

Data Science for Public Policy

Slava Mikhaylov
Professor of Public Policy and Data Science
Institute for Analytics and Data Science
Department of Government
University of Essex

Concept of Data Science

Data Scientist:

The Sexiest Job of the 21st Century

Meet the people who
can coax treasure out of
messy, unstructured data.
by Thomas H. Davenport
and D.J. Patil

W

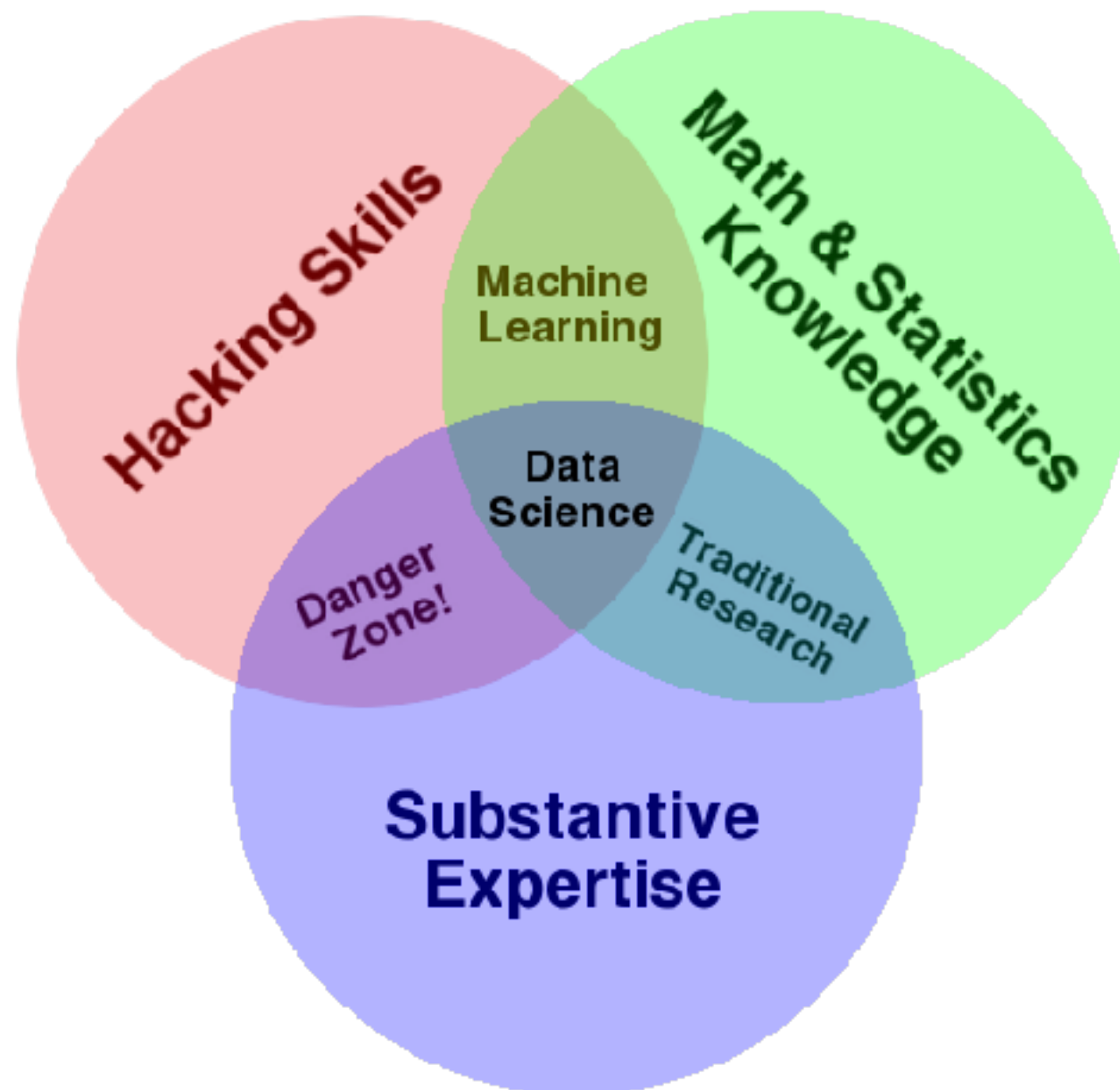
hen Jonathan Goldman arrived for work in June 2008 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know anyone. So you just stand in the corner sipping your drink—and you probably leave early."

98 Harvard Business Review October 2012



"I keep saying the sexy job in the next ten years will be statisticians. People think I'm joking, but who would've guessed that computer engineers would've been the sexy job of the 1990s?" Hal Varian (Chief Economist at Google, 2009).

What is Data Science?



System approach to data science

Quadruple Helix Innovation

Government, Academia, Industry and Citizens collaborating together to drive structural changes far beyond the scope of any one organization could achieve on it's own

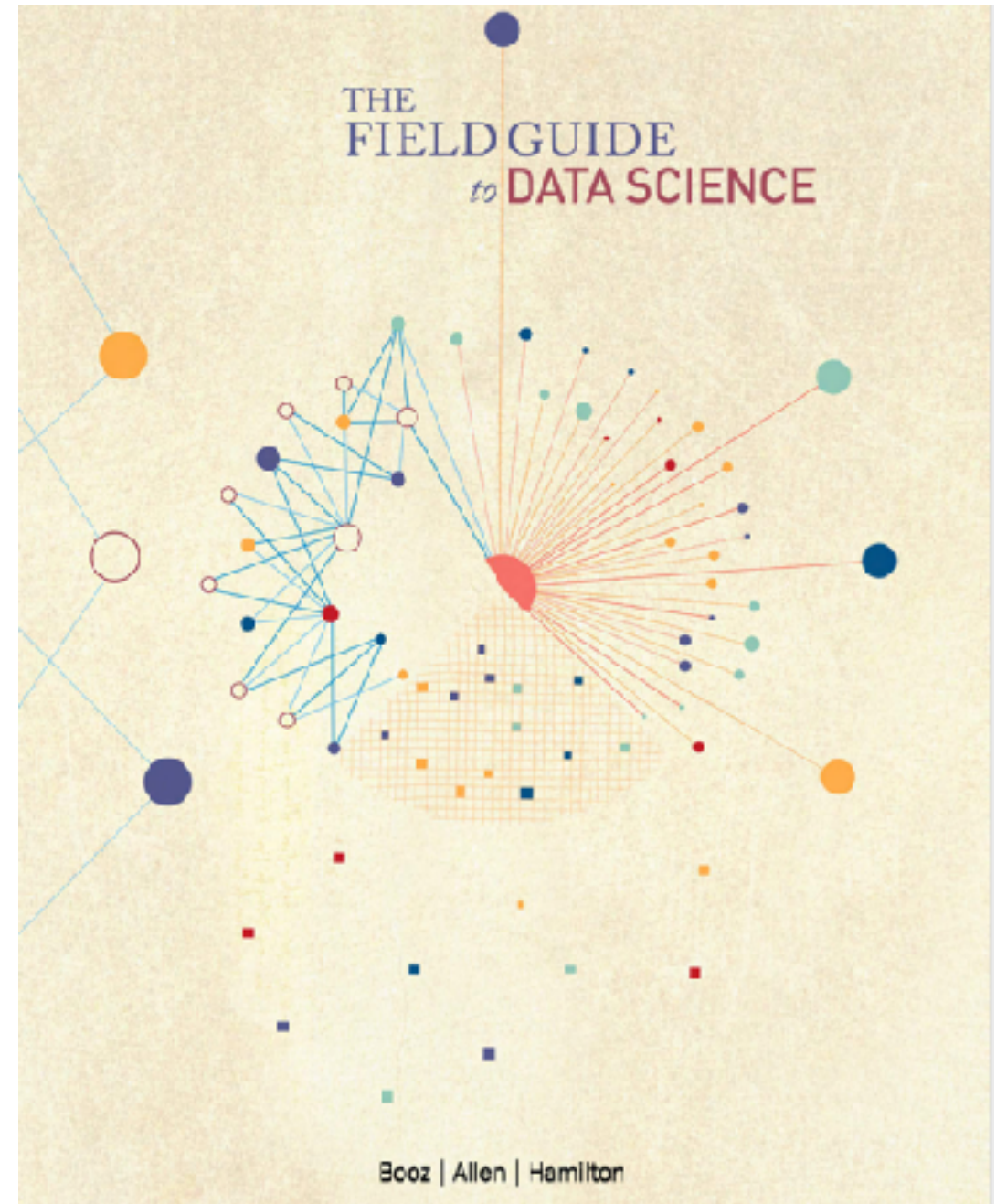


“Research in Big Data should be grounded in the quadruple helix model where civil society joins with business, academia, and government sectors to drive changes far beyond the scope of what any organization can do on their own.”

Intel Corp policy position paper on Big Data

Practical perspective

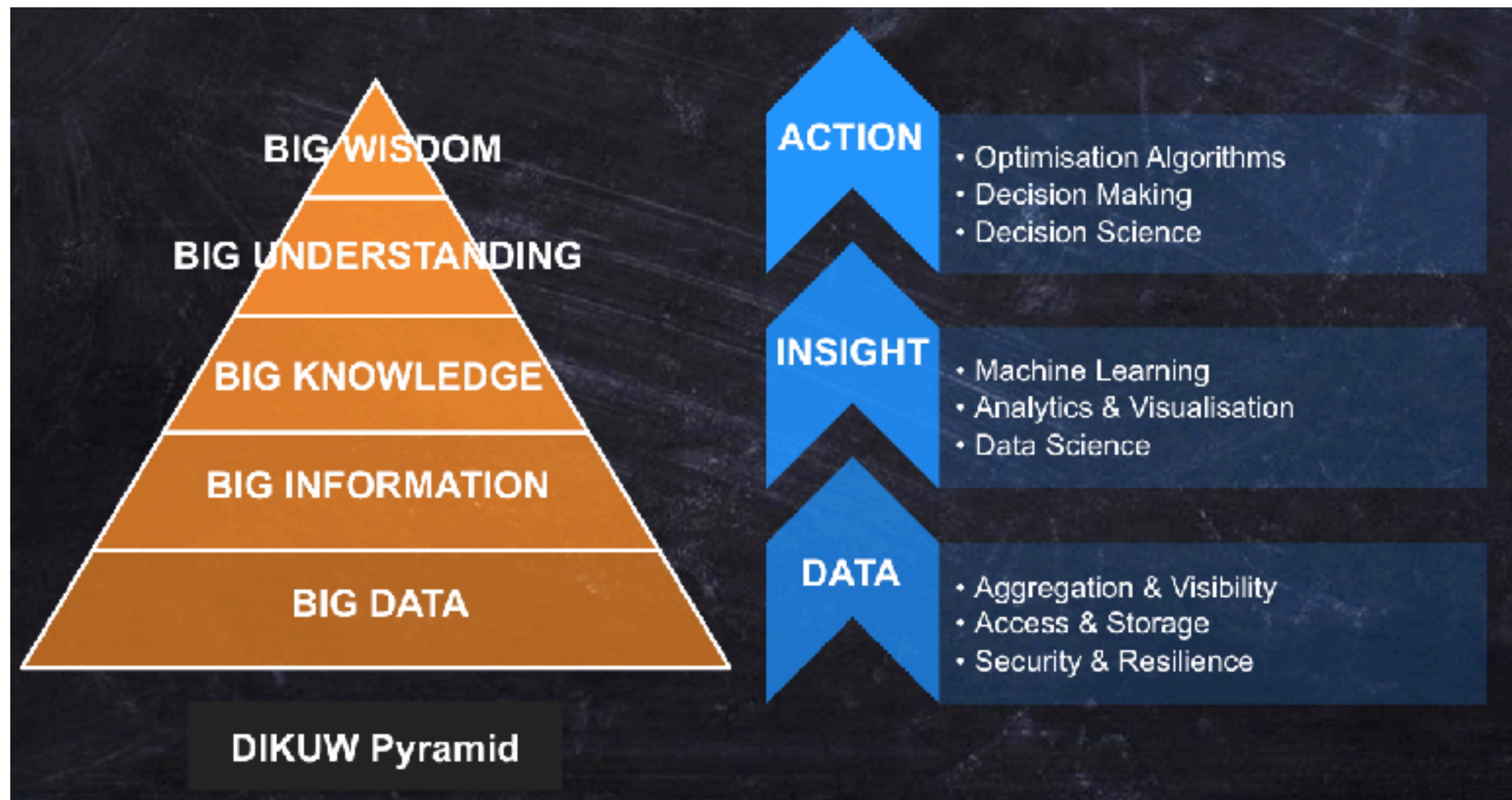
- Data Science is the art of turning data into actions
- It's all about the tradecraft.
- Tradecraft is the process, tools and technologies for humans and computers to work together to **transform data into insights**.



Practical perspective

- Data Science tradecraft creates data products
- Data products provide actionable information **without** exposing decision makers to the underlying data or analytics.





DIKUW Pyramid

Inductive and deductive reasoning

- Data Science supports and encourages shifting between deductive (hypothesis-based) and inductive (pattern-based) reasoning
- This is a fundamental change from traditional analysis approaches.
- Inductive reasoning and exploratory data analysis provide a means to form or refine hypotheses and discover new analytic paths.
- Models of reality no longer need to be static.
- They are constantly tested, updated and improved until better models are found.

LOOKING BACKWARD AND FORWARD

FIRST THERE WAS BUSINESS INTELLIGENCE

Deductive Reasoning

Backward Looking

Slice and Dice Data

Warehoused and Siloed Data

Analyze the Past, Guess the Future

Creates Reports

Analytic Output

NOW WE'VE ADDED DATA SCIENCE

Inductive and Deductive Reasoning

Forward Looking

Interact with Data

Distributed, Real Time Data

Predict and Advise

Creates Data Products

Answer Questions and Create New Ones

Actionable Answer

Practical perspective

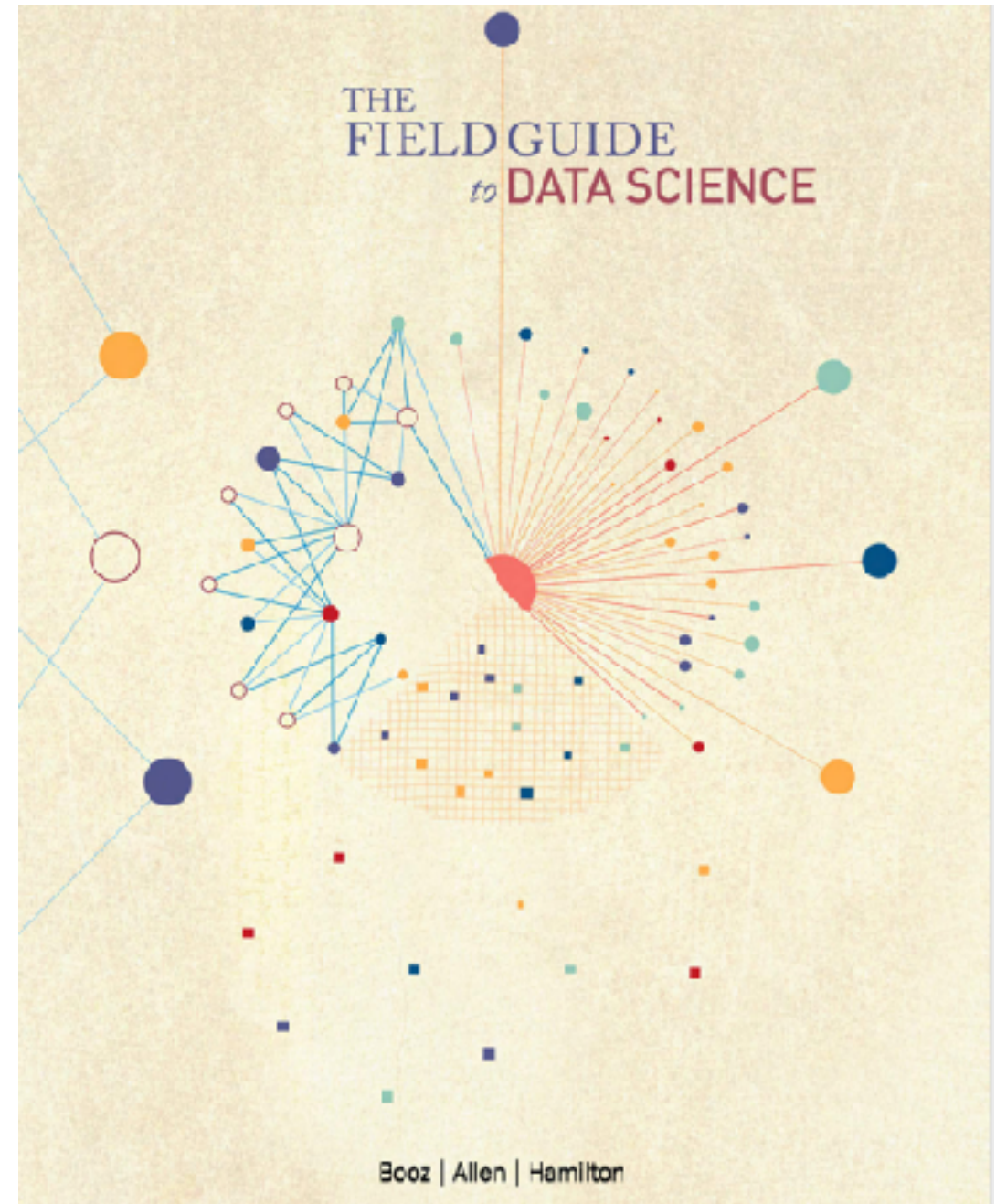
- Data Science is necessary for companies to stay with the pack and compete in the future.
- Data Science capabilities can be built over time.
- Data Science is a different kind of team sport.



Principles of Data Science

Data Science principles

- Be willing to fail.
- Fail often and learn **quickly**.
- Keep the goal in mind.
- Dedication and focus lead to success.



Learning quickly

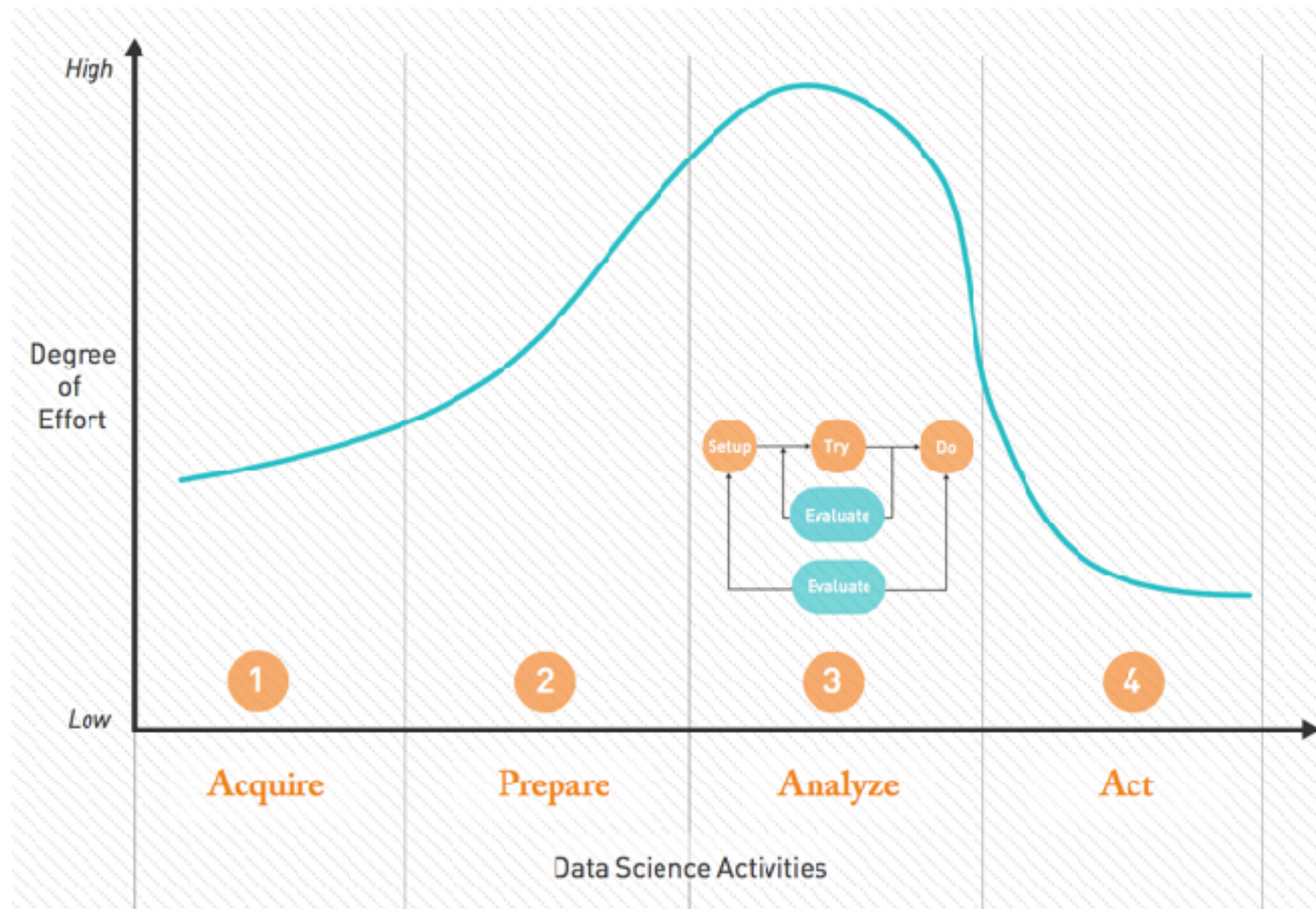
- Leon A. Gatys, Alexander S. Ecker, Matthias Bethge. "A Neural Algorithm of Artistic Style." arXiv:1508.06576. September 2015.



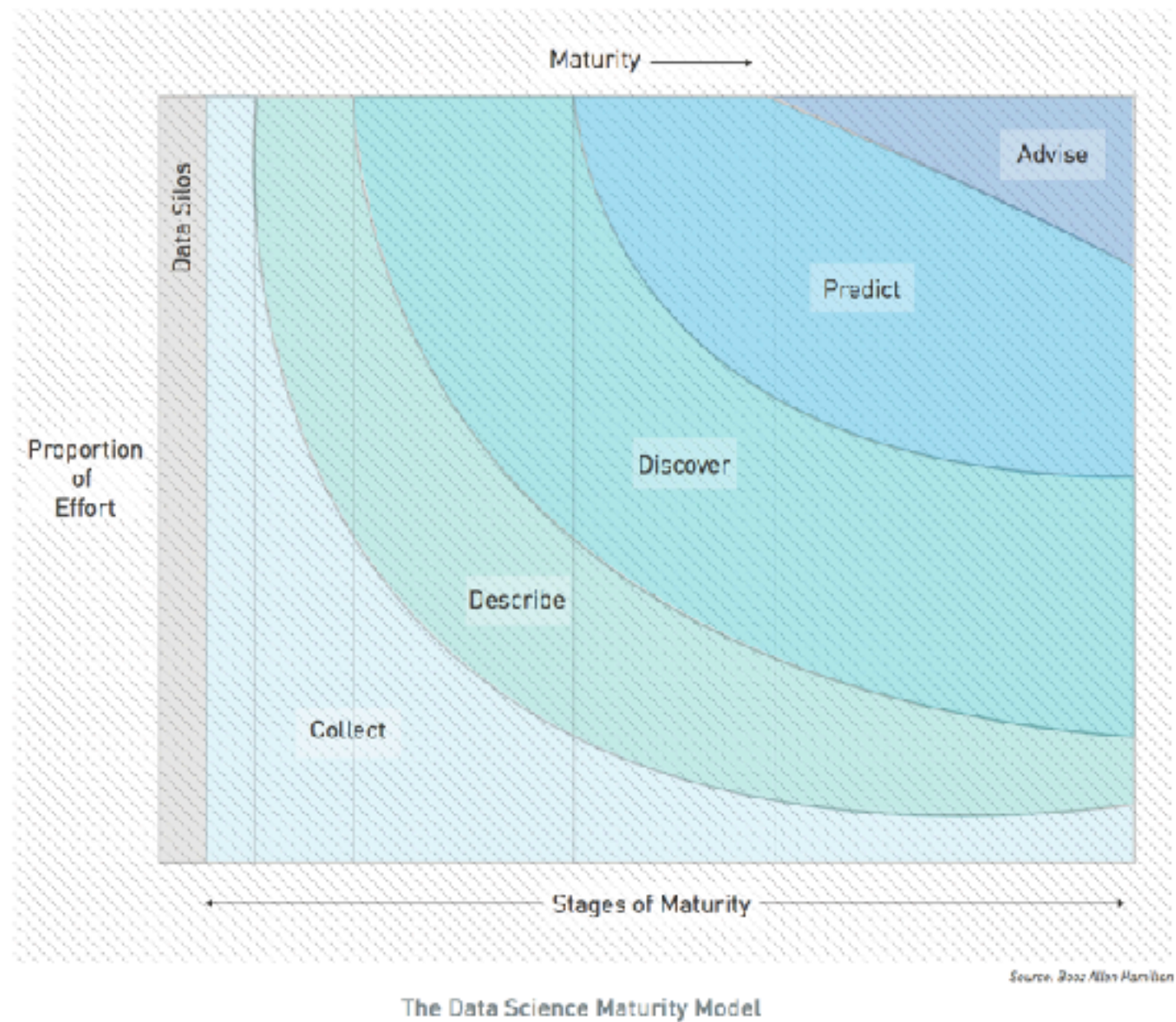
- Prisma and Convolutional Neural Networks: June 2016.



Practice of Data Science



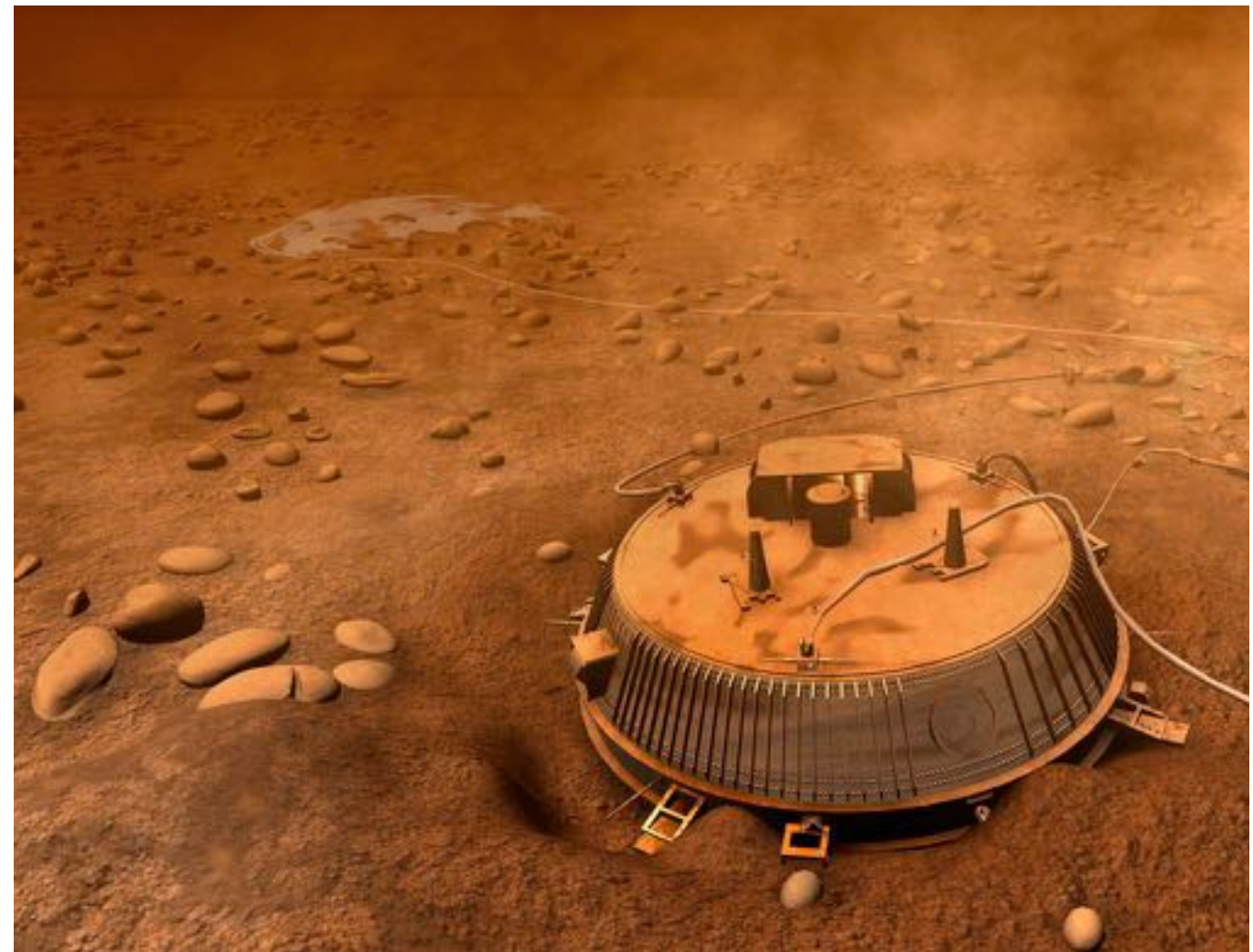
Data science workflow



Data science in organisations

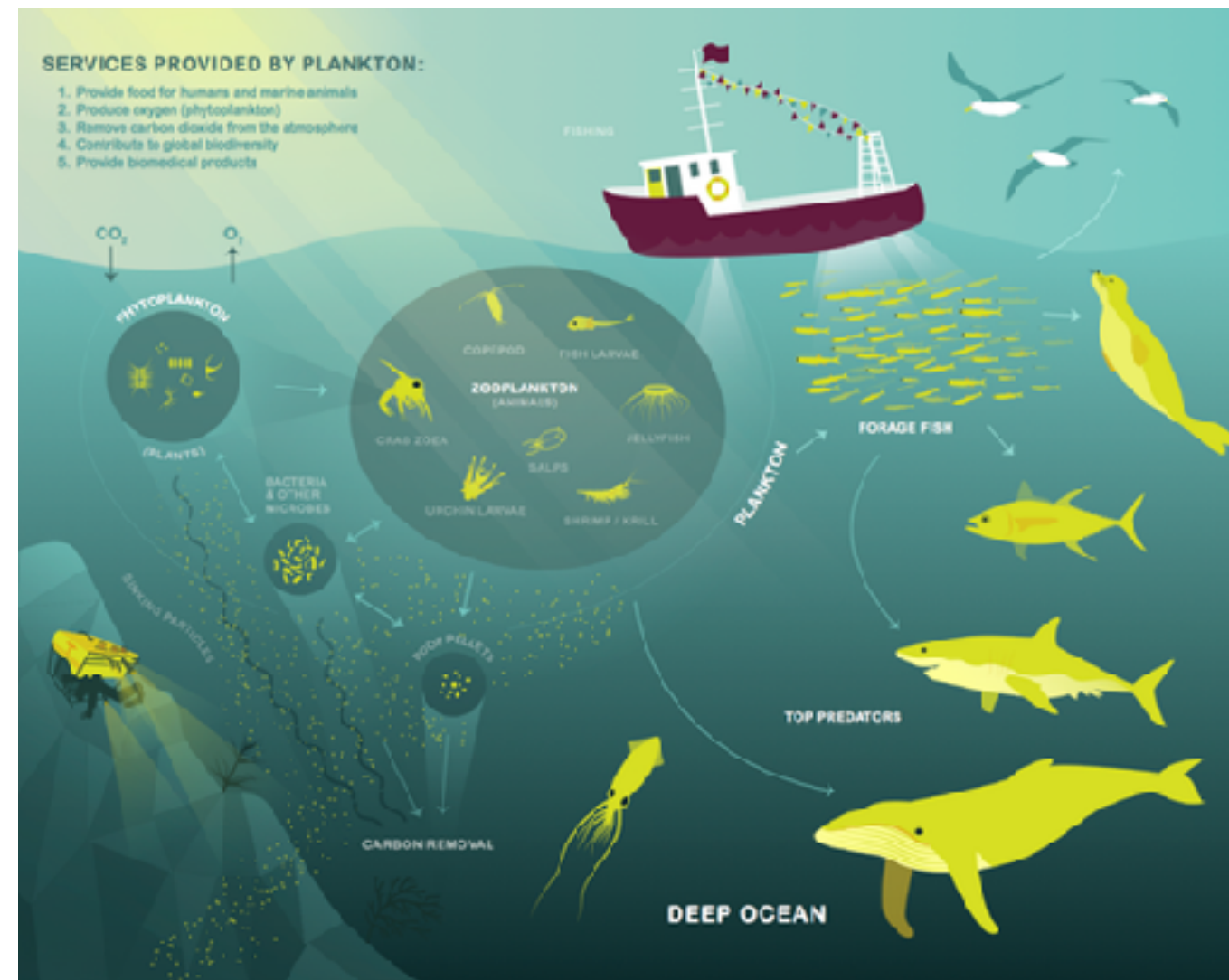
Unlocking value from data: Saturn and Cassini-Huygens mission

- "Bouncing on Titan: Motion of the Huygens Probe in the Seconds After Landing" by Stefan Schroeder et al (2 Feb 2017, ArXiv)
- Re-analyzing data from 2005. But also including (linking) data from a variety of other instruments. Get insights into surface structure.
- The way it bounced points to damp sand. Methane lakes.



Kaggle Competitions

- National Data Science Bowl - Predict ocean health, one plankton at a time.
- Passenger Screening Algorithm Challenge - Improve the accuracy of the US Department of Homeland Security's threat recognition algorithms
- Chicago Department of Public Health West Nile Virus Prediction - Predict West Nile virus in mosquitos across the city of Chicago.
- Redefining Cancer Treatment - Predict the effect of Genetic Variants to enable Personalized Medicine
- Genentech Flu Forecasting - Predict when, where and how strong the flue will be.
- Mercedes-Benz Greener Manufacturing - Can you cut the time a Mercedes-Benz spends on the test bench?



Governance



Cabinet Office

Data Science Ethical Framework

19 May 2016

Principles

1. Start with clear user need and public benefit
2. Use data and tools which have the minimum intrusion necessary
3. Create robust data science models
4. Be alert to public perceptions
5. Be as open and accountable as possible
6. Keep data secure

1 Start with clear user need and public benefit

- Data science offers huge opportunities to create evidence for policymaking, and make quicker and more accurate operational decisions.
- Being clear about the public benefit will help you justify the sensitivity of the data (principle 2) and the method that you want to use (principle 3).

2 Use data and tools which have the minimum intrusion necessary

- You should always use the minimum data necessary to achieve the public benefit.
- Sometimes you will need to use sensitive personal data.
- There are steps that you can take to safeguard people's privacy e.g. de-identifying or aggregating data to higher levels, querying against datasets or using synthetic data.

3 Create robust data science models

- Good machine learning models can analyse far larger amounts of data far more quickly and accurately than traditional methods.
- Think through the quality and representativeness of the data, flag if algorithms are using protected characteristics (e.g. ethnicity) to make decisions, and think through unintended consequences.
- Complex decisions may well need the wider knowledge of policy or operational experts.

4 Be alert to public perceptions

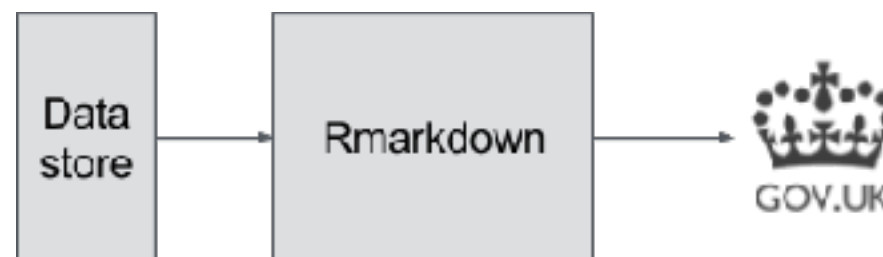
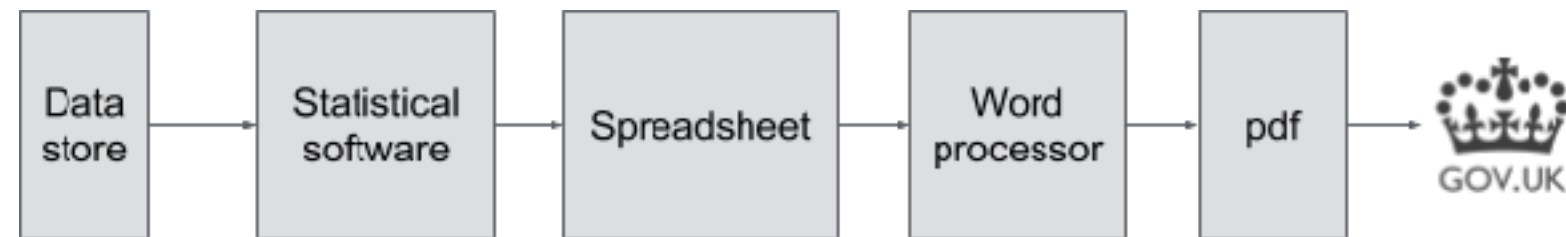
- The Data Protection Act requires you to have an understanding of how people would reasonably expect their personal data to be used.
- You need to be aware of shifting public perceptions.
- Social media data, commercial data and data scraped from the web allow us to understand more about the world, but come with different terms and conditions and levels of consent.

5 Be as open and accountable as possible

- Being open allows us to talk about the public benefit of data science.
- Be as open as you can about the tools, data and algorithms (unless doing so would jeopardise the aim, e.g. fraud).
- Provide explanations in plain English and give people recourse to decisions which they think are incorrectly made.
- Make sure your project has oversight and accountability built in throughout.

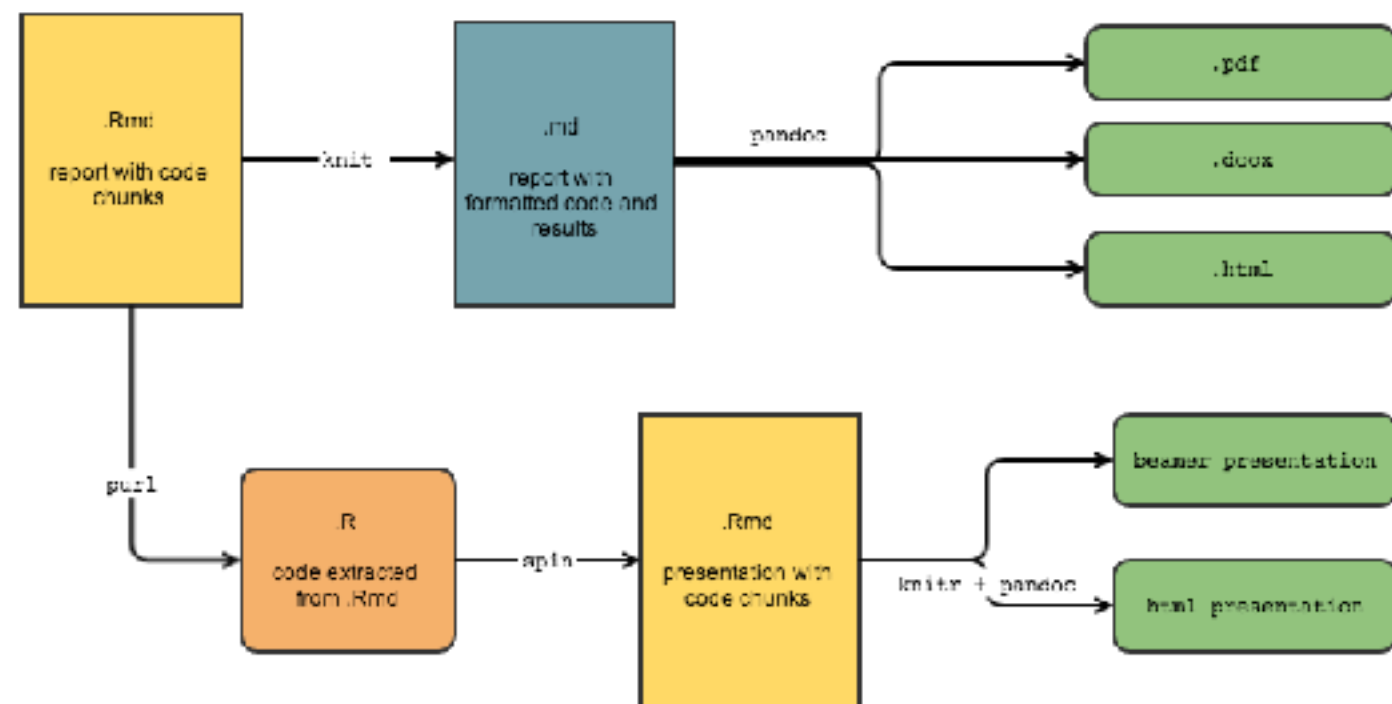
Reproducible Analytical Pipelines

- Open source rather than proprietary (R and Python)
- Version control, packaging code
- Procedural programming and unit testing
- Dependency management
- Data testing



Reproducible Analytical Pipelines and RMarkdown Reporting

- GDS RMarkdown reporting
- UK GDS on GitHub



6 Keep data secure

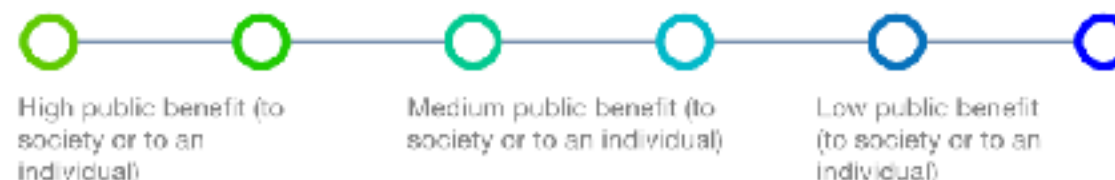
- We know that the public are justifiably concerned about their data being lost or stolen.
- Government has a statutory duty to protect the public's data and as such it is vital that appropriate security measures are in place.

Quick checklist

← - - - - Tick where you are on the scale - - - - →

1. Start with clear user need and public benefit

A. How does the department and public benefit?



2. Use data and tools which have the minimum intrusion necessary

B. How intrusive and identifiable is the data you are working with?



C. If identifying individuals, how widely are you searching personal data?



3. Create robust data science models

D. What is the quality of the data?



E. How automated are the decisions?



F. What is the risk that someone will suffer a negative unintended consequence as a result of the project?



4. Be alert to public perceptions

G. If personal data for operational purposes, how compatible was it with the reason collected?



H. Do the public agree with what you are doing?



5. Be as open and accountable as possible

I. How open can you be about the project?

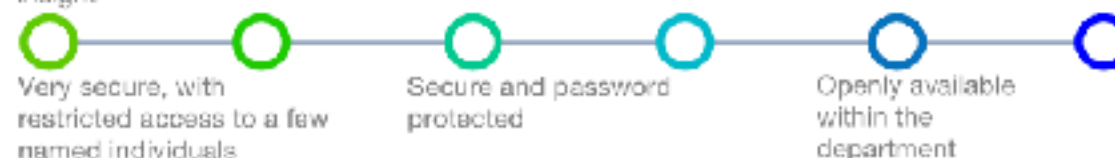


J. How much oversight and accountability is there throughout the project?



6. Keep data secure

K. How secure is your data?



*Not all may apply to your project

All fine?
Go forward!

Some issues?
Think carefully

Tricky issues?
Extreme care & oversight

Tricky issues?

Extreme care & oversight

Answering these questions will also act as your **Privacy Impact Assessment**

1. Start with clear user need and public benefit

How does the public benefit outweigh the risks to privacy and the risk that someone will suffer an unintended negative consequence? **(PIA Step 1)**

Brief description of the project, including data to be used, how it will be collected and deleted. **(PIA Step 2)**

2. Use data and tools which have the minimal intrusion necessary

What steps are you taking to maximise the benefit of the project outcome?

What steps are you taking to minimise risks to privacy? (for example using less intrusive data, aggregating data etc)?

3. Create robust data science models

What steps have you taken to make sure the insight is as accurate as possible and there are minimal unintended consequences? (for example thinking through quality of the data, human oversight, giving people recourse)

4. Be alert to public perceptions

How have you assessed what the public or stakeholders would think of the acceptability of the project? What have you done in addition to the above to address any concerns?

Risks (PIA Step 3) and mitigating steps (PIA Step 4)

5. Be as open and accountable as possible

How are you telling people about the project and how you are managing the risks?

Who has signed this off within your organisation? Who will make sure the steps are taken and how? **PIA Step 5**

6. Keep data secure

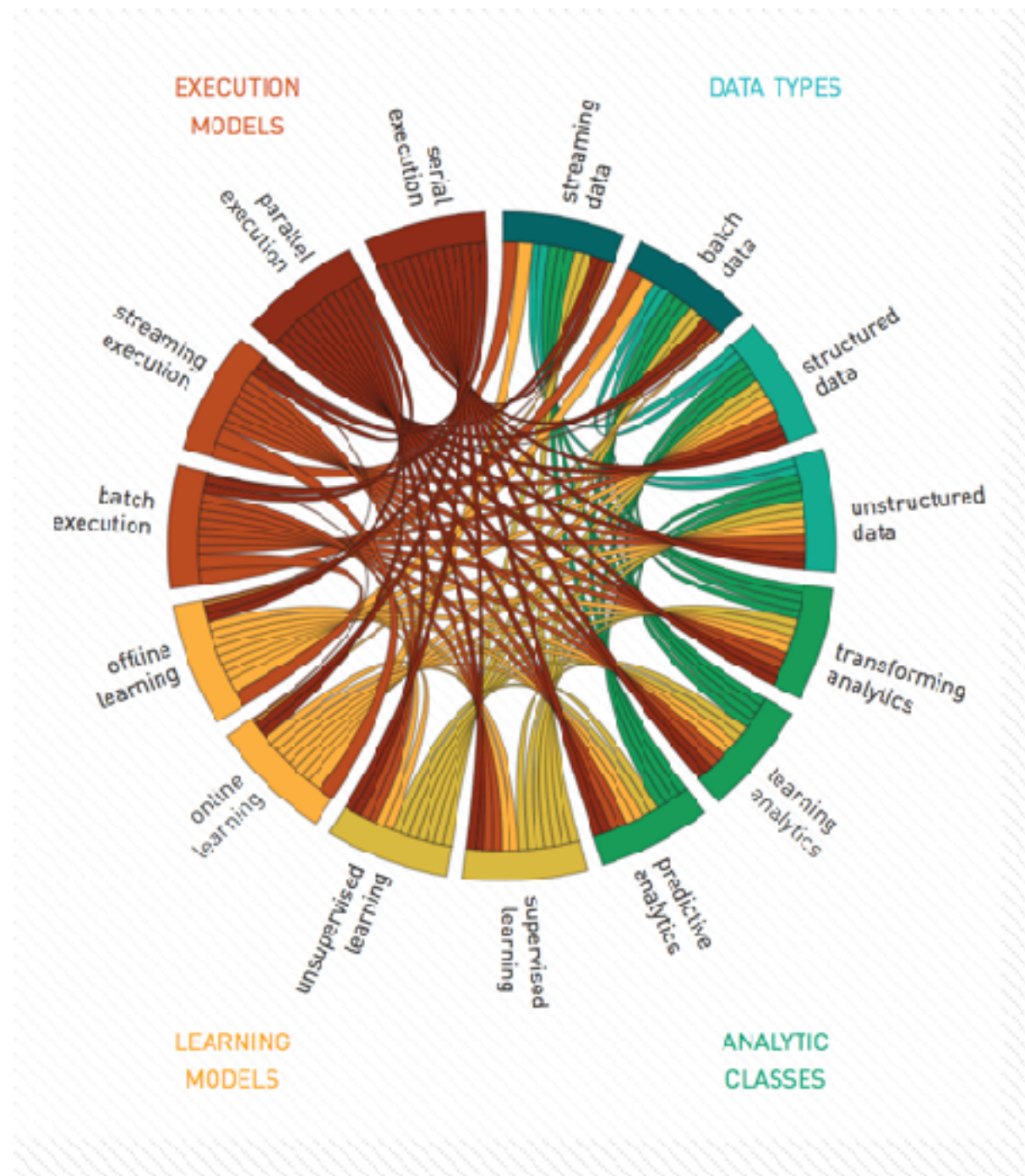
What steps are you taking to keep the data secure?

Case study

- Discussion of the case in the media
- ICO ruling
- Checking against Data Science Ethical Framework



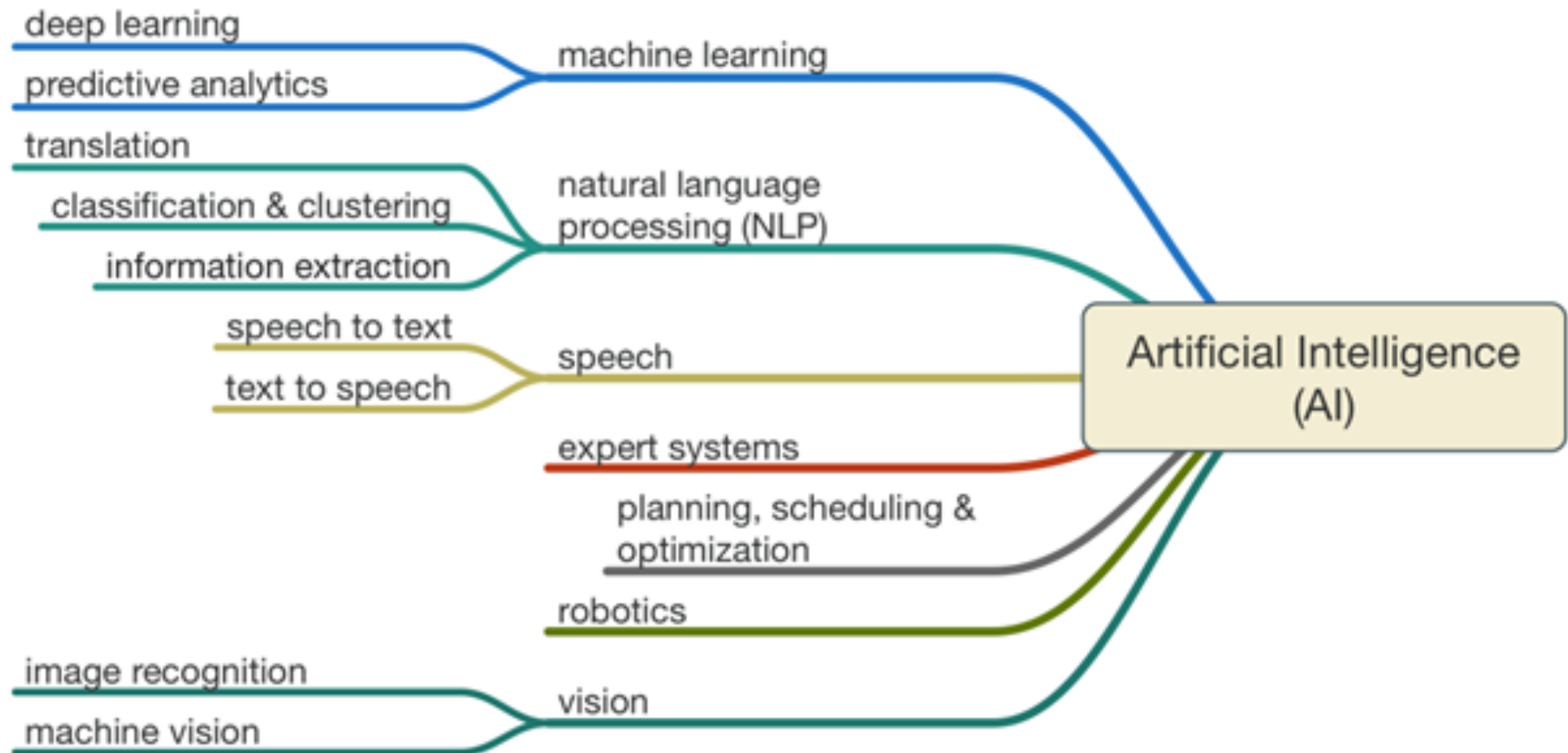
Analytics



Interconnection Among the Component Parts of Data Science

Source: Booz Allen Hamilton

Analytics components



Neota Logic

Main models and approaches