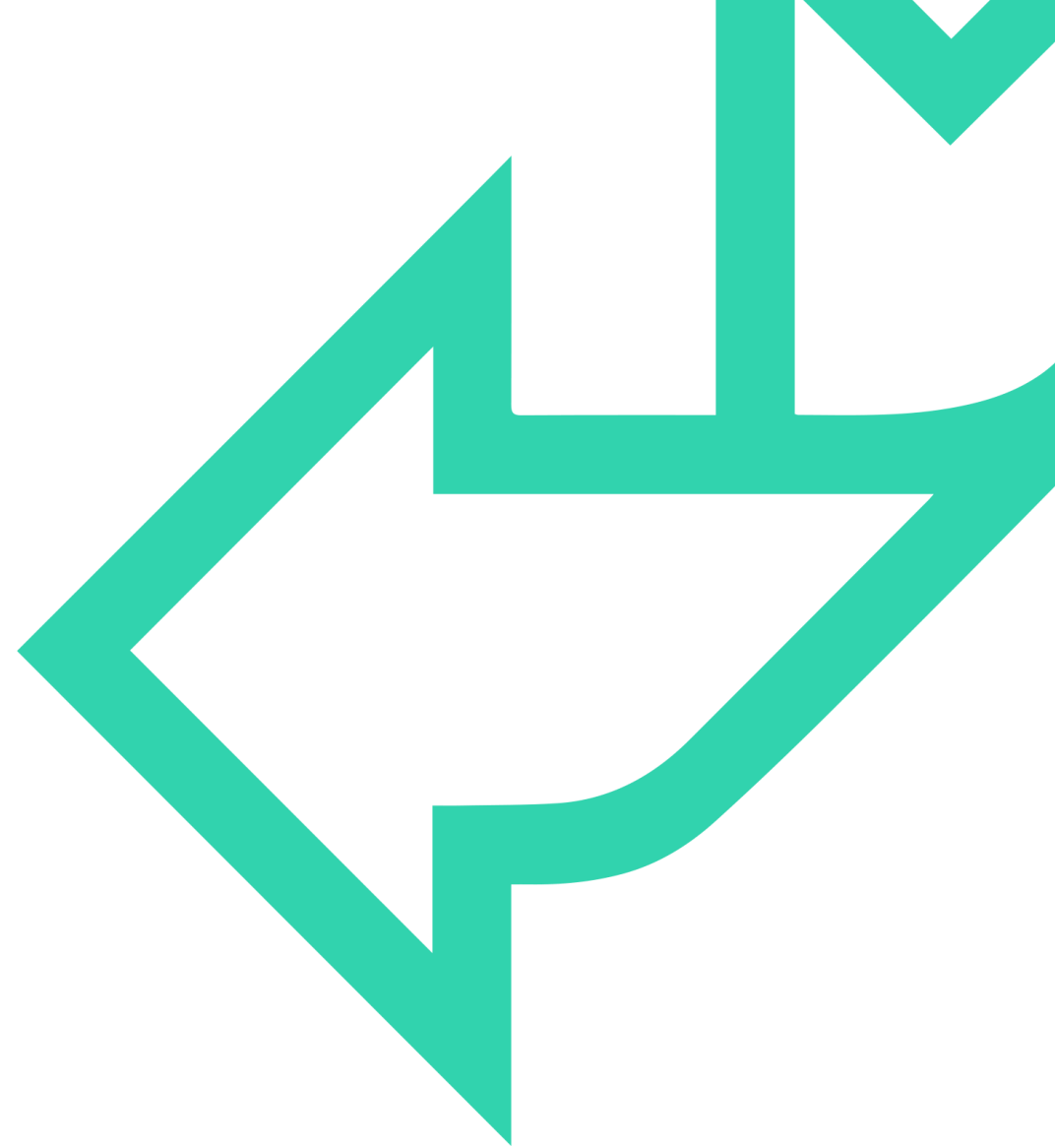




Data Science

the process of extracting knowledge from data automatically





DATA



What is data?

Information

Relevant

Raw

Data Types? (forms)

Structured Data

Tabular

SQL Databases

Unstructured Data

Photos

Sounds

Videos

Collected from where?

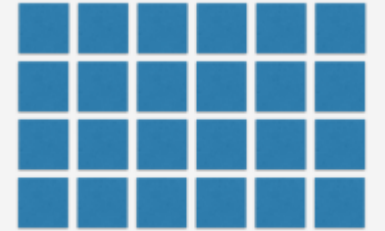
Social Media

Bank Accounts

..... everywhere

Smart devices

GPS



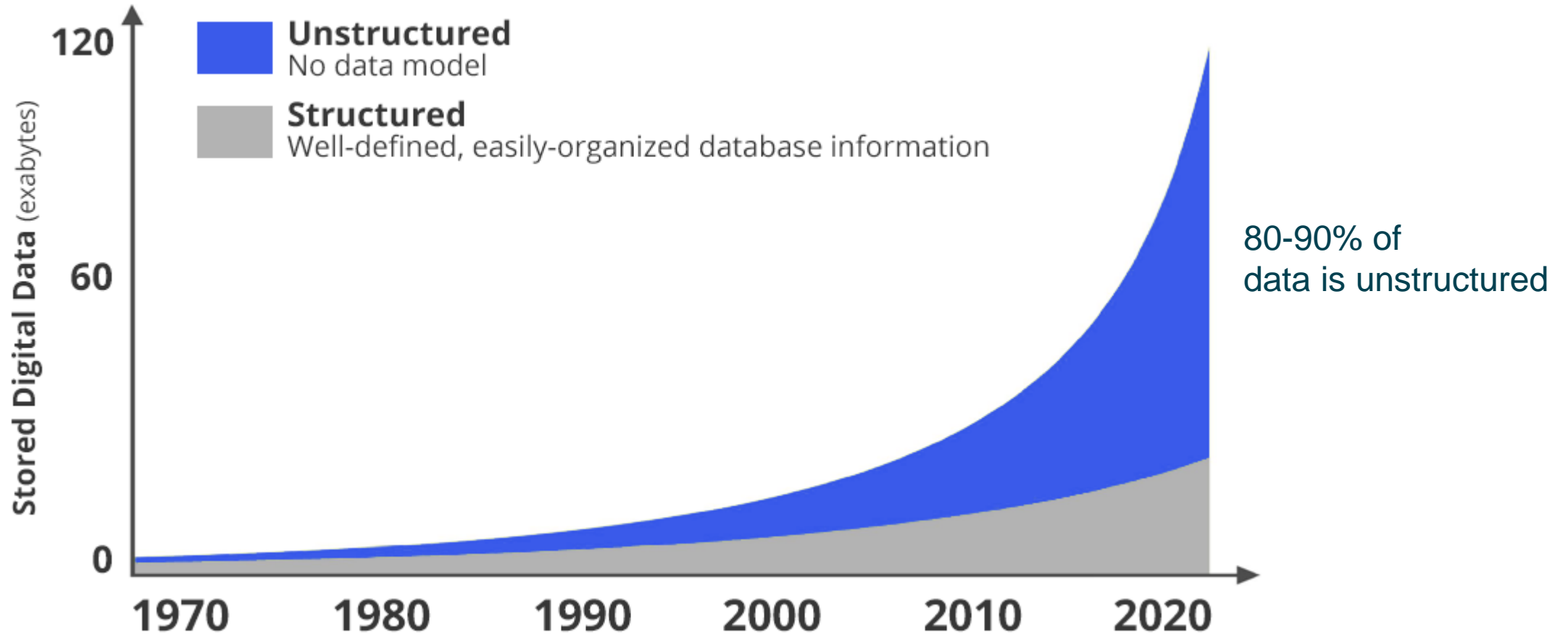
What you find in a DB
(typically)



What you find in the 'wild'
(text, images, audio, video)



volume of unstructured data is growing & that growth is accelerating





volume of unstructured data is growing
& that growth is accelerating





SCIENCE



Every baby knows the
scientific method!



- ✓ Observe
- ✓ Question
- ✓ Research



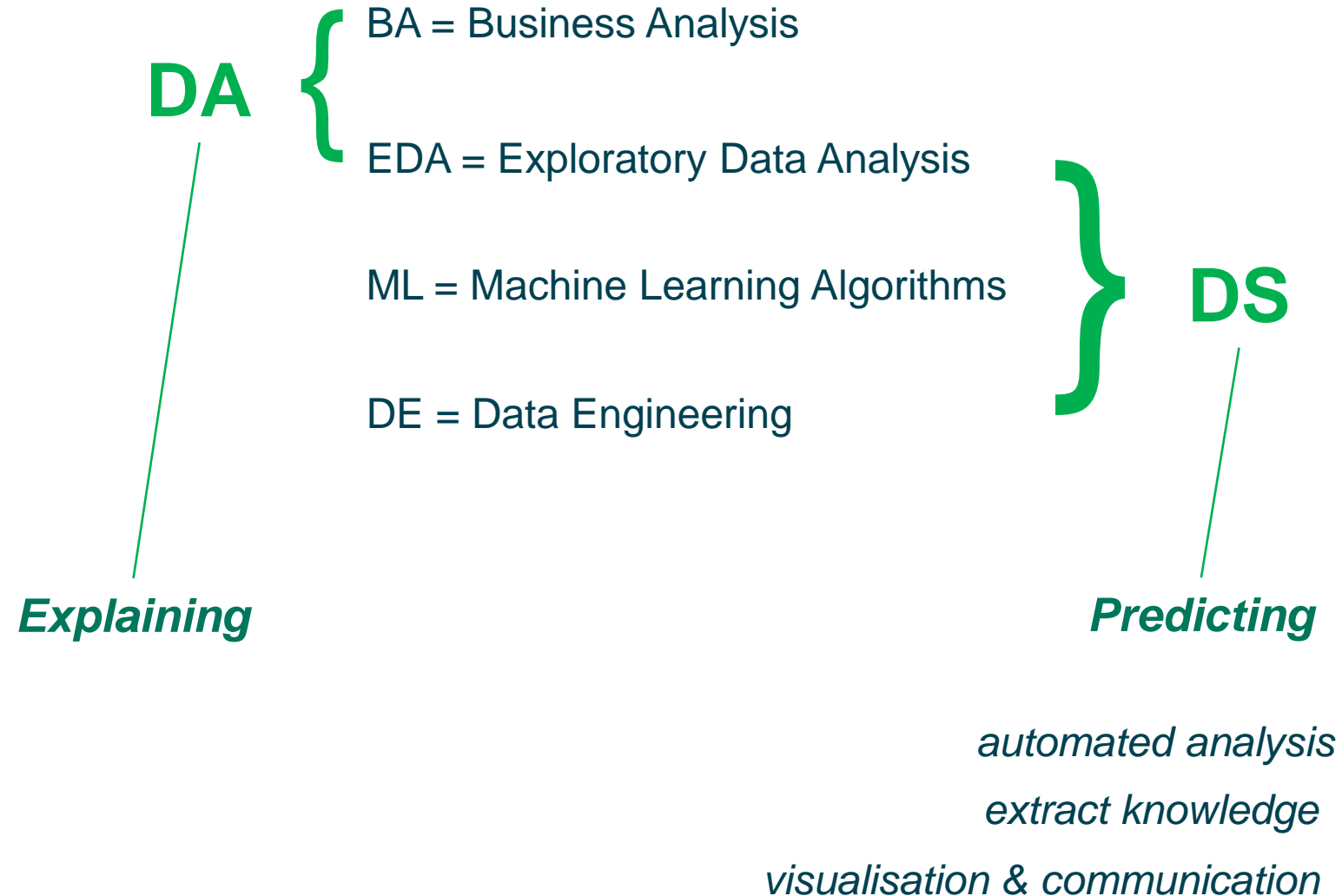


DATA SCIENTIST

DS = evidence-based methodology
for solving problems using
computational and analytical tools

Data Science = applied branch of statistics

data → valuable asset → insight → decisions → actions





DATA SCIENTIST

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

Wills (2012, Director of Data Science at Cloudera)

“... the sexiest job ... will be statisticians....

data ... understand it extract value from it ... visualise it ...

Varian (2008, Chief Economist at Google)

“A data scientist is somebody who is inquisitive, who can stare at data and spot trends... ... like a Renaissance individual who really wants to learn and bring change to an organization.”

Bhambhi (2012, VP of Big Data products at IBM)



DATA SCIENTIST

“effective data science should not be treated like just another business process, and can not be operationalized assembly-line style.

Data science -- as the name suggests -- is a mode of inquiry and exploration similar to “real” science.

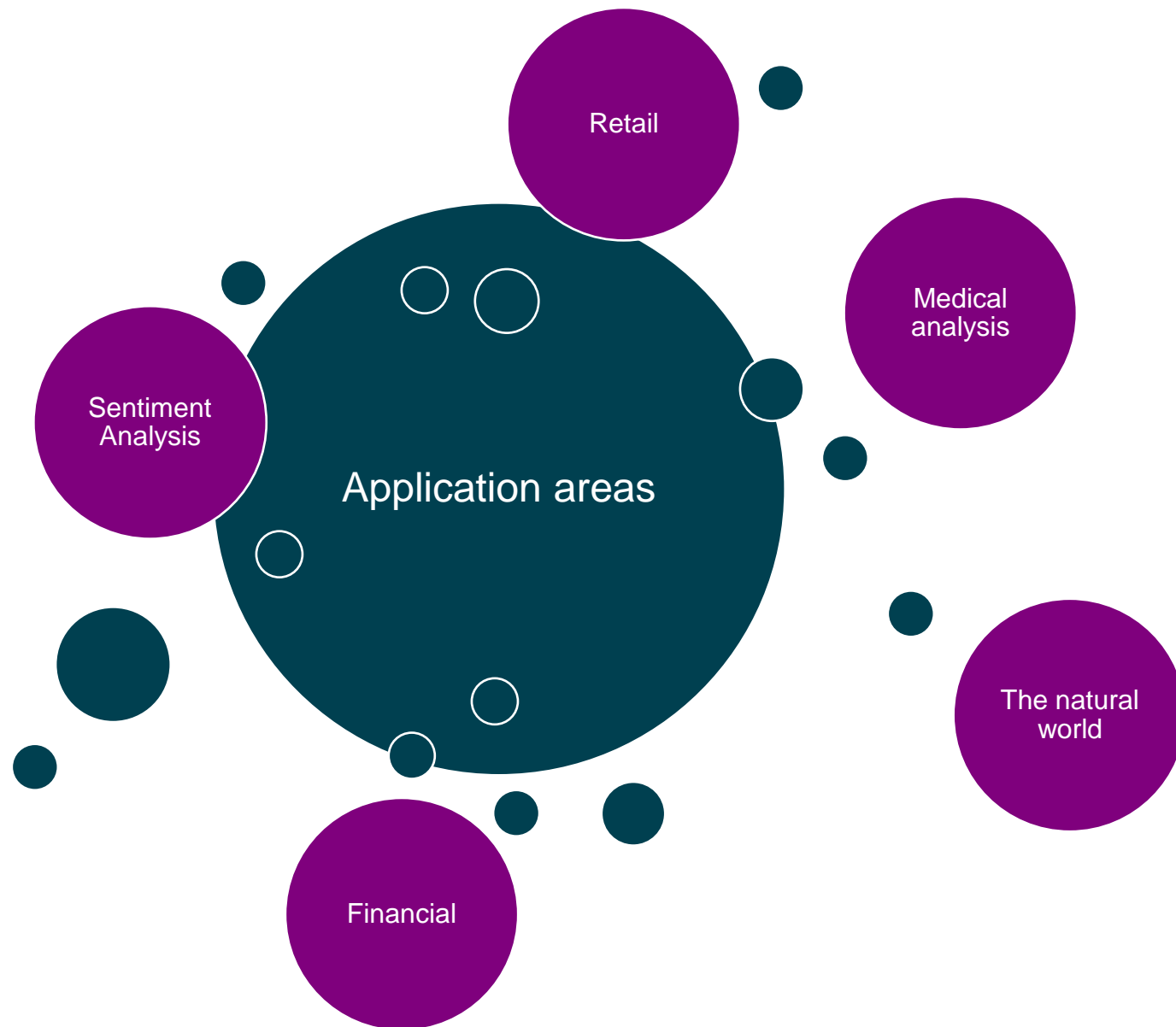
Just as a physicist uses math to reason about the natural world, data scientists harness mathematical and computational tools to reason about the business world.”

Peter Wang (2019, CEO Anaconda)





DATA SCIENCE APPLICATIONS





DATA SCIENCE APPLICATIONS



Marketing

Optimizing ads

Identifying the best ad to run online

Forecasting Churn

Predict which type of customers are most likely to leave your service

Segmenting Customer Base

Deeply understand customer base for better tailored products and services

Finance

Stock Price Prediction

Forecasting share indices and specific stock prices

Improved Customer Service

Using text data for chatbots and call center routing

Fraud Detection

Identify fraudulent transactions based on key patterns



DATA SCIENCE APPLICATIONS



Sentiment Analysis - Social media – can we gain intuition/vibe?

Medical analysis - Are you likely to get a certain illness based on your lifestyle (or other factors)?

The natural world - What is going on with the weather? snow?

Financial - Is the market going to go up or down?

Retail - Who is most likely to buy something?

Supermarkets - What stock to put in the front of the shops?
What discounts to give?

Personal - Why can't I hit 10,000 steps a day?



NOMENCLATURE

Data scientists

“example”

“feature”

“label”

Statistician

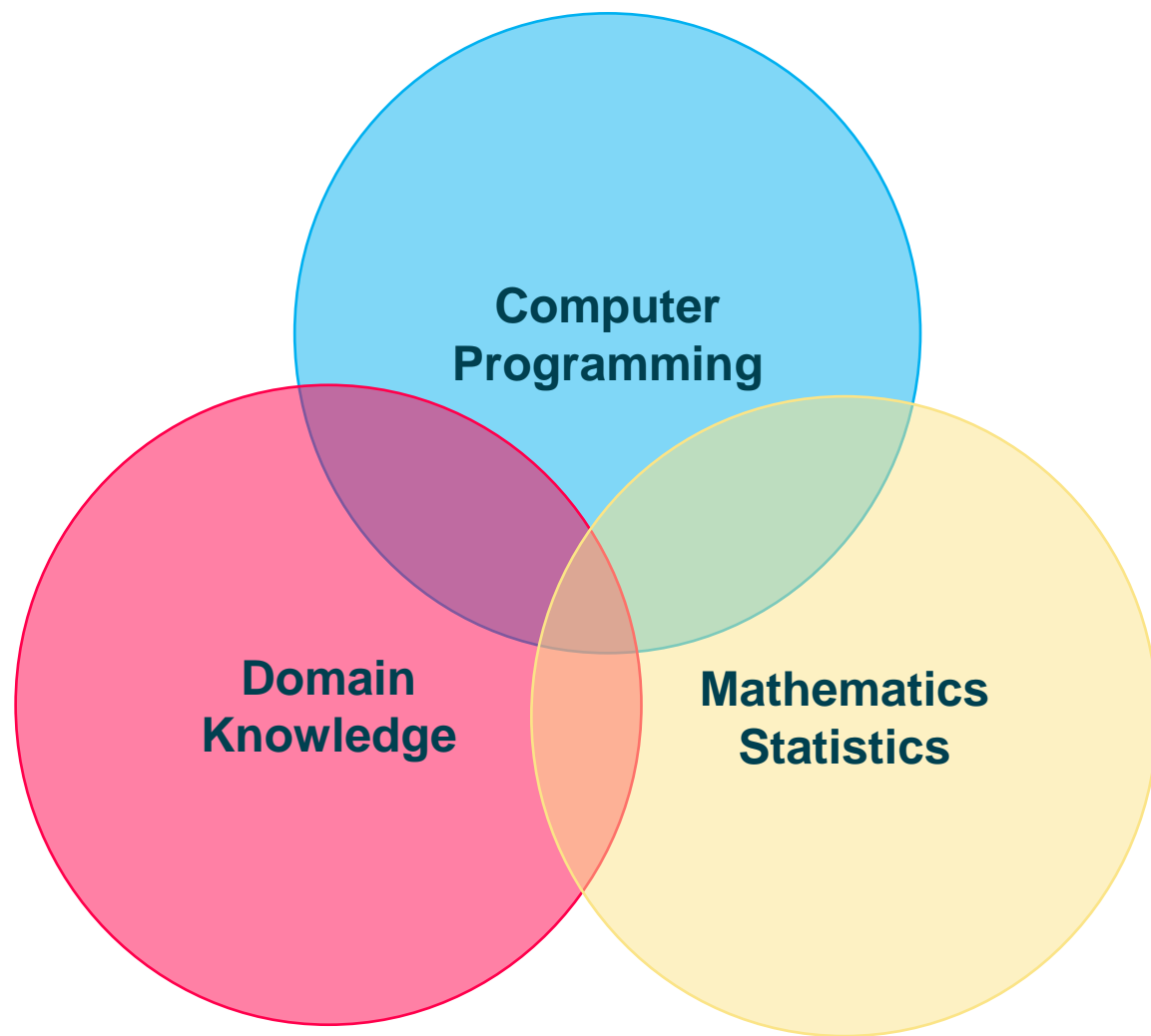
“observation”

“predictor”
“independent variable”

“response”
“dependent variable”

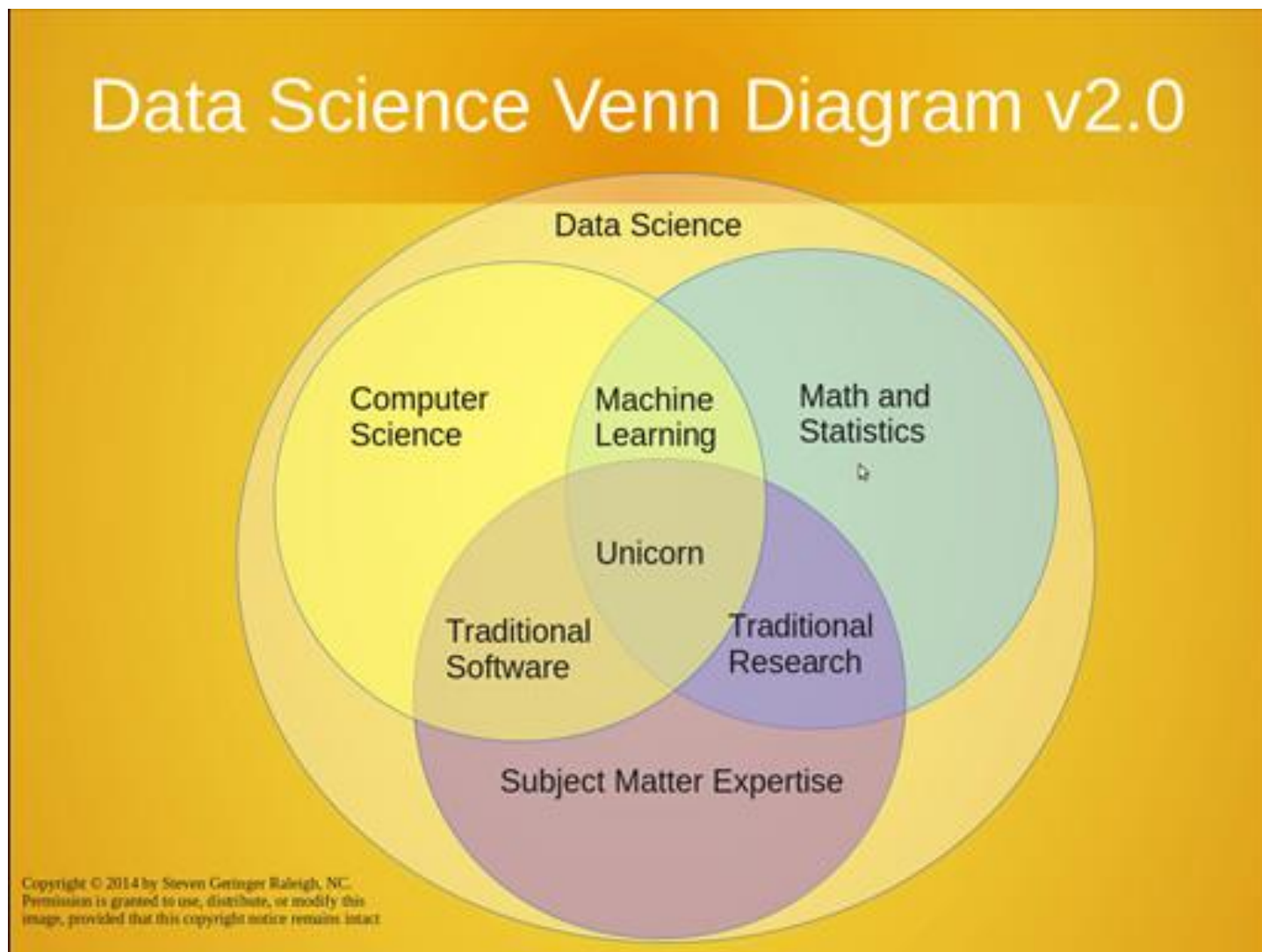


DATA SCIENCE SKILLSET





DATA SCIENCE SKILLSET





DATA SCIENCE SKILLSET

- ✓ **Domain Knowledge**
- ✓ **Programming**
- ✓ **Statistics**
- ☐ Curiosity
- ☐ Visualisation
- ☐ Communications
- ☐ Storytelling
- ☐ Project Management/DevOps
- ☐ Databases
- ☐ Data Mining
- ☐ Machine Learning





DATA SCIENCE SKILLSET

- ✓ Domain Knowledge
- ✓ Programming
- ✓ Statistics
- ☐ Curiosity
- ☐ Visualisation
- ☐ Communications
- ☐ Storytelling
- ☐ Project Management/DevOps
- ☐ Databases
- ☐ Data Mining
- ☐ Machine Learning

In many walks of life evolution ***selects against*** the kind of person who decides to find out what happens “if I push that button”.

In Data Science it ***selects for*** it.

2012, Ross
(Director Data Science at Teradata)

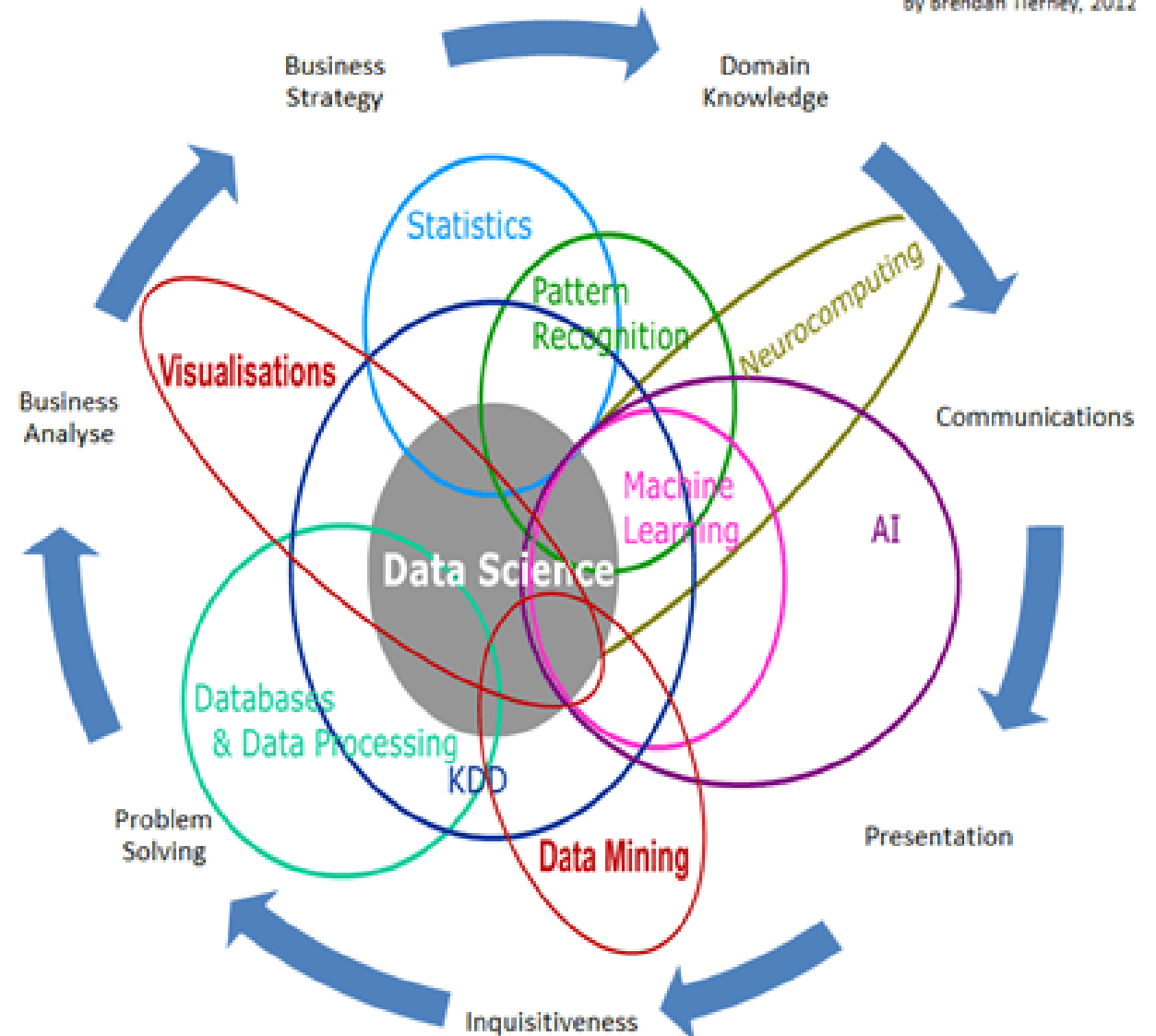


DATA SCIENCE SKILLSET



Data Science Is Multidisciplinary

By Brendan Tierney, 2012





ANALYTICS

Descriptive Analytics

Making data visible and getting it in the right hands

Predictive Analytics

Predict the future with data

Prescriptive Analytics

Making data-driven decisions





MATHEMATICS & STATISTICS



Probability

- Belief of event occurrence
- Bayes rule

Statistics

- Random Variables (Features)
- Samples vs Populations
- Distributions
- Inference

Linear Algebra

- Transformations on data sets
- Weighted Sums
- Models (Tensor, Matrix, Vector)

Calculus

- Areas and Rates
- Rates of error functions
- Parameter Optimisation



PYTHON LIBRARIES

NumPy

- fast numerical arrays
- optimized fortran and C extensions

Pandas

- numpy wrapper
- provides "data frames"
- tabular model over numpy arrays

matplotlib

- visualization and plotting

seaborn

- convenience matplotlib wrapper

Bokeh

- alternative graphing library (for the web)
- especially useful for geoplots and other complex plots

SciKit Learn

- comprehensive machine learning library
- provides good-enough implementations of most key algorithms

Tensorflow

- fast (concurrent, distributed, gpu) numerical computing library
- describes computations as optimizable graphs

Keras

- tensorflow (et al.) wrapper providing neural network abstractions



TOOLS

Programming Languages



Databases



Command Line Tools



Spreadsheets



Business Intelligence Tools



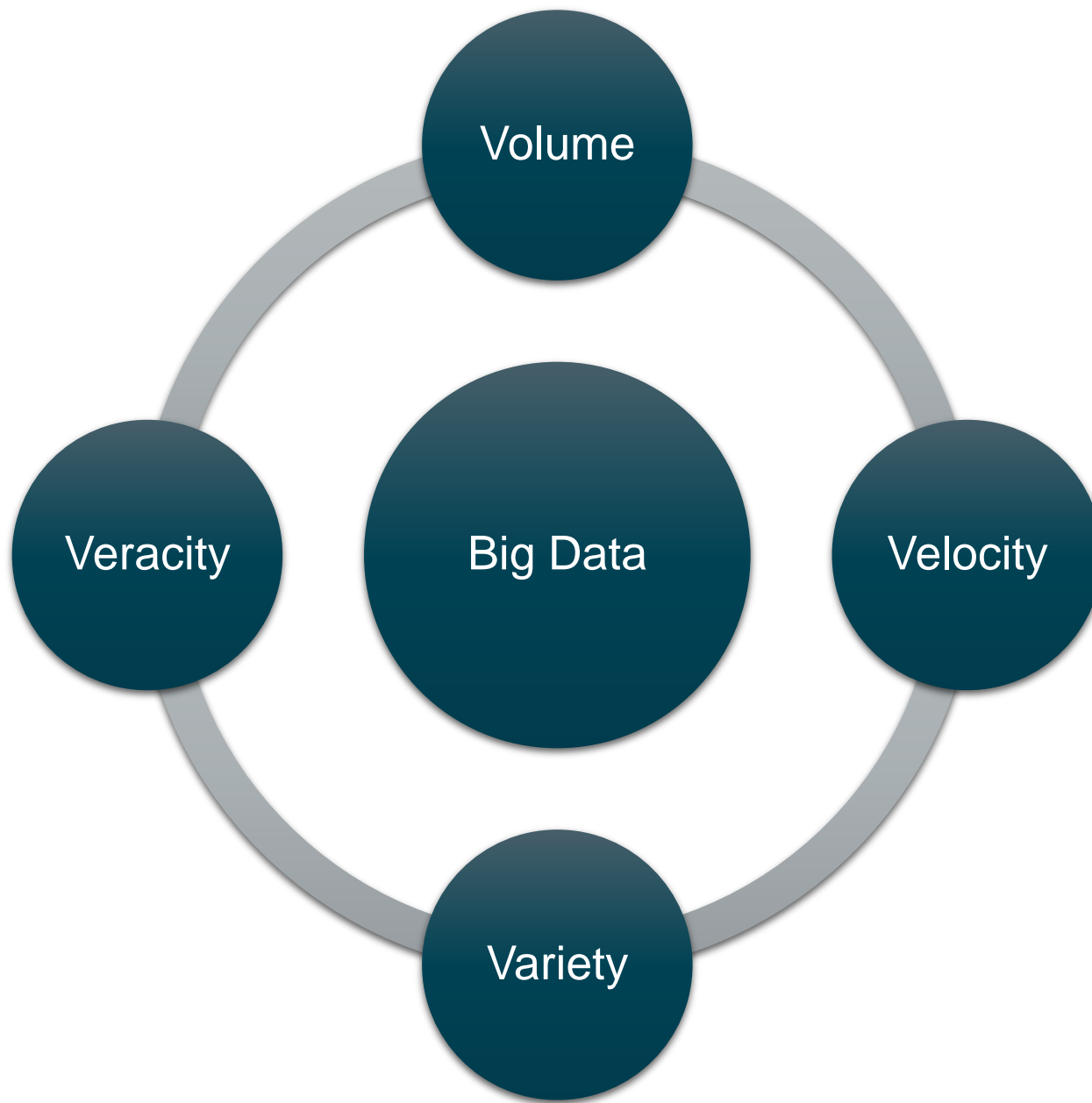
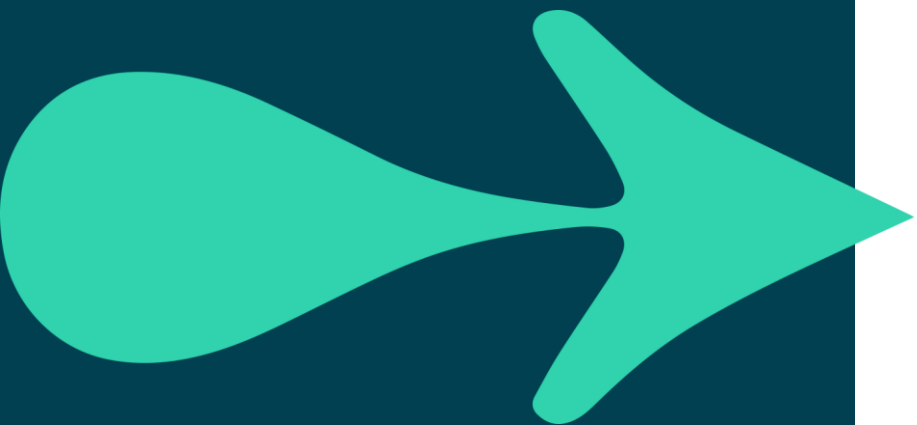


Big Data

collection, ingestion, processing and analysis of datasets **too large**, and generated **too quickly**, to be analyzed effectively by traditional analytical tools and methodologies



BIG DATA



VOLUME

Huge amount of data



VERACITY

Inconsistencies and uncertainty in data



VARIETY

Different formats of data from various sources



VELOCITY

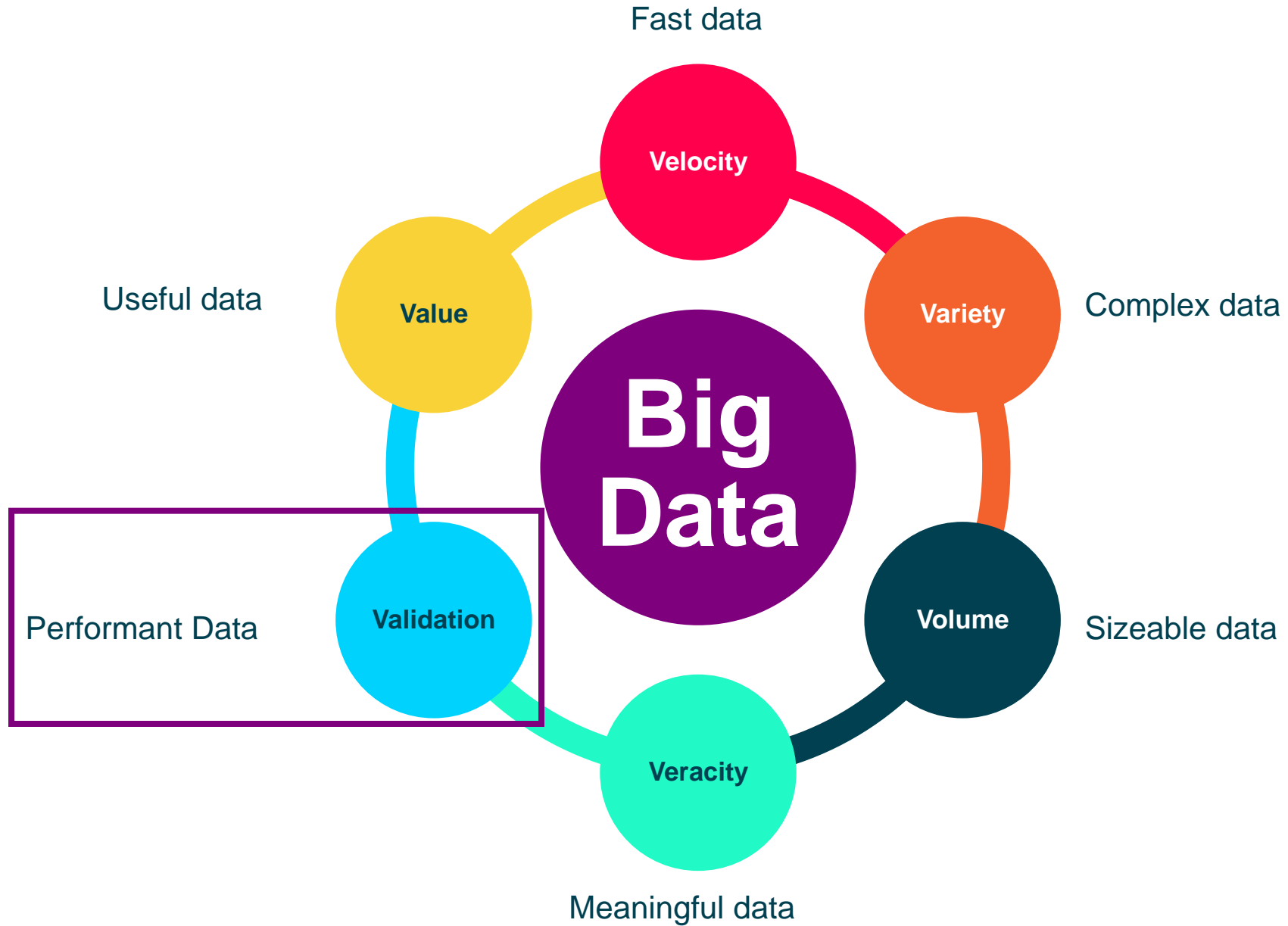
High speed of accumulation of data

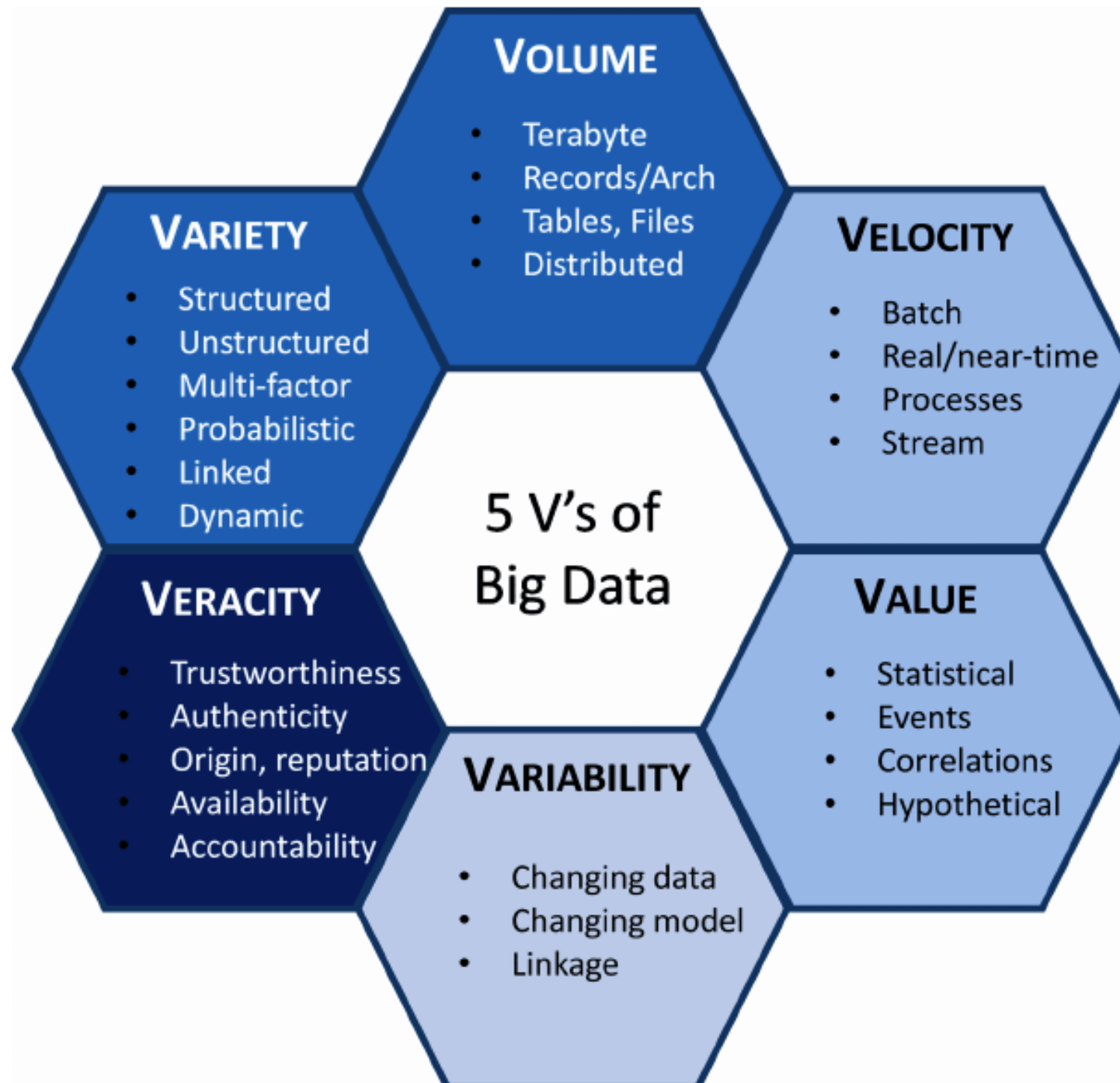


VALUE

Extract useful data











BIG DATA: SOURCES

Human Generated

- Photos, videos
- Text
- Click Like
- Web search
- Emails & SMS
- Online Purchases

Machine Generated

- Cell Phones , GPS
- Industrial Process Monitoring
- Climate monitoring
- Medical Devices
- IoT



BIG DATA VS SMALL DATA

Small Data

Goals: Specific

Location: In one place

Structure: structured

Preparation: by end user for their own purpose

Longevity: kept for specific time

Reproducibility: can be reproduced

Stakes: data loss cost is limited

Analysis: can be analysed at once in one PC

Big Data

Goals: Evolving

Location: distributed

Structure: unstructured & multi formats

Preparation: by many, used by many

Longevity: kept for long period of time

Reproducibility: may not be possible

Stakes: data loss cost is high

Analysis: Multiple steps as files can be in different places and formats

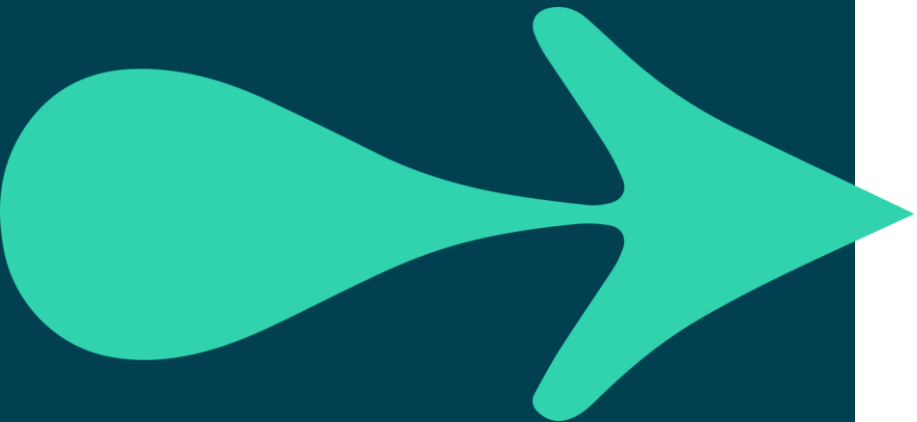


APPENDIX: Jargon





ARTIFICIAL INTELLIGENCE



Neural Networks

- A particular pattern-finding algorithm

Deep Neural Networks

- A particular generalization of the neural network algorithm

Artificial Intelligence

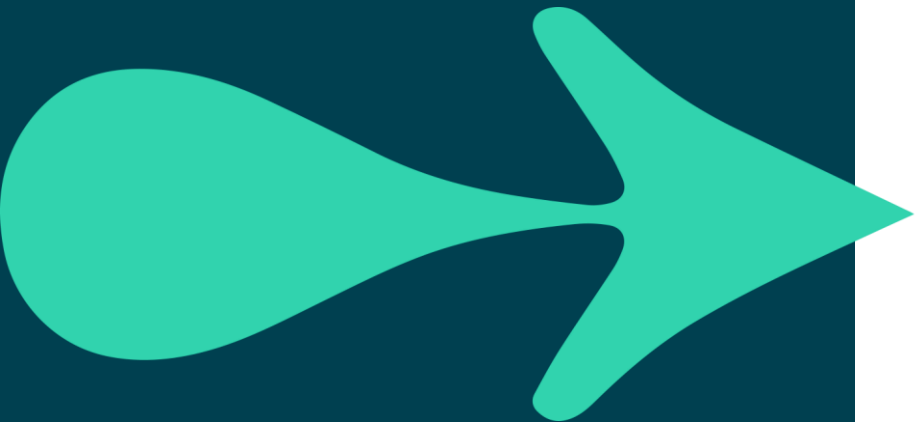
- (broad def.) the project of simulating animal intelligence

Artificial Intelligence

- whatever artificial system is the best at decision making
- 50s: computers running simple programs
- 80s: computers running expert programs
- 00s: computers running machine learning programs
- 10s: computers running machine learning with neural networks



DATA SCIENCE



Statistics

- Describing and finding patterns in data

Patterns

- Correlations between variables

Statistical Inference

- Finding novel patterns / making predictions / generalizing from observation

Machine Learning

- Computational Statistical Inference = Statistical Inference with Computers

Learning

- Finding the “best” patterns



ROLES

Data Analysis

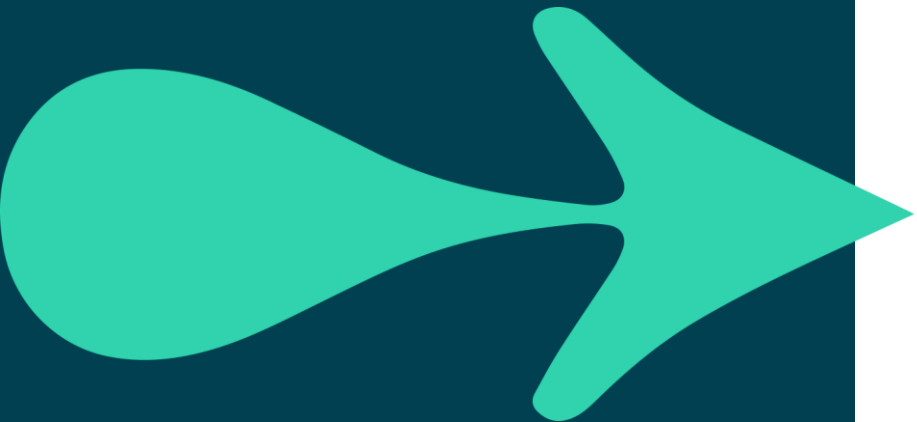
- Historical Trends, Data Fact Finding

Data Science

- Future Trends, Inference & Prediction
- Domain Understanding
- Data Understanding
- Data Exploration, Preparation
- Statistical Modelling
- Evaluation of Models
- Deploying solutions



DATA INSIGHT



Data Scientist

- Machine Learning Developer
- Researcher

Data Engineer

- Big Data Engineer
- DataOps
- Data Science with DevOps

Data Leader

- Project Manager