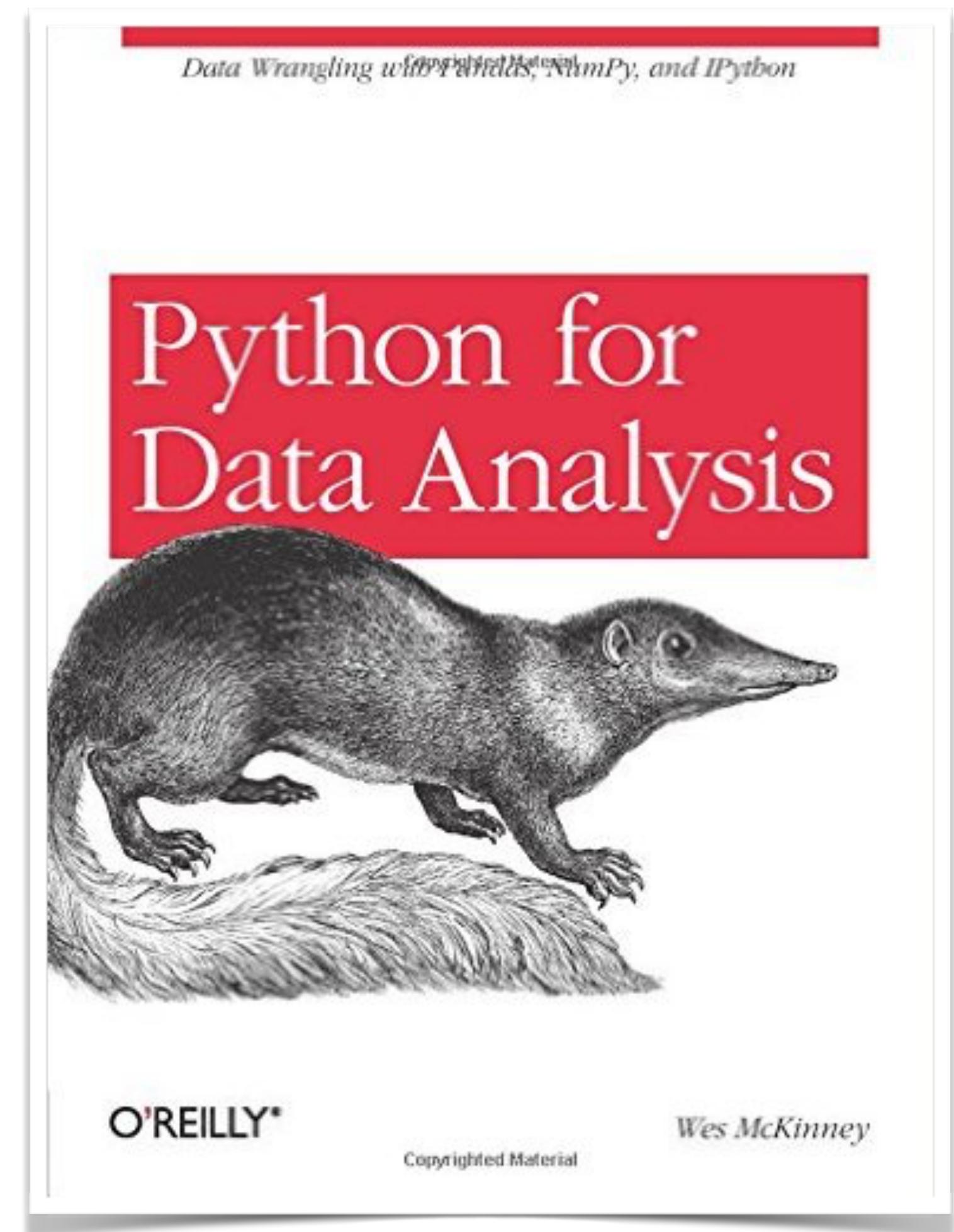
By The End of the Class, You Will Be Able To:

- Quickly and effectively prepare data for analysis
- Apply machine learning techniques to enhance security

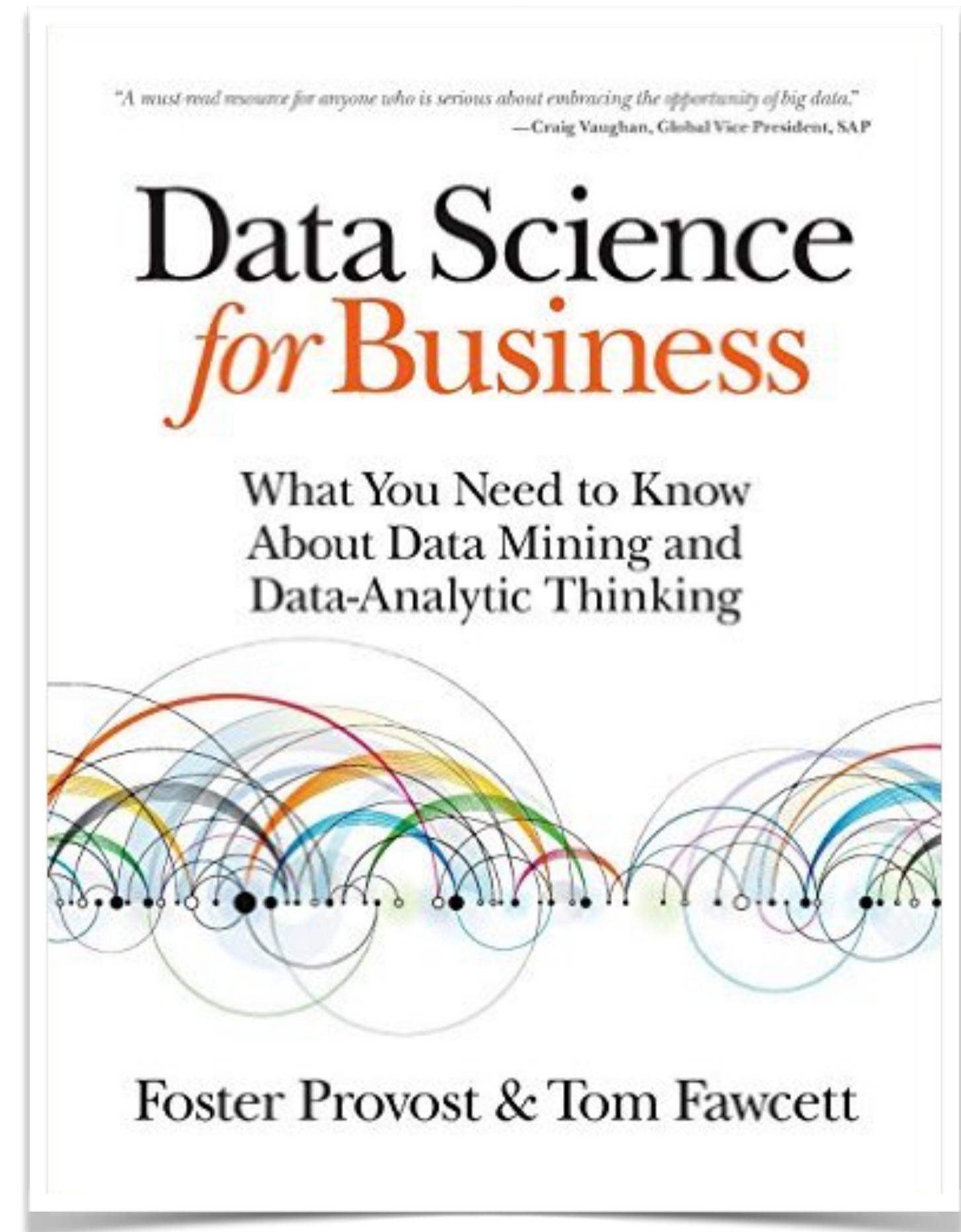




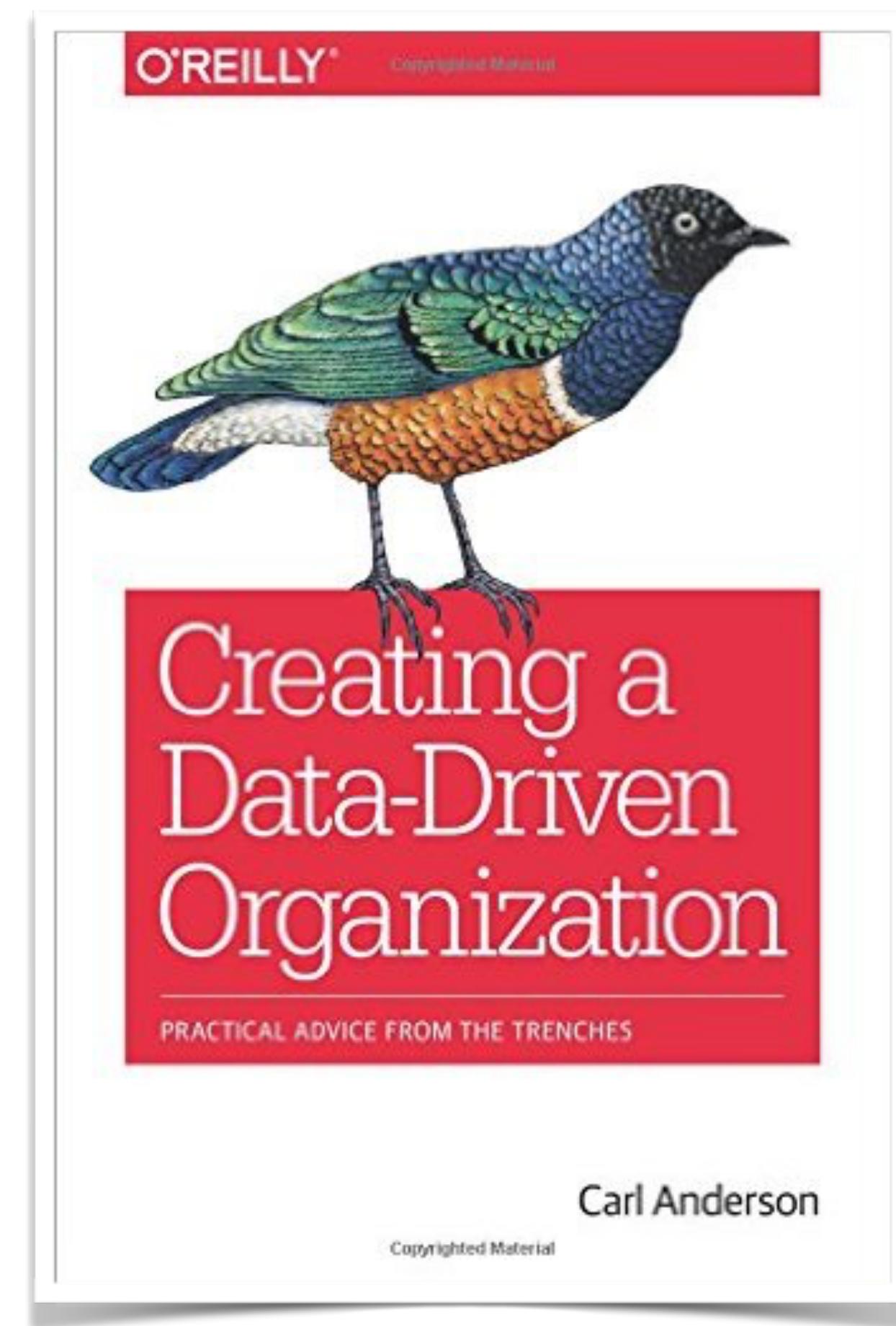
Recommended Reading



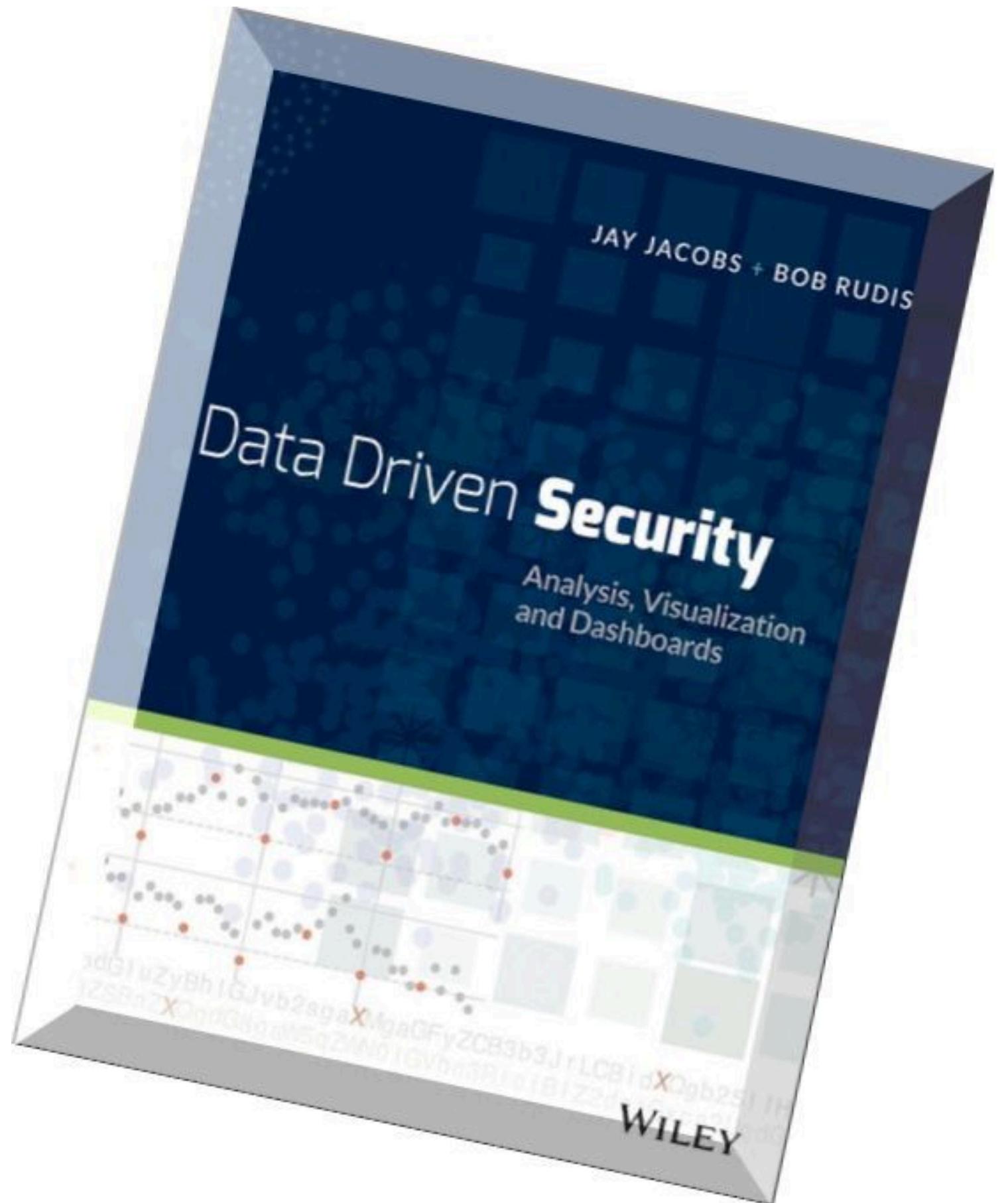
Recommended Reading



Recommended Reading



Recommended Reading



<http://datadrivensecurity.info>

Machine Learning

"Type a Quote Here"

-- *Bill Gates*

What is Machine Learning (ML) Artificial Intelligence (AI)

“Machine Learning is the science of getting computers to act without being explicitly programmed.”

– <https://www.coursera.org/course/ml>

“A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.”

–*Tom Mitchell, Carnegie Mellon University*

“Machine learning explores the construction and study of algorithms that can learn from and **make predictions on data**. Such algorithms operate by building a model from example inputs in order to make data-driven predictions or decisions, **rather than following strictly static program instructions.**”

–https://en.wikipedia.org/wiki/Machine_learning



GTK Cyber





- Blacklists
- Simple keyword matching
- Naive Bayesian Classifiers
- Deep Learning

Artificial Intelligence

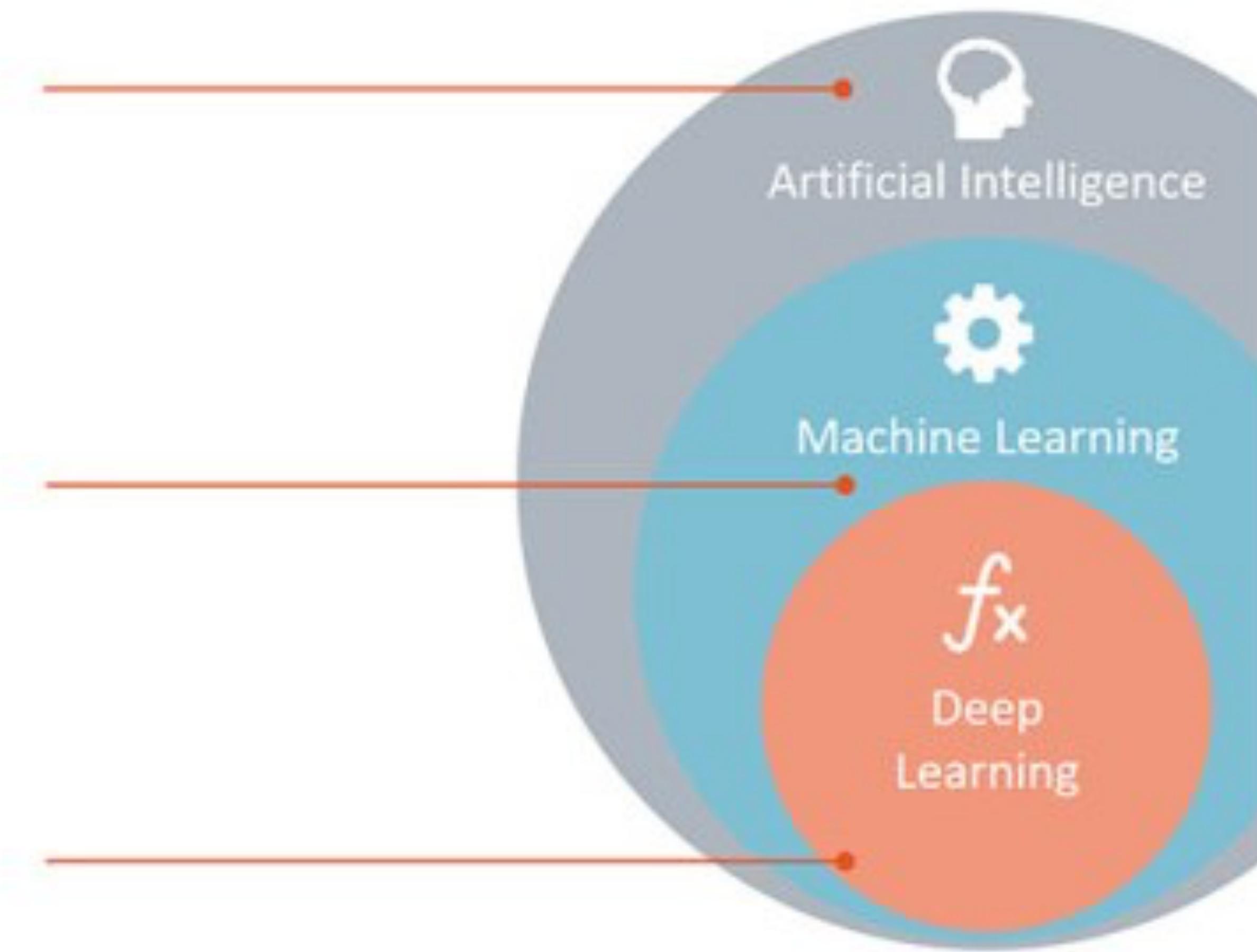
Any technique which enables computers to mimic human behavior.

Machine Learning

Subset of AI techniques which use statistical methods to enable machines to improve with experiences.

Deep Learning

Subset of ML which make the computation of multi-layer neural networks feasible.



@katherinebailey Because marketing? Every time someone calls simple linear regression “AI” Gauss turns over in his grave.

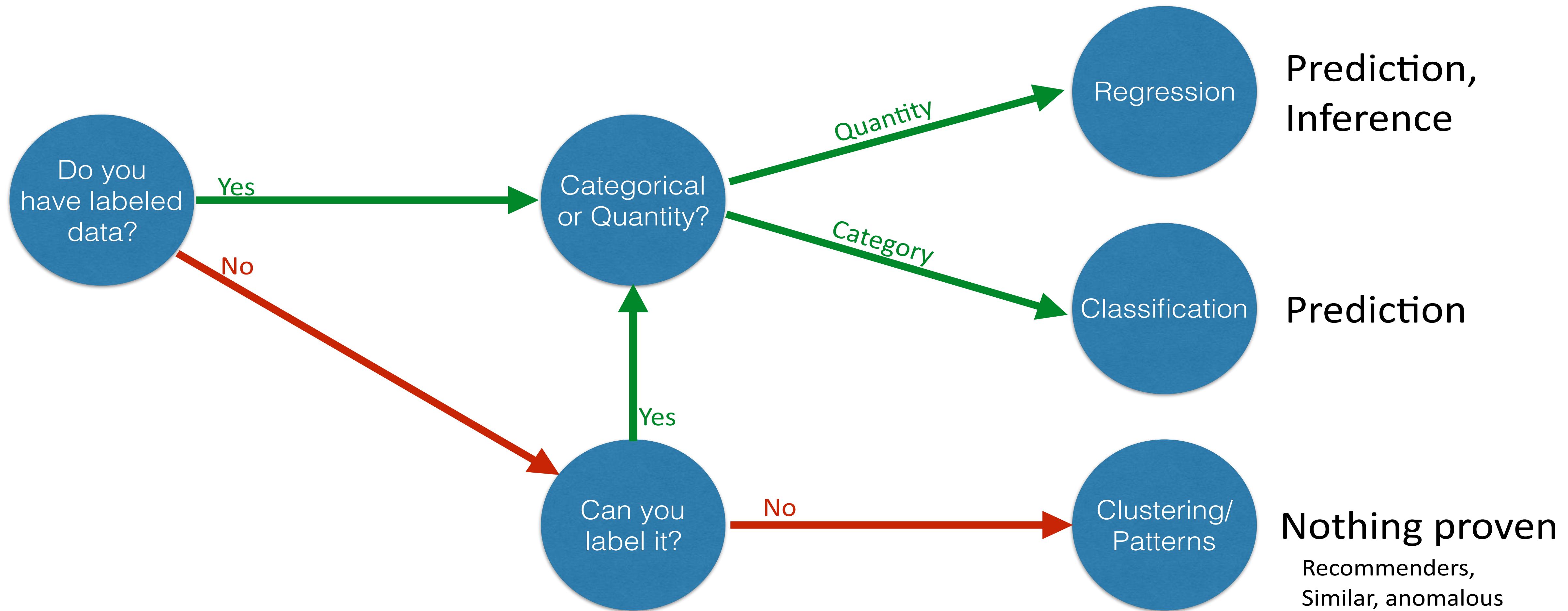
Machine Learning Problems

- **Supervised Learning:** Supervised Learning is a class of Machine Learning in which a model is "trained" using a set of pre-existing labeled data.
- **Unsupervised Learning:** A class of Machine Learning algorithms in which a model is built without the use of labeled data.

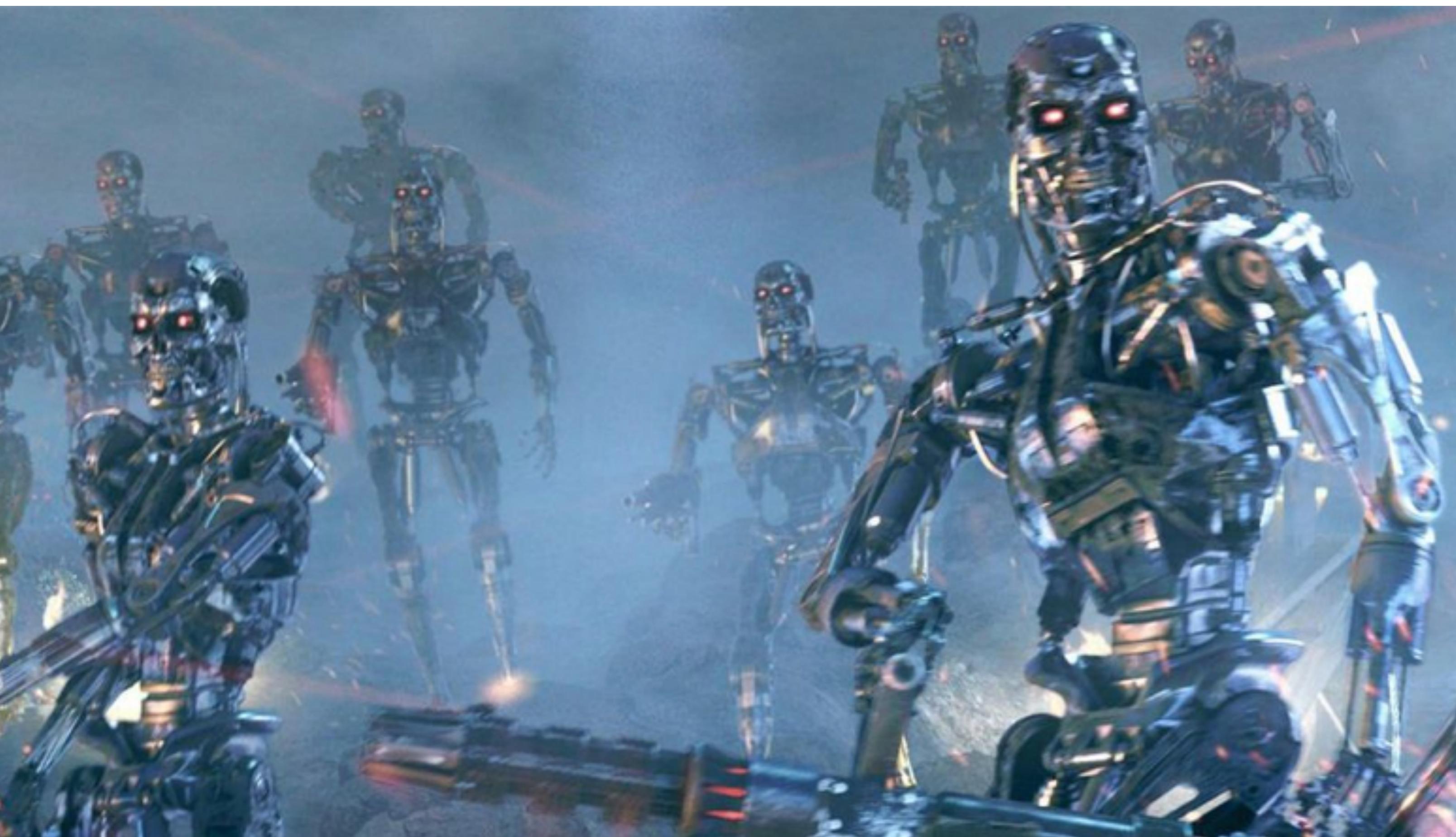
Machine Learning Problem Types

- **Classification:** Assigning or predicting a observation's membership in discrete class
- **Regression:** Predicting a continuous value based on the observations' features
- **Clustering:** Identifying groupings within a dataset
- **Dimensionality Reduction:** Reducing the number of variables in a feature set

What Problem am I solving?



what it is Not



Applications to Security

Regression Example

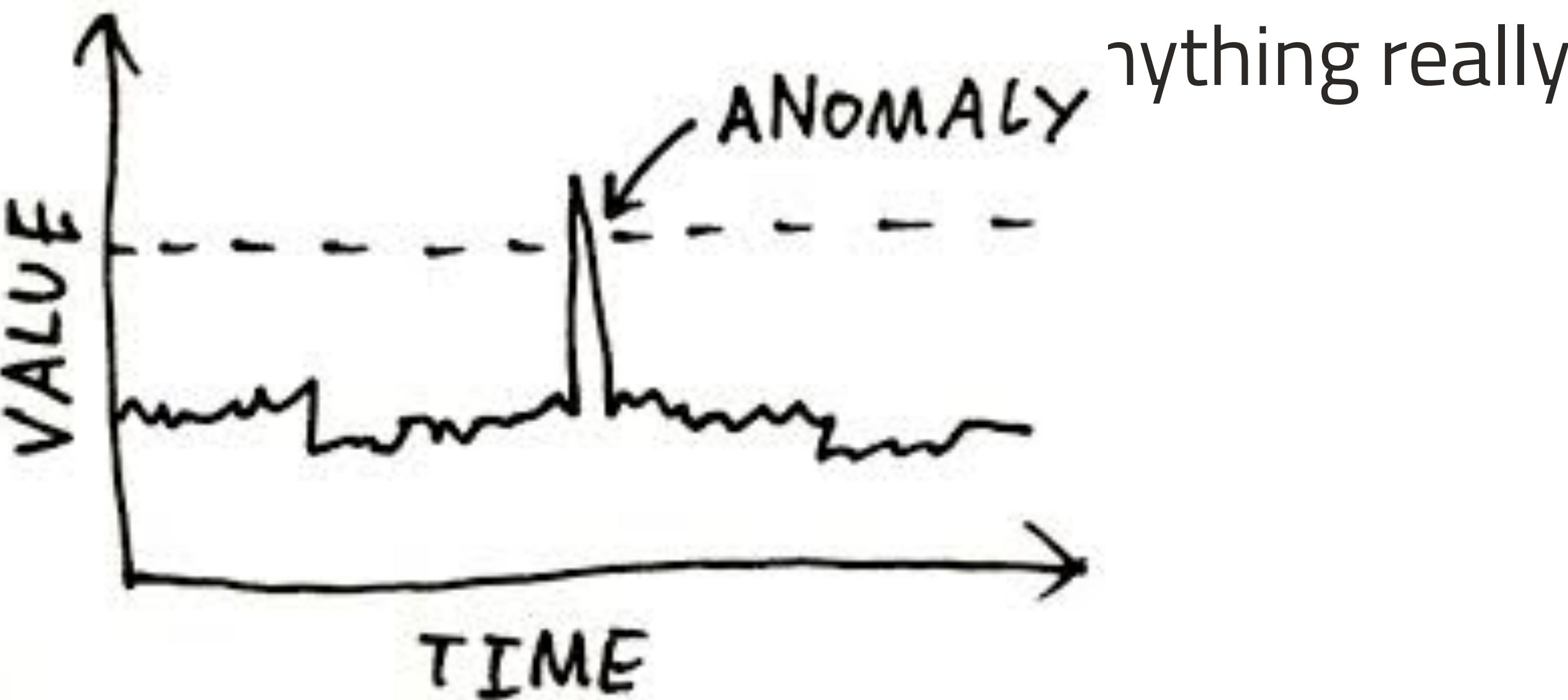
Server Capacity Prediction: Regression analysis can be used to predict a server's capacity (or CPU usage) based on the server's historical performance.



https://www.researchgate.net/publication/256645877_LiRCUP_Linear_Regression_based_CPU_Usage_Prediction_Algorithm_for_Live_Migration_of_Virtual_Machines_in_Data_Centers

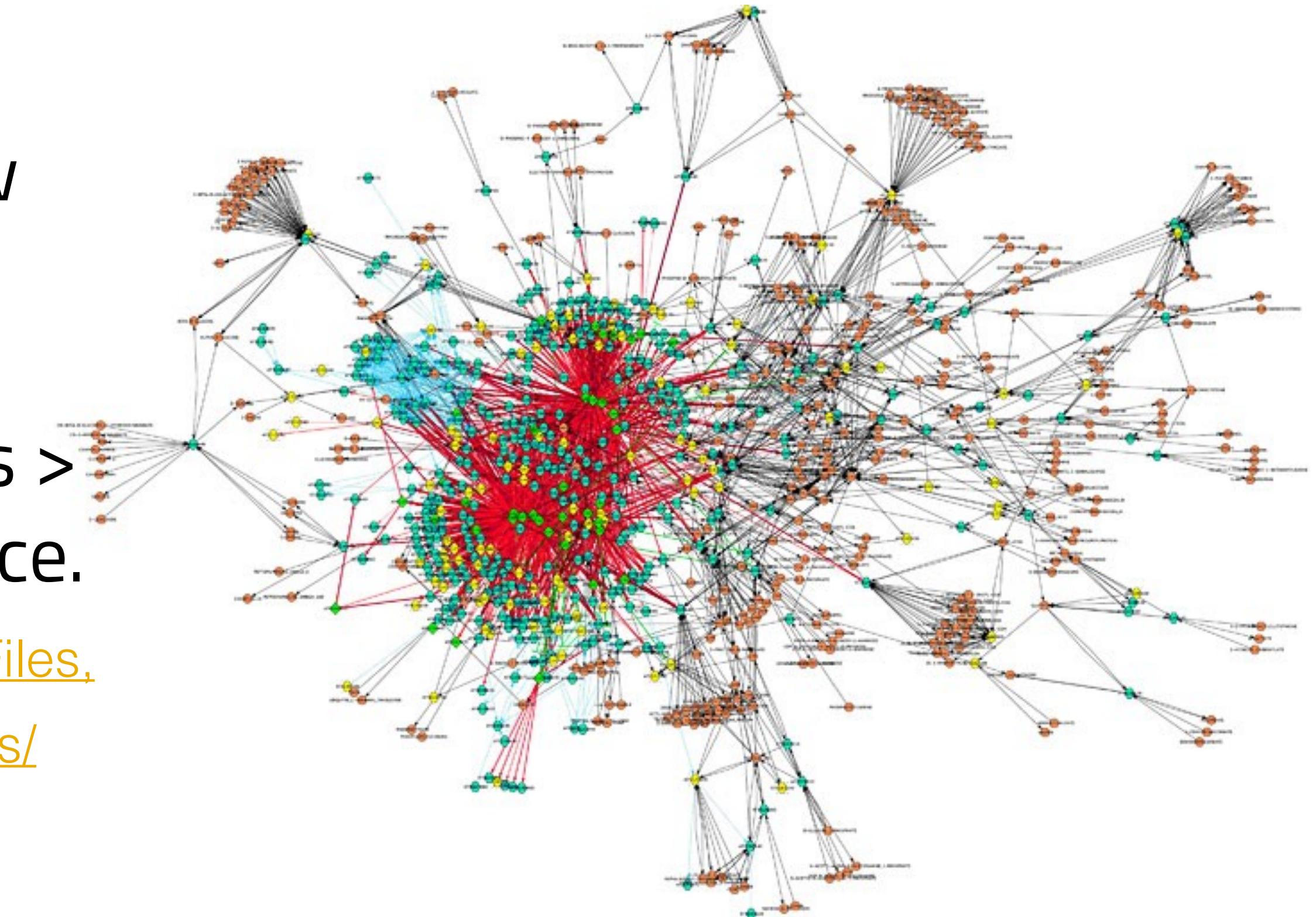
Clustering Example

Anomaly Detection: Clustering techniques can be used to detect



Network-Based Intrusion Detection

- Derive Features from Network Traffic
Captures “pcap” at packet level or NetFlow level (tools: tshark, tcpdump, bro...)
- Example Features based on header information: 2s-windowing of connections > duration, protocol, src and dat bytes, service.
- Get data sets: <http://www.netresec.com/?page=PcapFiles>,
<https://maccdc.org/>, <http://www.westpoint.edu/crc/SitePages/DataSets.aspx> <http://www.unb.ca/cic/research/datasets/>,



Malware Detection/Classification

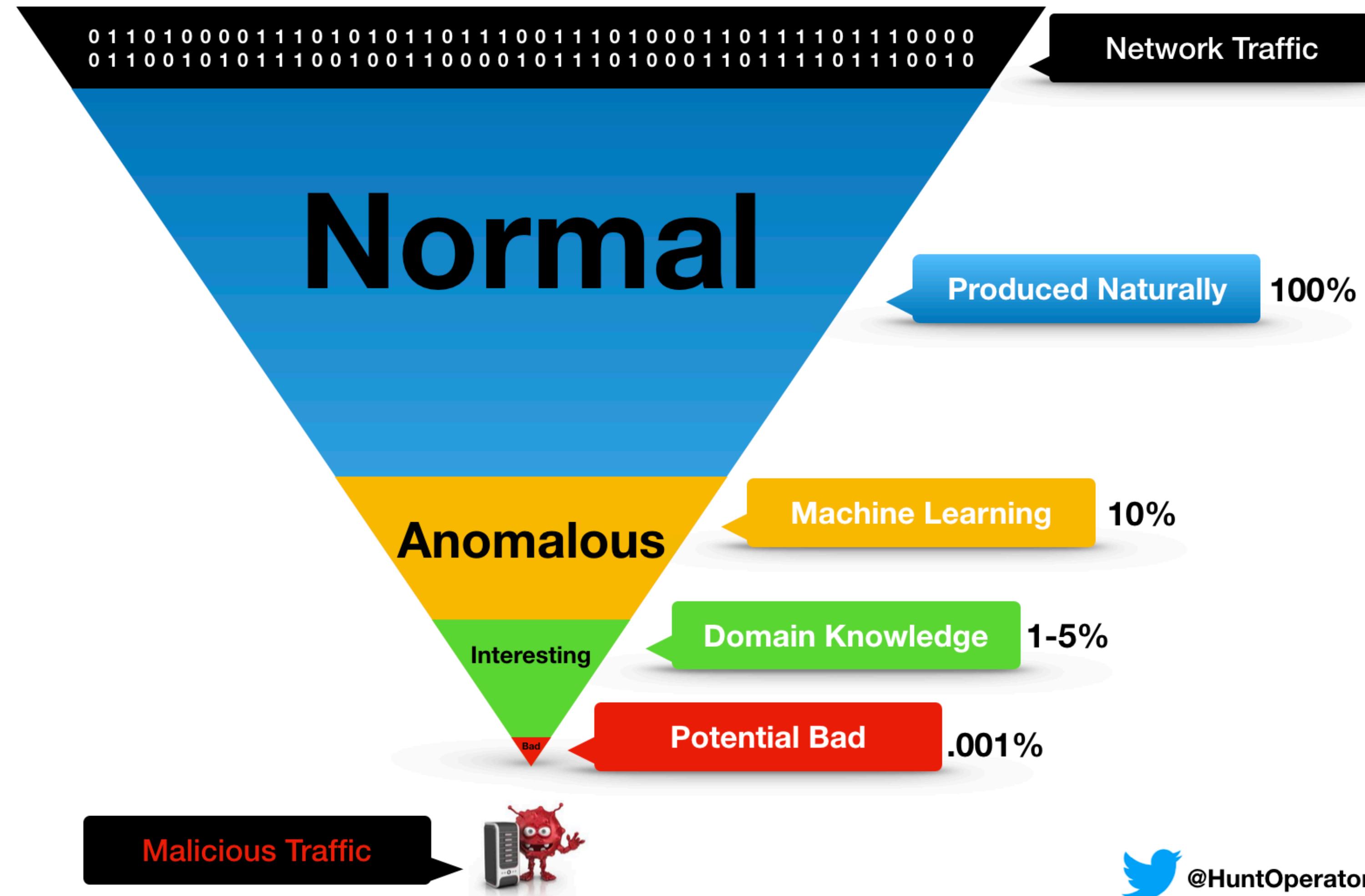
- Derive Features from Binary Content and metadata manifest (function calls, string obtained from IDA Disassembler)
- Example Features: opcode count (n-grams), segment count, asm pixel intensity, n-gramming of bytes, function name.
- Featureless Deep Learning with word2vec embedding
- Get open source malware samples: Vx Heaven, Virus Share, Maltrieve, Open Malware



Security Applications of Machine Learning

- Domain Generation Algorithm (DGA) Detection (Classification)
- Malicious URL Detection (Classification)
- Network Traffic: Beaconing Detection (Classification/Clustering)
- Detection of new classes of malware (Classification/Clustering)
- General Network Traffic Anomaly Detection (Classification/Clustering)
- Log Analysis - Anomaly Detection (Classification/Clustering)
- Phishing Detection (Classification)
- Identifying SQL Injection (Classification)
- Identifying XSS cross-site scripting (Classification)
- DOS/DDOS Detection (Classification)
- Authentication (Classification)

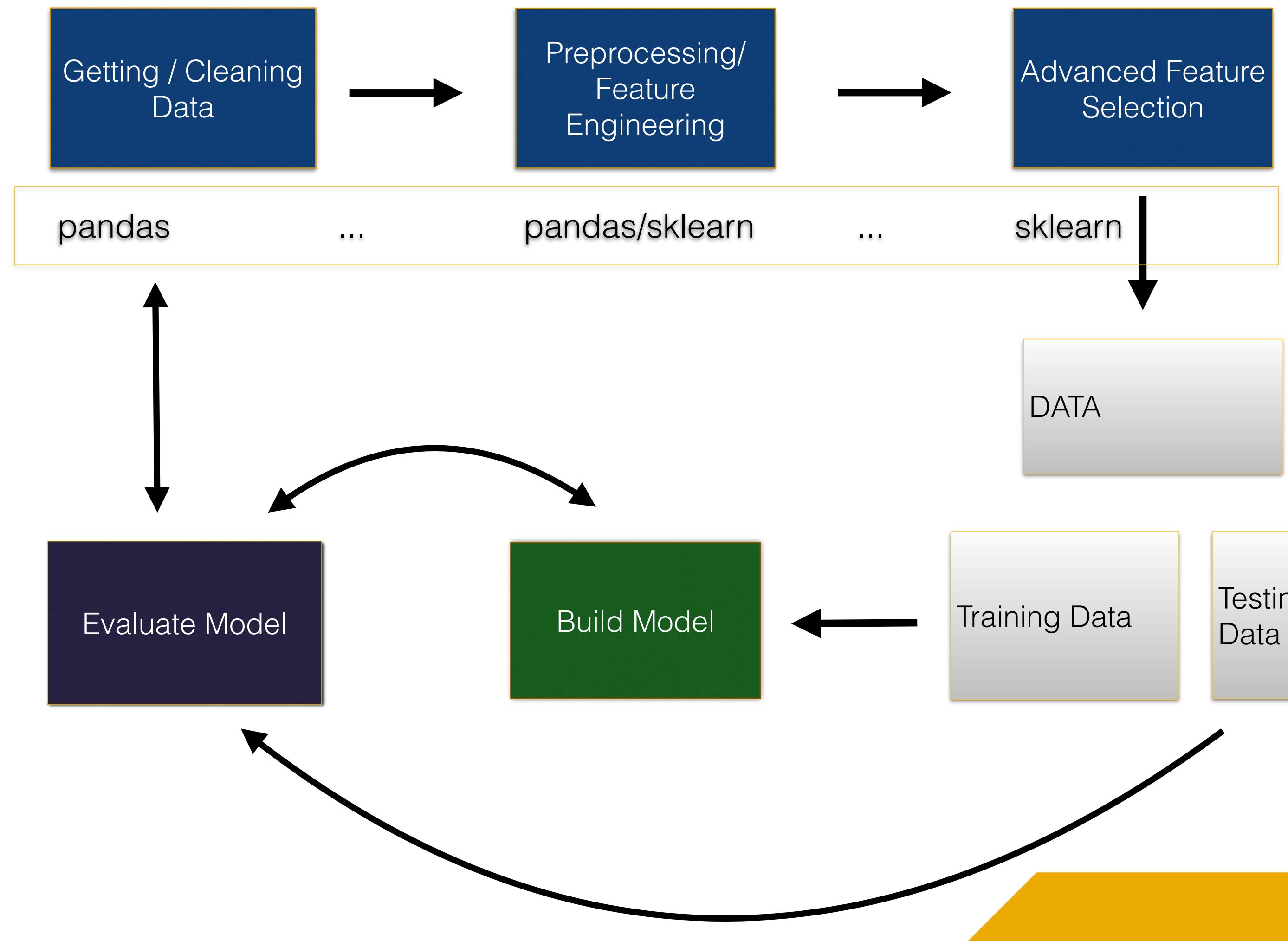
Data Science Hunting Funnel



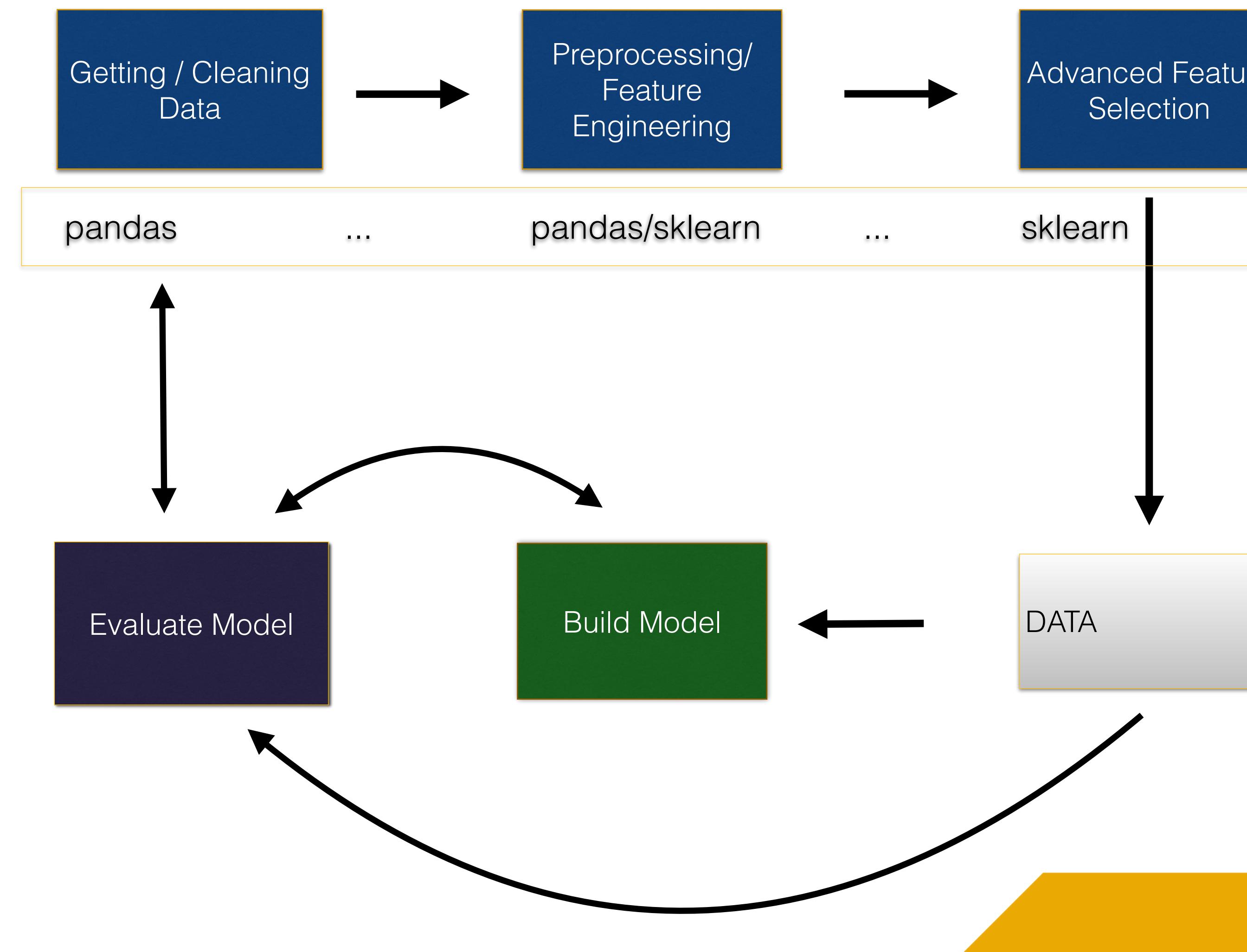
<http://www.austintaylor.io/network/traffic/threat/data/science/hunting/funnel/machine/learning/domain>

The Machine Learning Process

Supervised Machine Learning Process



Unsupervised Machine Learning Process



**First, define your analytic
question.**

what are you trying to do?

**How do you define success?
What are you measuring?**

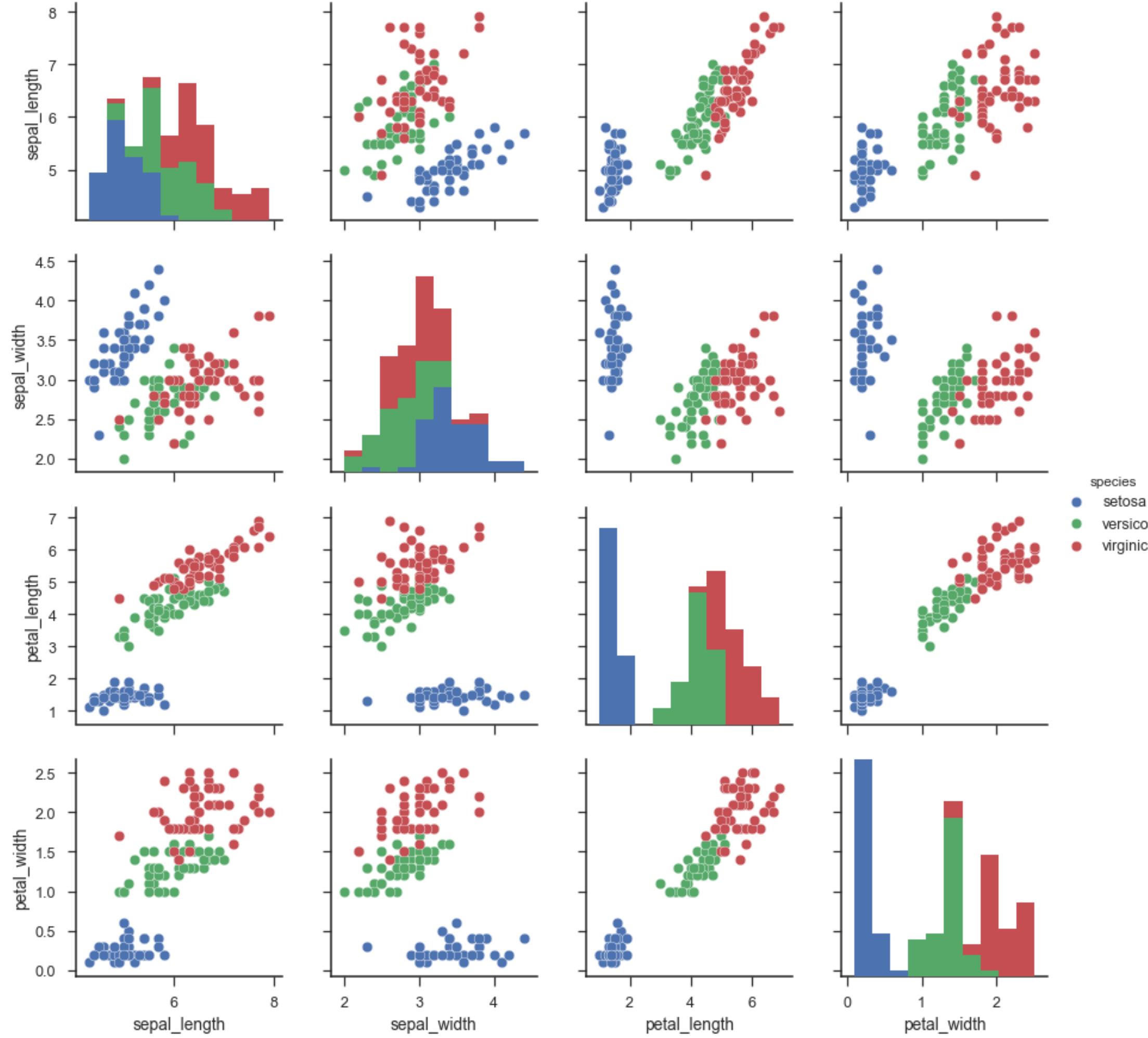
Choose data sources

- What is available?
- Is it enough?
- Is the data reliable/clean/consistent?
- What other data could you use?

Other Considerations

- Policies
- Legal constraints
- Biases in Data
- Latency
- Data size

Gather and Explore Your Data



Is the data good enough?

What are the rules governing its use?

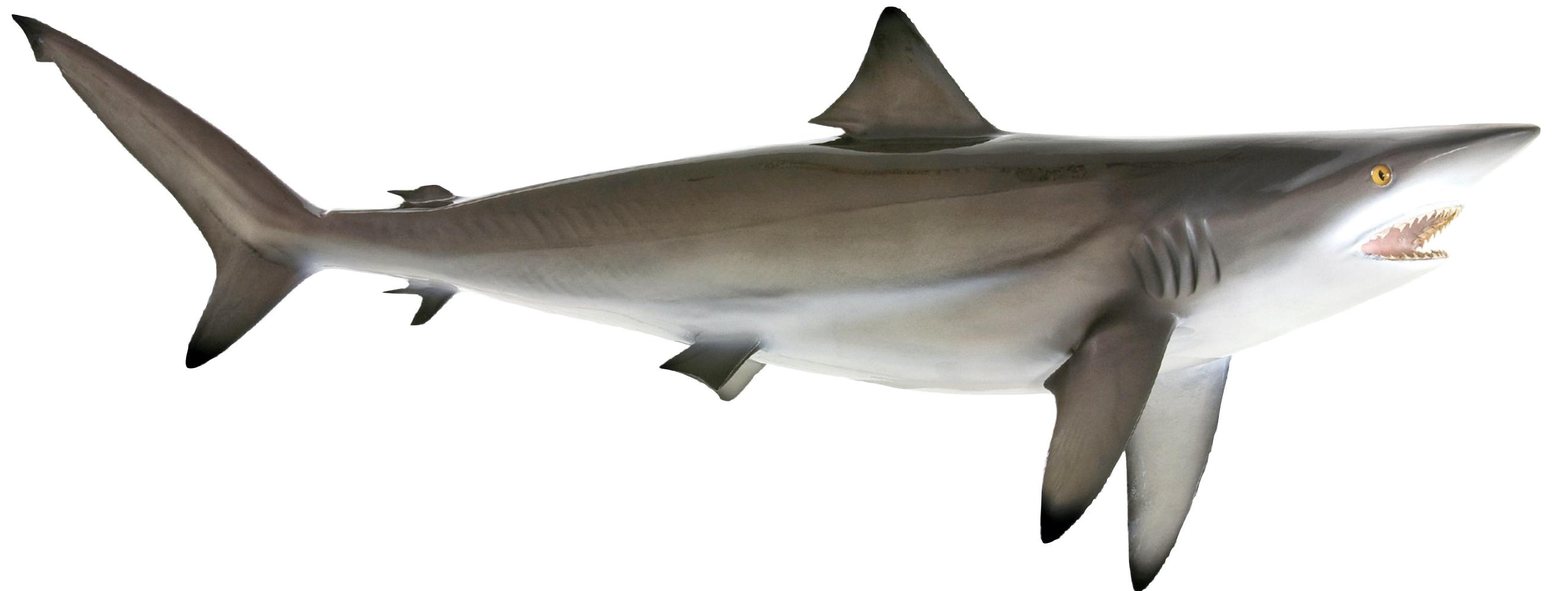
Do I have enough?

Do problems or biases exist in the data
that could cause problems?

Feature Engineering

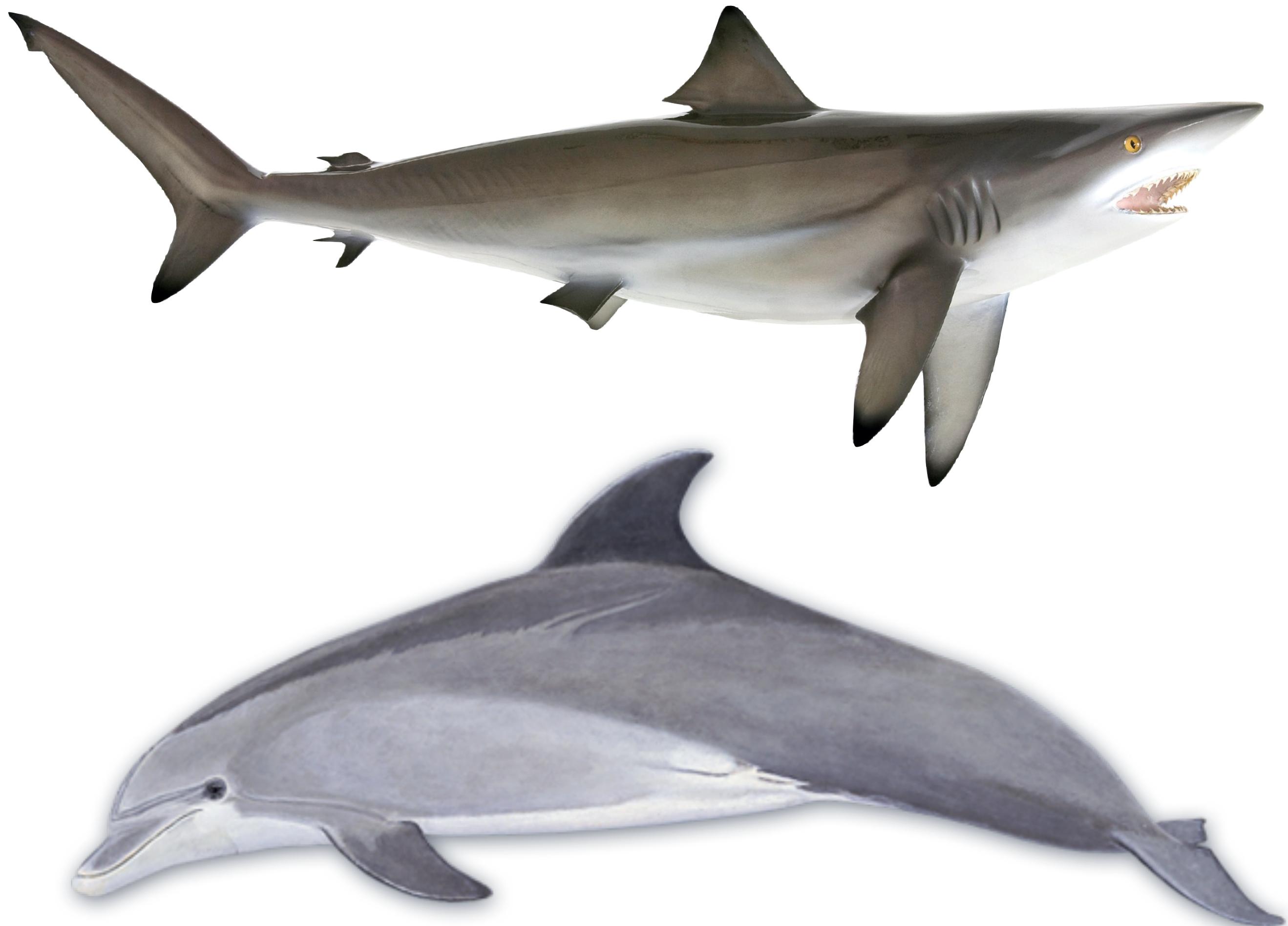
- Define what you are trying to measure. These will become the observations or rows of your final dataset
- Define how you will mathematically represent your data. This will be come the features or columns of your final dataset.

Feature Engineering



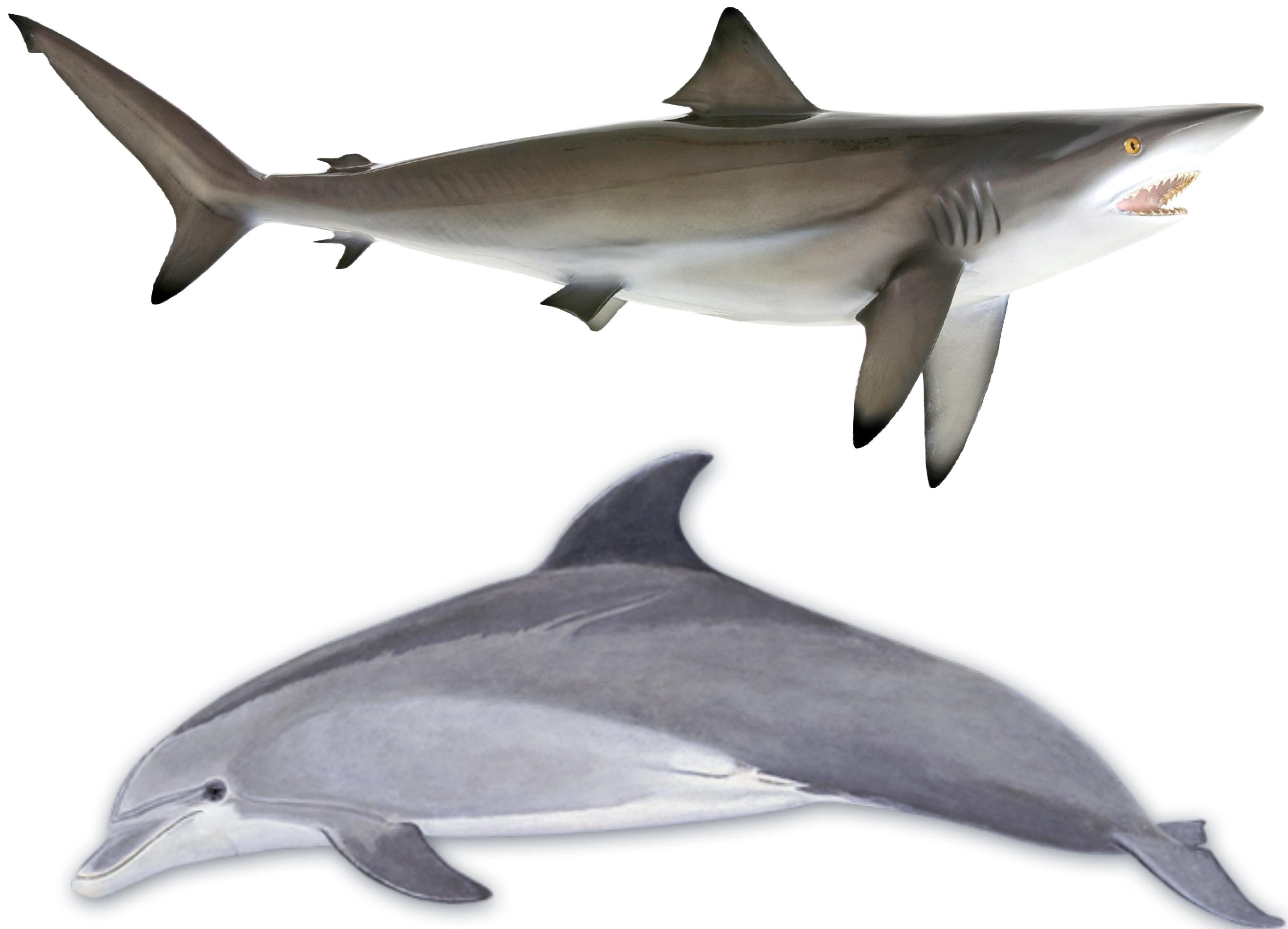
Feature	Value
Color	Gray
Fins	7
Predator	TRUE

Feature Engineering



Feature	Value
Color	Gray
Fins	7
Predator	TRUE

Feature Engineering



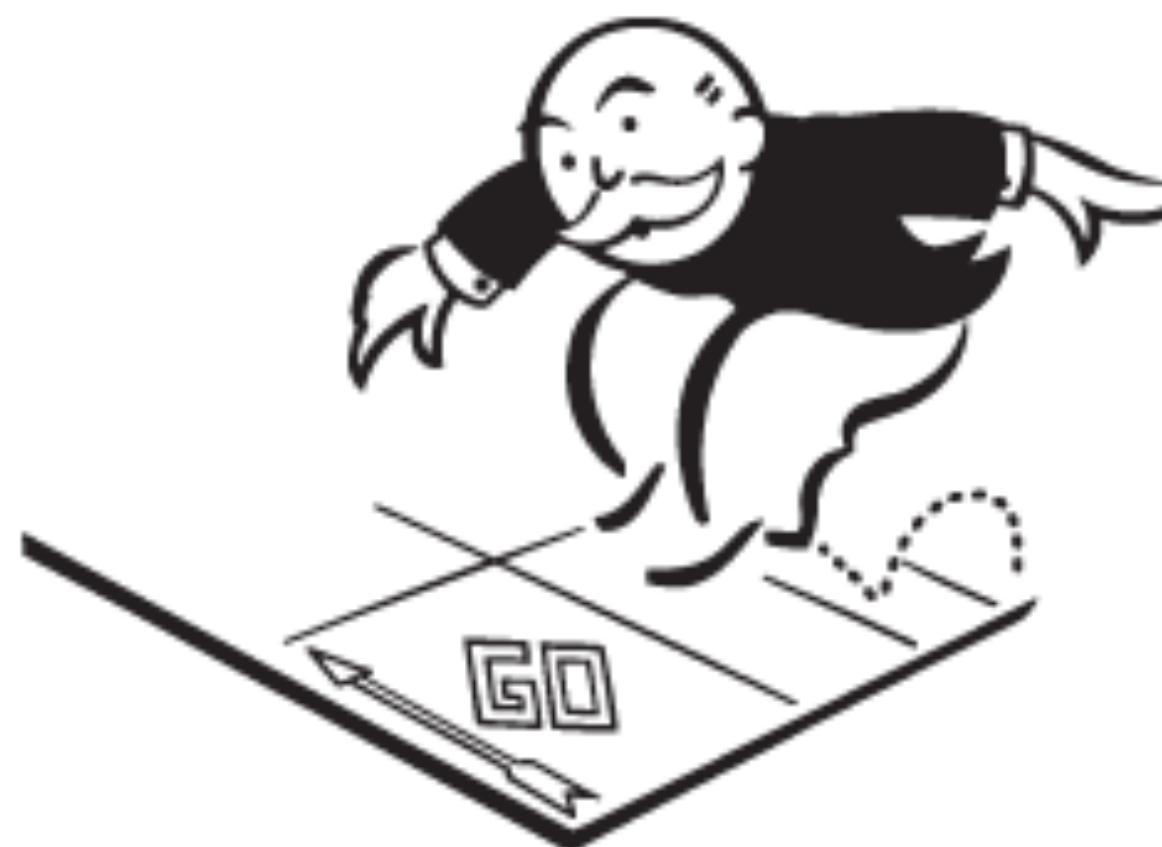
Feature	Value
Color	Gray
Fins	7
Predator	TRUE
Mammal	TRUE

Build and Tune your Model

- Believe it or not, this is the easy part.
- Most of this is **done using libraries** like scikit-learn, mllib, tensorflow, caret or keras, and **many steps can be automated**.
- You can even do it in Splunk or Elasticsearch.

Evaluate Performance

- Use various scoring methods, or write your own to determine model performance.
- Go back to step 1 and repeat! (Do not pass go, do not collect \$200)



Group Discussion

Consider that you are building a system to identify fraudulent credit card transactions. In your groups, try to answer the following questions:

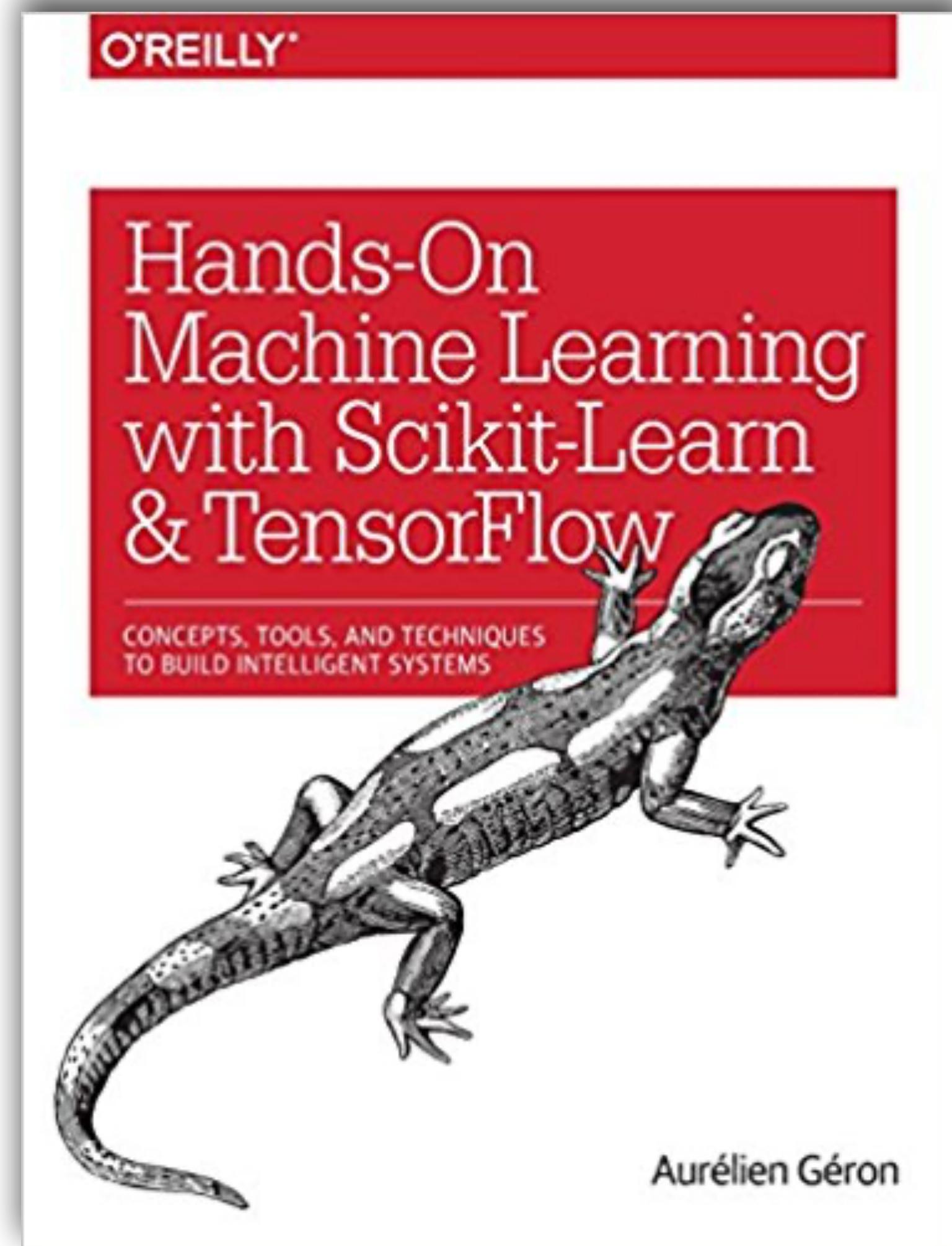
1. What are some features that you would want to capture?
2. What data sets will you need?
3. What legal and policy challenges might you face?
4. What other challenges you could foresee in this problem?
5. How will you define success?
6. How can you articulate the value of this model to stakeholders?

The Python Data Science Ecosystem

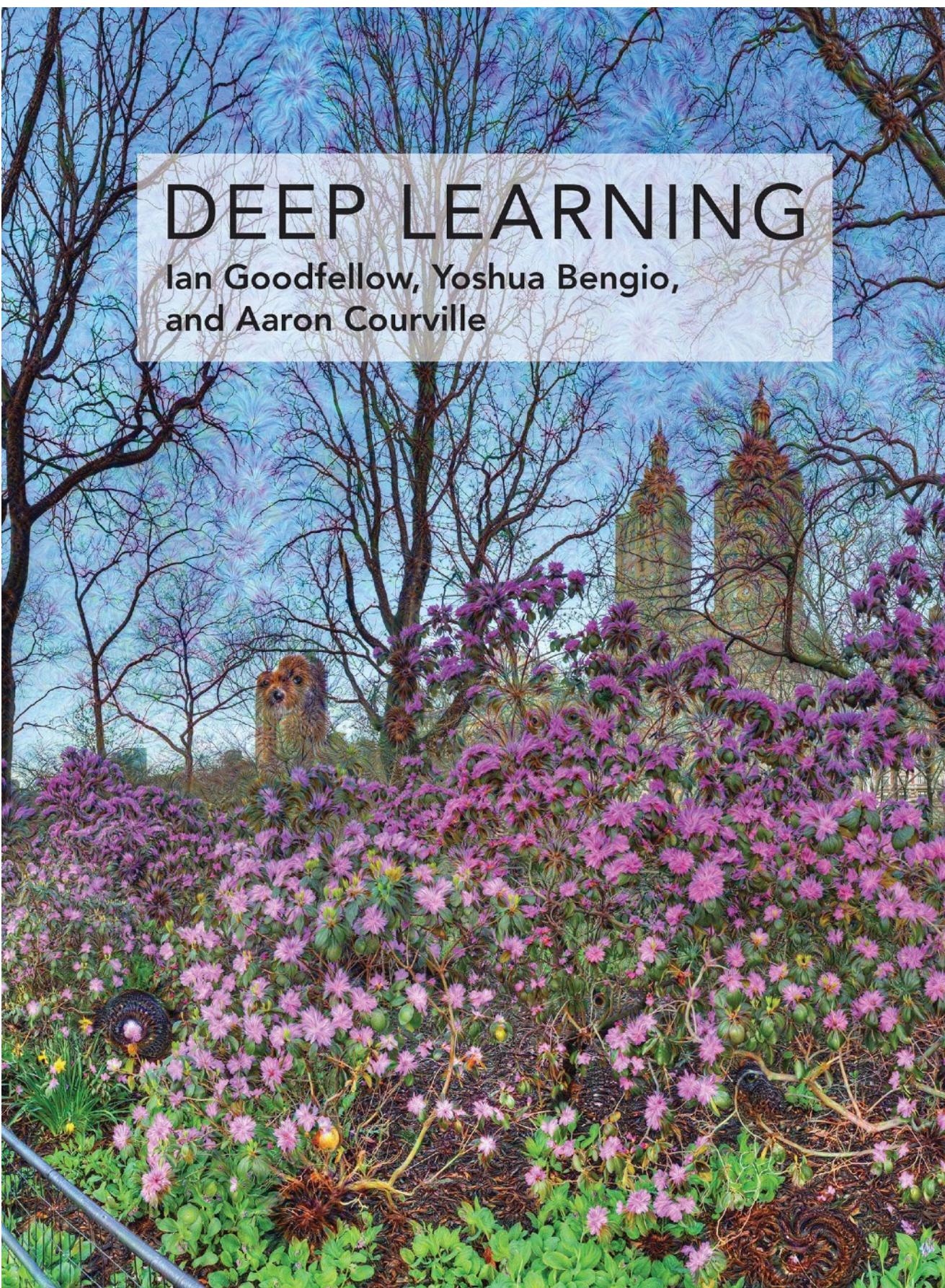
Machine Learning Ecosystem

- **Data Gathering:** Pandas, Drill, BeautifulSoup, PyDBAPI, PyDAL, Boto3
- **Feature Extraction:** Pandas, NumPy, Featuretools
- **Machine Learning**
 - **"Regular" ML:** Scikit-learn (sklearn), h2o, mllib (PySpark)
 - **Deep Learning:** Tensorflow, Keras, Theano, Caffe, PyTorch
- **Visualization:** Matplotlib, Seaborn, Yellowbrick, LIME, ggplot, plot.ly,

Recommended Reading



Recommended Reading



<http://www.deeplearningbook.org/>

O'REILLY®

Machine Learning & Security

PROTECTING SYSTEMS WITH DATA AND ALGORITHMS



Clarence Chio & David Freeman

GTK Cyber

O'REILLY®



Feature Engineering for Machine Learning

PRINCIPLES AND TECHNIQUES FOR DATA SCIENTISTS

Alice Zheng & Amanda Casari

GTK Cyber

O'REILLY



Learning Apache Drill

QUERY AND ANALYZE STRUCTURED DATA

Charles Givre & Paul Rogers

GTK Cyber

The Virtual Machine: Griffon

```
File Edit View Search Terminal Help
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/hbase/hbase-1.1.3/lib/slf4j-log4j12-1
.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop/share/hadoop/common/lib/slf4
j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
2016-05-16 13:04:54,887 WARN  [main] util.NativeCodeLoader: Unable to load nativ
e-hadoop library for your platform... using builtin-jar classes where applicabl
e
2016-05-16 13:05:11,827 ERROR [main] zookeeper.RecoverableZooKeeper: ZooKeeper exists faile
d after 4 attempts
2016-05-16 13:05:11,828 WARN  [main] zookeeper.ZKUtil: hconnection-0x46a145ba0x0, quorum=lo
calhost:2181, baseZNode=/hbase Unable to set watcher on znode (/hbase/hbaseid)
org.apache.zookeeper.KeeperException$ConnectionLossException: KeeperErrorCode = ConnectionL
oss for /hbase/hbaseid
        at org.apache.zookeeper.KeeperException.create(KeeperException.java:99)
        at org.apache.zookeeper.KeeperException.create(KeeperException.java:51)
        at org.apache.zookeeper.ZooKeeper.exists(ZooKeeper.java:1045)
        at org.apache.hadoop.hbase.zookeeper.RecoverableZooKeeper.exists(RecoverableZooKeep
er.java:221)
        at org.apache.hadoop.hbase.zookeeper.ZKUtil.checkExists(ZKUtil.java:541)
        at org.apache.hadoop.hbase.zookeeper.ZKClusterId.readClusterIdZNode(ZKClusterId.jav
a:65)
        at org.apache.hadoop.hbase.client.ZooKeeperRegistry.getClusterId(ZooKeeperRegistry.
java:105)
```

Do Data Science, Not Sysadmin

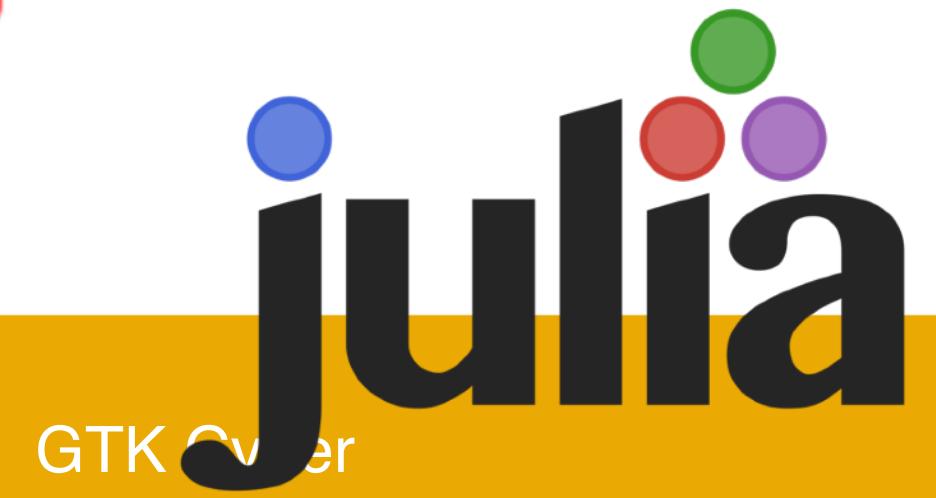
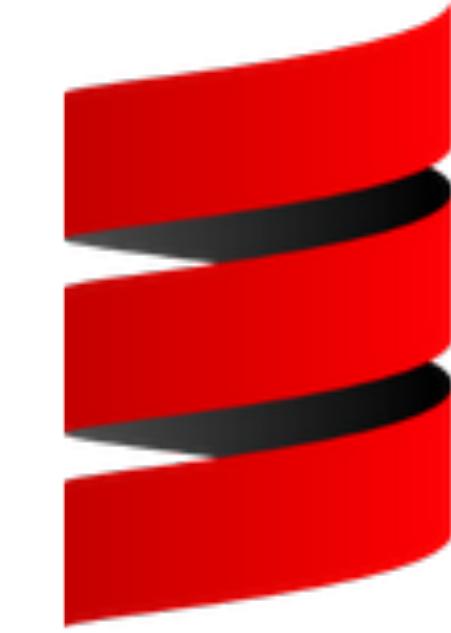
Built on Ubuntu MATE

Easy to Use

Programming Languages

- Languages
- Libraries
- Editors and Notebooks
- Databases + Administrative tools
- Big Data Tools
- Machine Learning Libraries
- Data Visualization

Scripting Languages



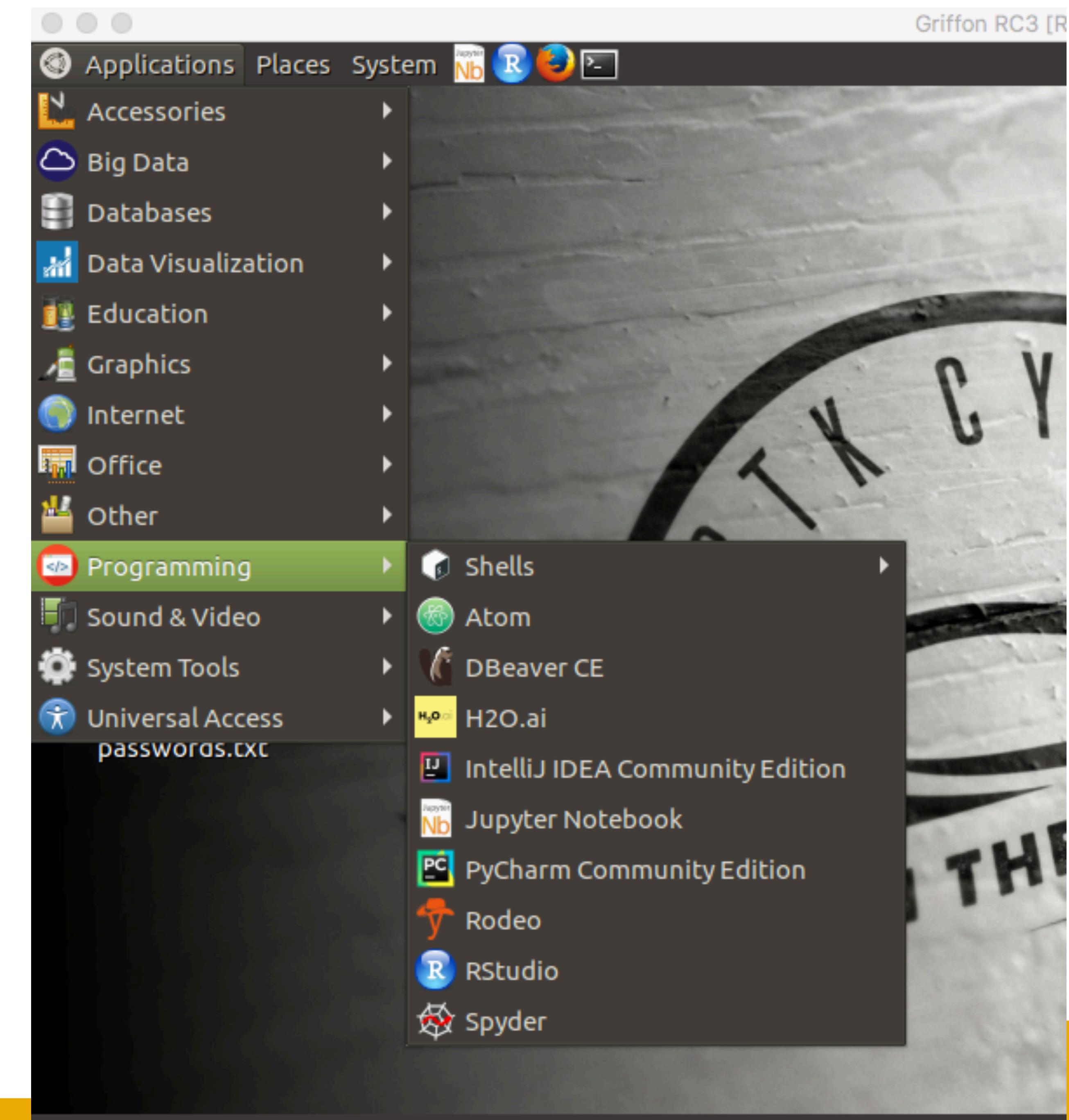
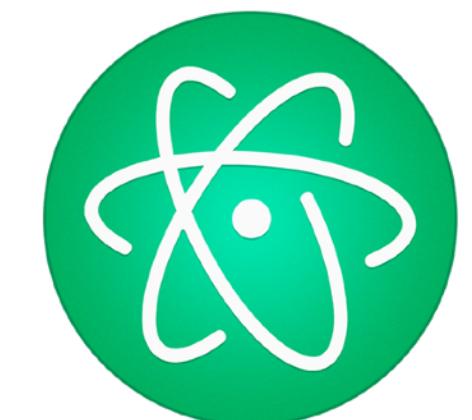
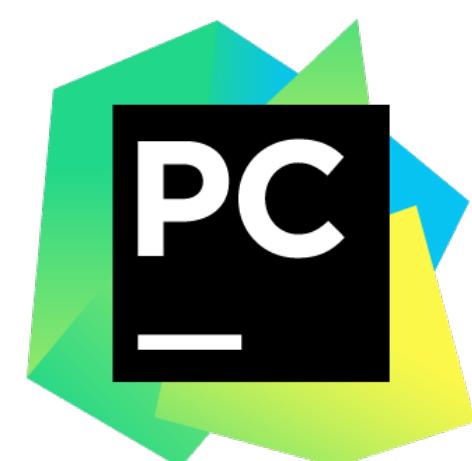
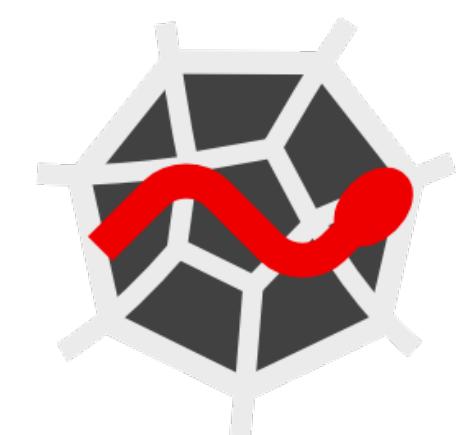
GTK Overlay

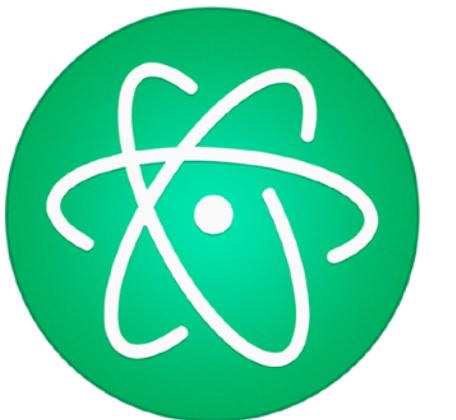
Applications Places System

- Nb
 - R
 - Firefox
- Accessories
- Big Data
- Databases
- Data Visualization
- Education
- Graphics
- Internet
- Office
- Other
- Programming
- Sound & Video
- System Tools
- Universal Access
- passwords.cxc

- Shells
 - Julia
 - NodeJS
 - PHP Shell
 - Pig Shell
 - Python (v2.7)
 - Python (v3.5)
 - R
 - Ruby
 - Scala
- H2O.ai
- IntelliJ IDEA Community Edition
- Jupyter Notebook
- PyCharm Community Edition
- Rodeo
- RStudio
- Spyder

Editors & Notebooks





Atom

Project Merlin Beta 3 (Small updates) [Running]

Applications Places System R Jupyter Nb Mon Aug 15, 12:26

python_demo.py — /usr/share/atom/resources — Atom

File Edit View Selection Find Packages Help

resources

python_demo.py

```
1 import pandas as pd
2 df = pd.DataFrame([2,3,4,5,6,7,8])
```



Project Merlin Development Version (Alpha 0.2) [Running]

Tue May 10, 23:47

Applications Places System Firefox R Nb

Home - Mozilla Firefox

File Edit View History Bookmarks Tools Help

Home localhost:8888/tree Search

jupyter

Files Running Clusters

Select items to perform actions on them.

Upload New

- Text File
- Folder
- Terminal

- Notebooks
- Julia 0.4.2
- Python 2
- Python 3
- R
- Ruby 2.1.5
- Scala 2.11
- pySpark (Spark 1.6.0)

anaconda3

Desktop

Documents

Downloads

metastore_db

Music

node_modules

A screenshot of a Linux desktop environment showing a Jupyter Notebook interface. The window title is "Home - Mozilla Firefox" and the tab title is "localhost:8888/tree". The main content area displays the Jupyter logo and navigation tabs for "Files", "Running", and "Clusters". Below the tabs, a message says "Select items to perform actions on them." A sidebar on the right lists file operations like "Upload" and "New", along with a list of kernel options: Text File, Folder, Terminal, Notebooks, Julia 0.4.2, Python 2, Python 3, R, Ruby 2.1.5, Scala 2.11, and pySpark (Spark 1.6.0). The file tree view shows standard directory names like "anaconda3", "Desktop", "Documents", "Downloads", "metastore_db", "Music", and "node_modules".



Jupyter Notebook

- Python 2 & 3
- R
- Ruby
- Scala
- PySpark



mongoDB®



neo4j



SQLite
GTK Cyber





MySQL®

Project Merlin Beta 3 (Before beaker 1.6 update) [Running]

Fri Aug 19, 00:14

MySQL Workbench

Local instance 3306 Local instance 3306

File Edit View Query Database Server Tools Scripting Help

SQL SQL Data Modeler Data Definition Data Manipulation Data Comparison Data Transformation Data Migration Data Cleaning Data Mining Data Science Data Visualization Data Integration Data Quality Data Governance Data Privacy Data Security Data Compliance Data Governance Data Privacy Data Security Data Compliance

MANAGEMENT

- Server Status
- Client Connections
- Users and Privileges
- Status and System Variables
- Data Export
- Data Import/Restore

INSTANCE

- Startup / Shutdown
- Server Logs
- Options File

PERFORMANCE

- Dashboard
- Performance Reports
- Performance Schema Setup

SCHEMAS

Filter objects

- phpmyadmin
- test
 - Tables

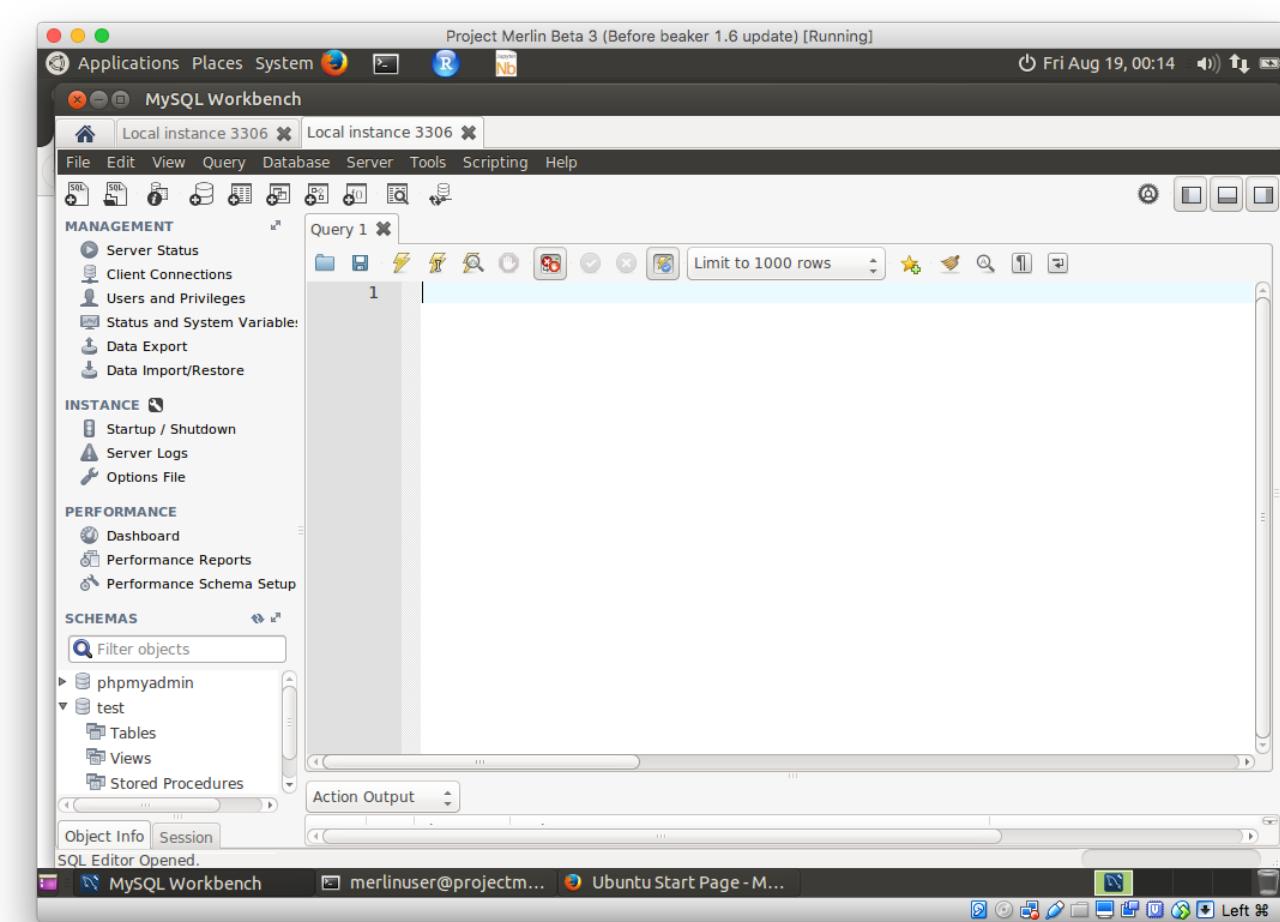
Query 1

Limit to 1000 rows

The screenshot shows the MySQL Workbench interface. The title bar indicates it's running on a local instance at port 3306. The main window has tabs for 'Local instance 3306' and 'Local instance 3306'. The menu bar includes File, Edit, View, Query, Database, Server, Tools, Scripting, and Help. The toolbar contains various icons for management, data manipulation, and analysis. On the left, there's a sidebar with sections for MANAGEMENT, INSTANCE, PERFORMANCE, and SCHEMAS. The SCHEMAS section shows the 'test' database expanded, revealing its tables. A query editor titled 'Query 1' is open, showing a single digit '1' in the results pane. A limit of 1000 rows is set for the query.



MySQL®



Project Merlin Beta 3 (Before beaker 1.6 update) [Running]

Fri Aug 19, 00:30

localhost / localhost | phpMyAdmin 4.4.13.1deb1 - Mozilla Firefox

localhost / localhost... +

localhost/phpmyadmin/index.php?token=62fb6ffb630fba16f23788ee21dc4e

Search

Applications Places System R Nb

phpMyAdmin

Server: localhost

Databases SQL Status Users Export Import Settings More

General Settings

- Change password
- Server connection collation: utf8mb4_unicode_ci

Database server

- Server: Localhost via UNIX socket
- Server type: MySQL
- Server version: 5.6.31-0ubuntu0.15.10.1 - (Ubuntu)
- Protocol version: 10
- User: merlinuser@localhost
- Server charset: UTF-8 Unicode (utf8)

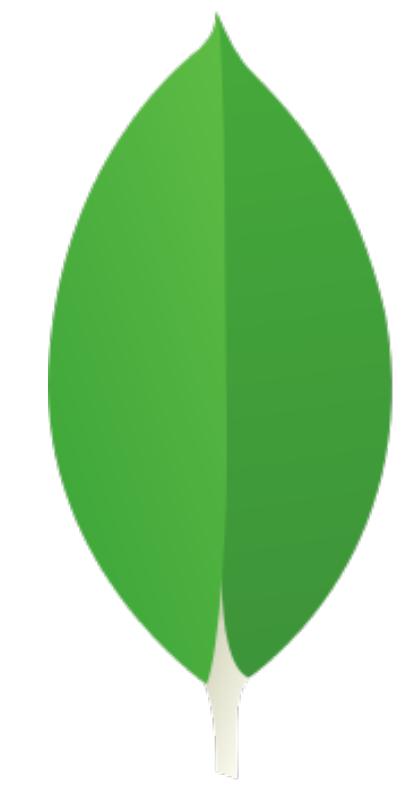
Appearance Settings

- Language: English
- Theme: pmahomme
- Font size: 82%

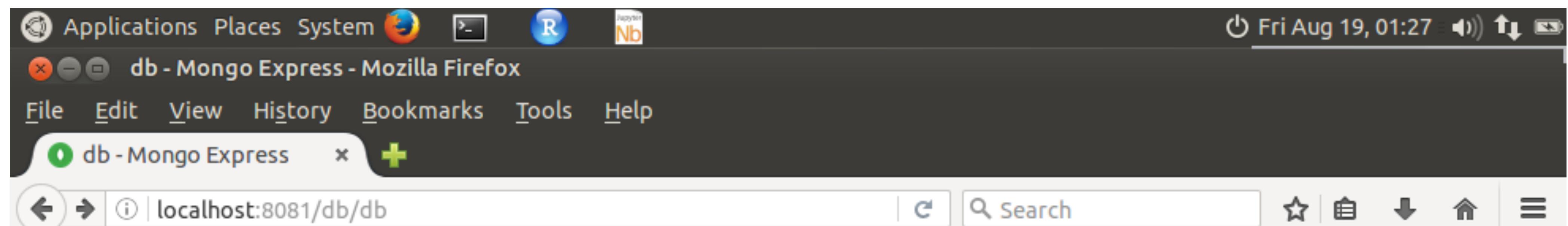
Web server

- nginx/1.10.1
- Database client version: libmysql - 5.6.31
- PHP extension: mysqli
- PHP version: 5.6.11-1ubuntu3.4

phpMyAdmin



mongoDB®



Mongo Express Database: db

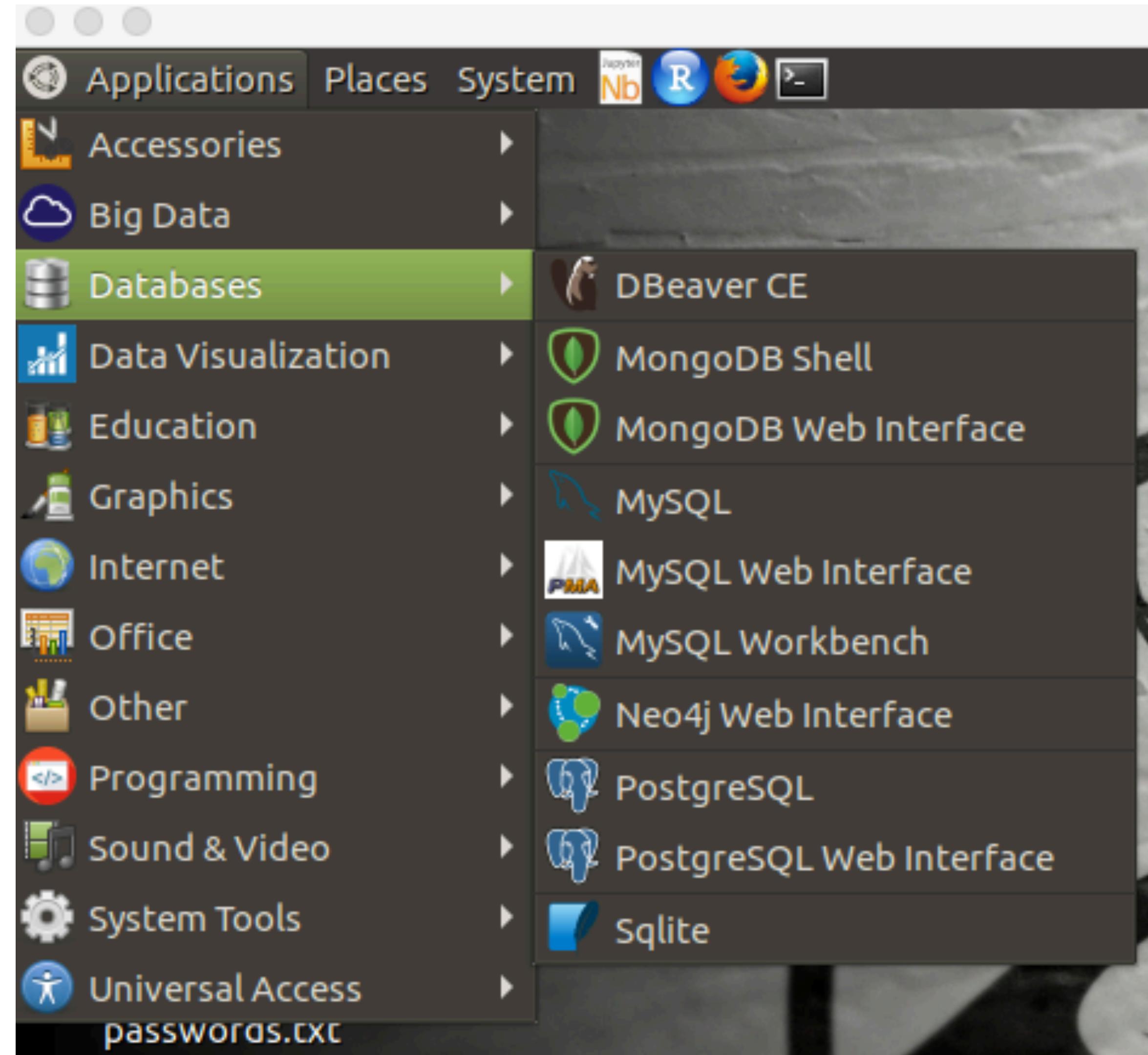
Collection "test" deleted!

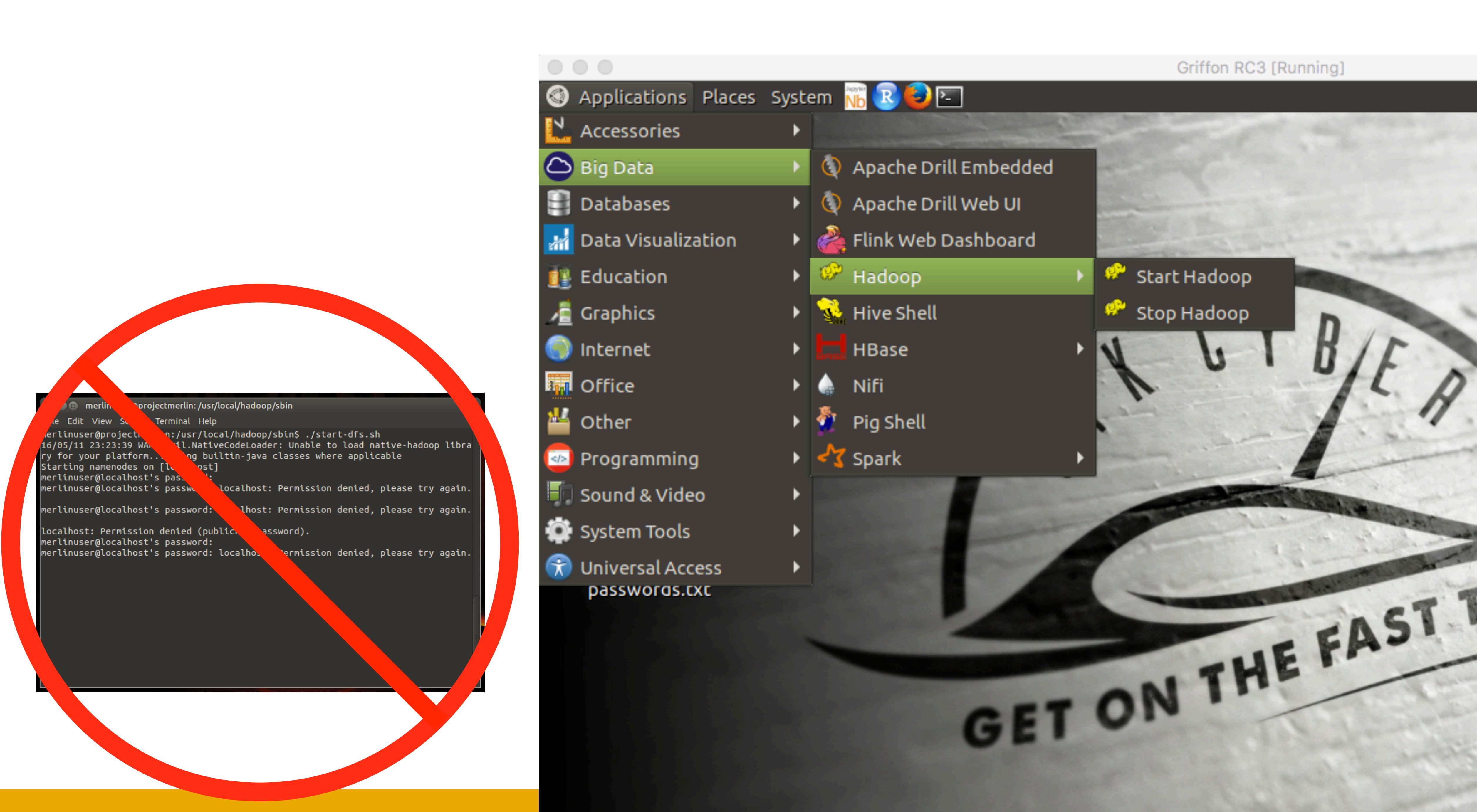
Collections

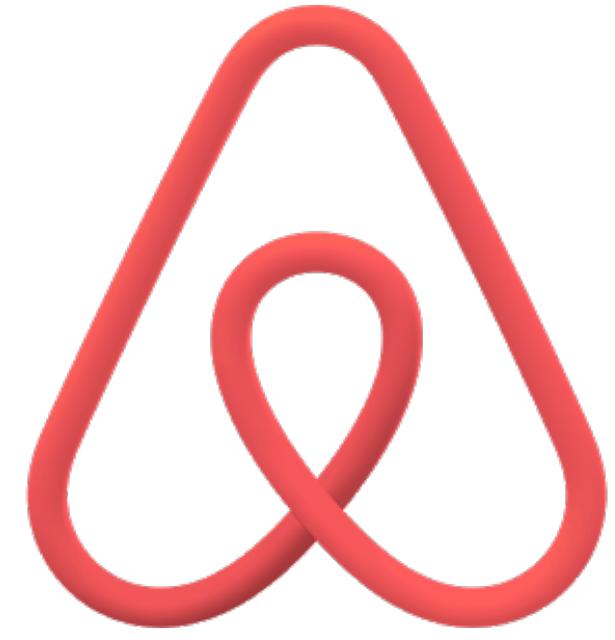
Collection Name [+ Create collection](#)

[View](#) [Export](#) [\[JSON\]](#) system.indexes [Del](#)

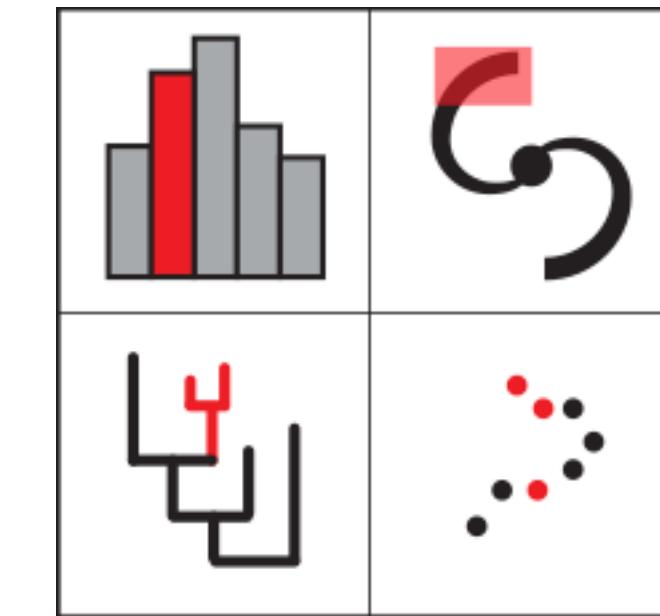
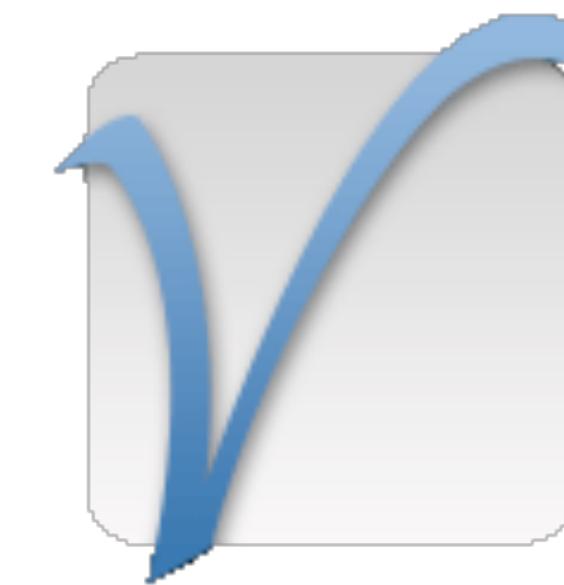
This screenshot shows the Mongo Express interface for the "db" database. A green notification bar at the top indicates that the "test" collection has been deleted. The main area is titled "Collections" and features a search bar for "Collection Name" with a "+ Create collection" button. Below this are three orange buttons labeled "View", "Export", and "[JSON]". To the right of these buttons is the collection name "system.indexes" in blue text. A red button with a trash can icon and the label "Del" is positioned next to it, indicating the option to delete the collection.







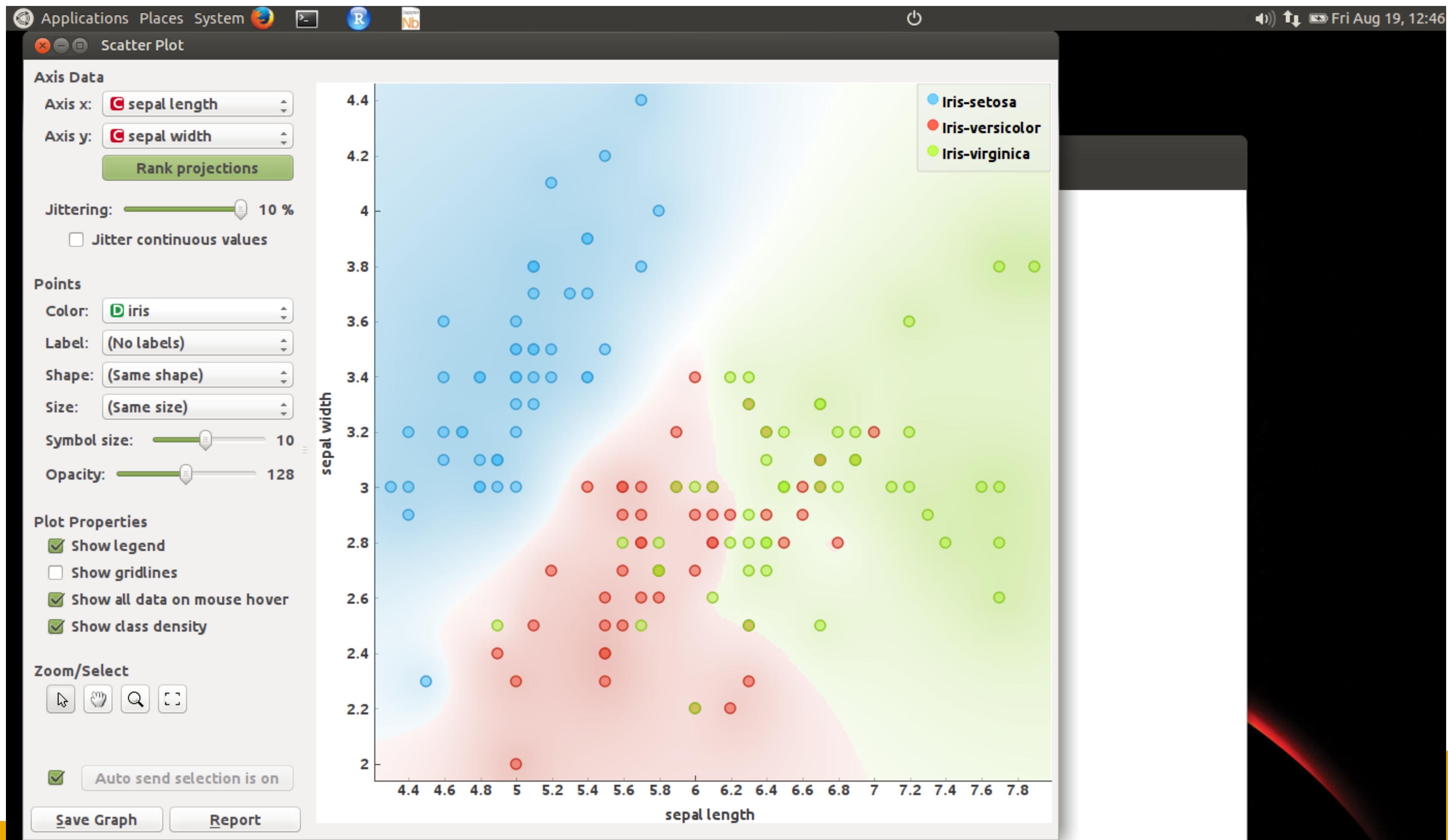
orange
DATA MINING
FRUITFUL&FUN





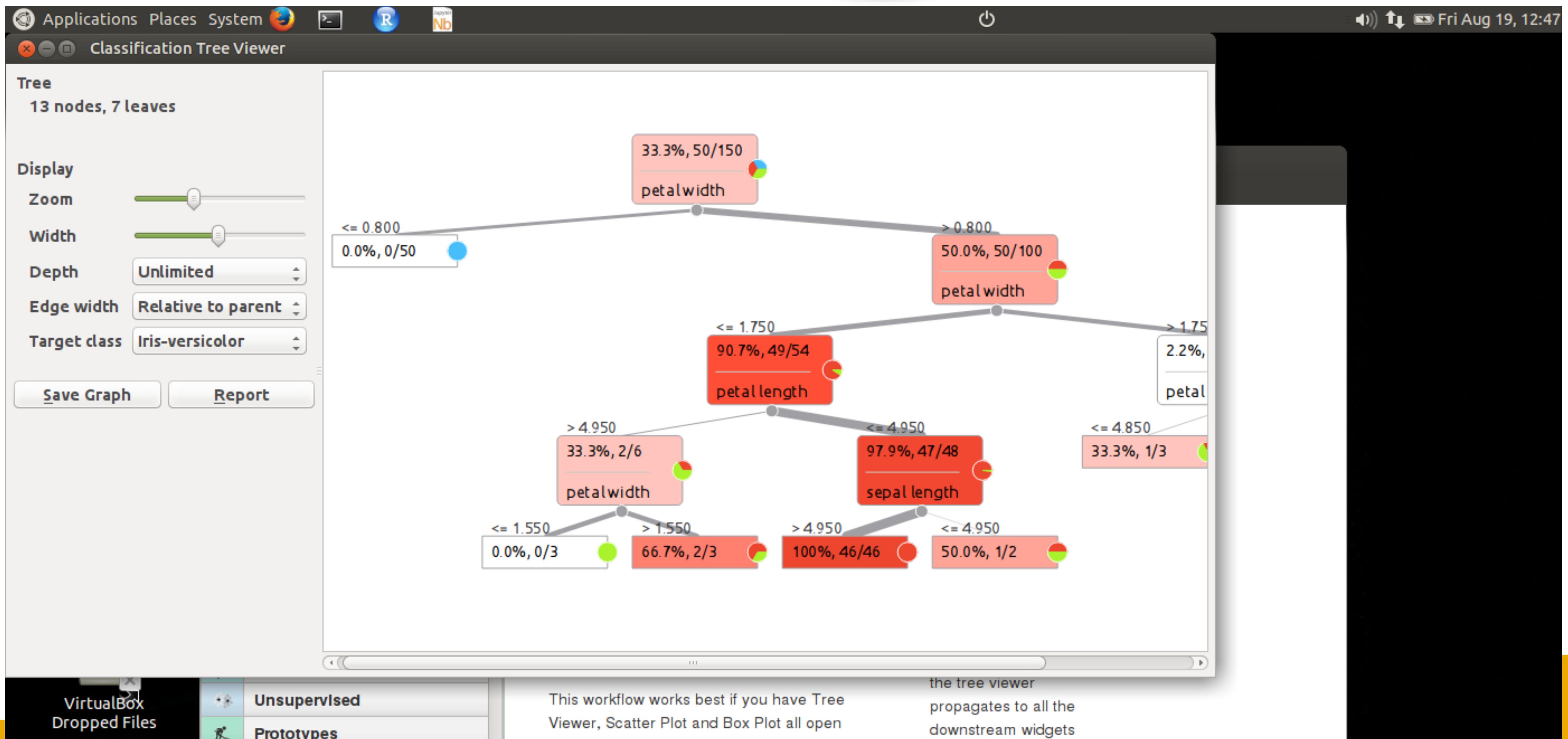
orange

DATA MINING
FRUITFUL&FUN



orange

DATA MINING
FRUITFUL&FUN



Questions?

Using Jupyter Notebook

Using Jupyter Notebook

The screenshot shows the Jupyter Notebook web interface. At the top, there's a navigation bar with links for Files, Running, Clusters, Formgrader, Assignments, BeakerX, and Nbextensions. On the far right of the header are Quit and Logout buttons. Below the header is a toolbar with Upload, New (with a dropdown arrow), and a refresh icon. A large arrow points from the text "Open a notebook" to the New button. To the right of the toolbar is a sidebar titled "Notebook:" with a list of kernel options: Bash, Clojure, Groovy, Java, Javascript (Node.js), Kotlin, PHP, Python 3, R, Ruby 2.5.1, SQL, Scala, Other:, Text File, Folder, and Terminal. At the bottom of the sidebar, it says "14 days ago". On the left side of the main content area is a file browser showing a directory structure with items like anaconda3, Desktop, Documents, Downloads, drill, metastore_db, Music, Pictures, Public, snap, sqldpad, Templates, and Videos. A "Select items to perform actions on them." message is displayed above the file list. The main content area is currently empty.

Open a notebook

Upload New

Notebook:

- Bash
- Clojure
- Groovy
- Java
- Javascript (Node.js)
- Kotlin
- PHP
- Python 3
- R
- Ruby 2.5.1
- SQL
- Scala

Other:

- Text File
- Folder
- Terminal

14 days ago

Using Jupyter Notebook

jupyter Untitled Last Checkpoint: 13 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Navigate Widgets Help Snippets Trust

Run

Demo Notebook

This is a markdown cell. It contains the instructions as to how to do the exercises.

In [1]:

```
1 #This is an executable cell
2 print( "This is a python cell")
```

This is a python cell

In []:

```
1 |
```

A screenshot of the Jupyter Notebook interface. The title bar shows 'jupyter Untitled Last Checkpoint: 13 minutes ago (unsaved changes)'. The toolbar includes standard file operations like File, Edit, View, Insert, Cell, Kernel, Navigate, Widgets, Help, and Snippets, along with a 'Trust' button. Below the toolbar is a toolbar with various icons for file operations, code execution, and data visualization. A large black arrow points from the text 'Run a cell' down to the 'Run' button in the toolbar, which is circled in black. The main workspace displays a 'Demo Notebook' section with a markdown cell containing instructions. Below it is an executable cell with Python code. The output of the executable cell is 'This is a python cell'. A new input cell, 'In []:', is shown at the bottom, with the number '1' entered. The entire interface is set against a light gray background.

Using Jupyter Notebook

The screenshot shows the Jupyter Notebook interface with the following elements:

- Header:** jupyter Untitled Last Checkpoint: 18 minutes ago (unsaved changes), Logout, Python 3.
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Navigate, Widgets, Help, Snippets.
- Icon Bar:** Includes icons for file operations (New, Open, Save, etc.), Run, Cell Types (Code, Markdown, etc.), and various notebook-related functions.
- Main Area:** Demo Notebook. A markdown cell contains the text: "This is a markdown cell. It contains the instructions as to how to do the exercises." An executable cell (In [1]) contains the code:

```
1 #This is an executable cell
2 print( "This is a python cell")
```

 and outputs "This is a python cell".
- Variable Inspector Panel:** Shows a table with one entry: X x list 88 [4, 5, 6].
- Variable Explorer Icon:** Located in the toolbar, circled with a black oval and connected by a vertical arrow to the Variable Inspector panel.
- Text Labels:** "Variable Explorer" is written below the main area.

Questions?