# Data Science in an hour

**Jarlath Quinn – Analytics Consultant**

Just waiting for all attendees to join…

# Data Science in an hour

**Jarlath Quinn – Analytics Consultant**

A SELECT INTERNATIONAL COMPANY

# FAQ's

- Is this session being recorded? Yes

- Can I get a copy of the slides? Yes, we'll email links to download materials after the session has ended.

- Can we arrange a re-run for colleagues? Yes, just ask us.

- How can I ask questions? All lines are muted so please use the chat panel – if we run out of time we will follow up with you.

- Gold accredited partner to IBM, Predictive Solutions and DataRobot specialising in advanced analytics & big data technologies

- Work with open-source technologies (R, Python, Spark etc.)

- Team each has 15 to 30 years of experience working in the advanced and predictive analytics industry

- Deep experience of applied advanced analytics applications across sectors
  - Retail
  - Healthcare/Pharma
  - Finance/Insurance
  - Media/Telecoms
  - Utilities
  - FMCG
  - Charity/Housing/Government



A SELECT INTERNATIONAL COMPANY

# How did we get here?

# Statistical Analysis to AI



DATA SCIENCE TIMELINE v. 2.0

Illustration by Héizel Vázquez

@faviovaz
@heizelvazquez

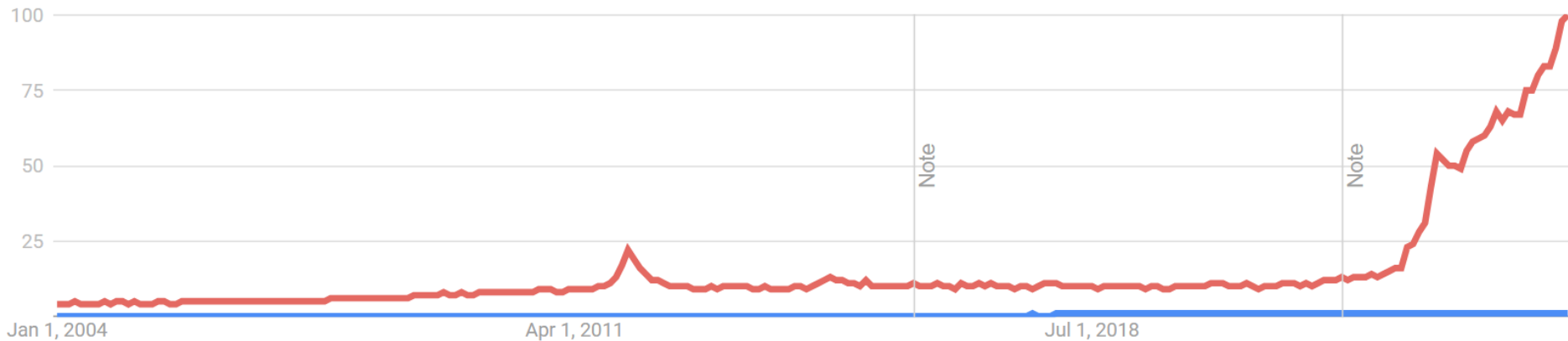Credit: https://medium.com/towards-data-science/the-roots-of-data-science-77c71115229

A SELECT INTERNATIONAL COMPANY

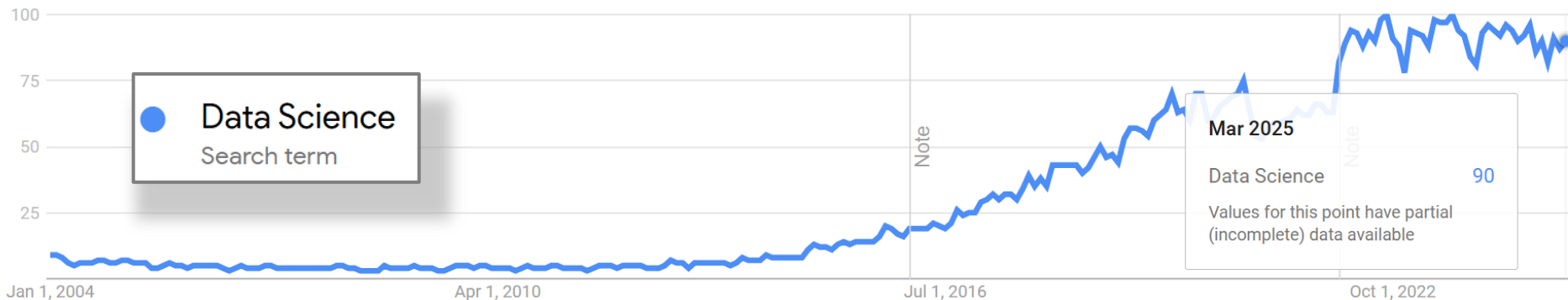# The new frontiers of data analysis

# Not forgetting of course…the elephant in the room



Data Science ● AI

SMART VISION
Europe

The term **Data Science** was first proposed by Peter Naur in 1974 as an alternative name for computer science.

But it wasn't until 2008 that Patil and Hammerbacher popularized the term **Data Scientist** to describe professionals who combine programming skills with statistical knowledge to extract insights from data.



Data Science
Search term

Mar 2025
Data Science                90
Values for this point have partial (incomplete) data available

100
75
50
25

Jan 1, 2004          Apr 1, 2010          Jul 1, 2016          Oct 1, 2022

**Data Science** is using data, statistics and computing power to answer questions and make better decisions

SMARTVISION
Europe

A SELECT INTERNATIONAL COMPANY

# Disciplines



Approximate Academic Backgrounds of Data Scientists

- Other (Social Sci & Humanities) — 12%
- Business & Economics — 10%
- Other Engineering — 10%
- Physics & Natural Sciences — 18%
- Computer Science & Software Eng — 20%
- Maths / Stats / Econometrics / OR / ML — 30%

**Skills**



Approximate Emphasis of Skill Sets in Data Science

# Tools

## Approximate Usage of Tool Categories in Data Science



- AI / LLM tools & APIs (7%)
- Cloud computing platforms (8%)
- Data management / DBs & warehouses (13%)
- Programming environments (12%)
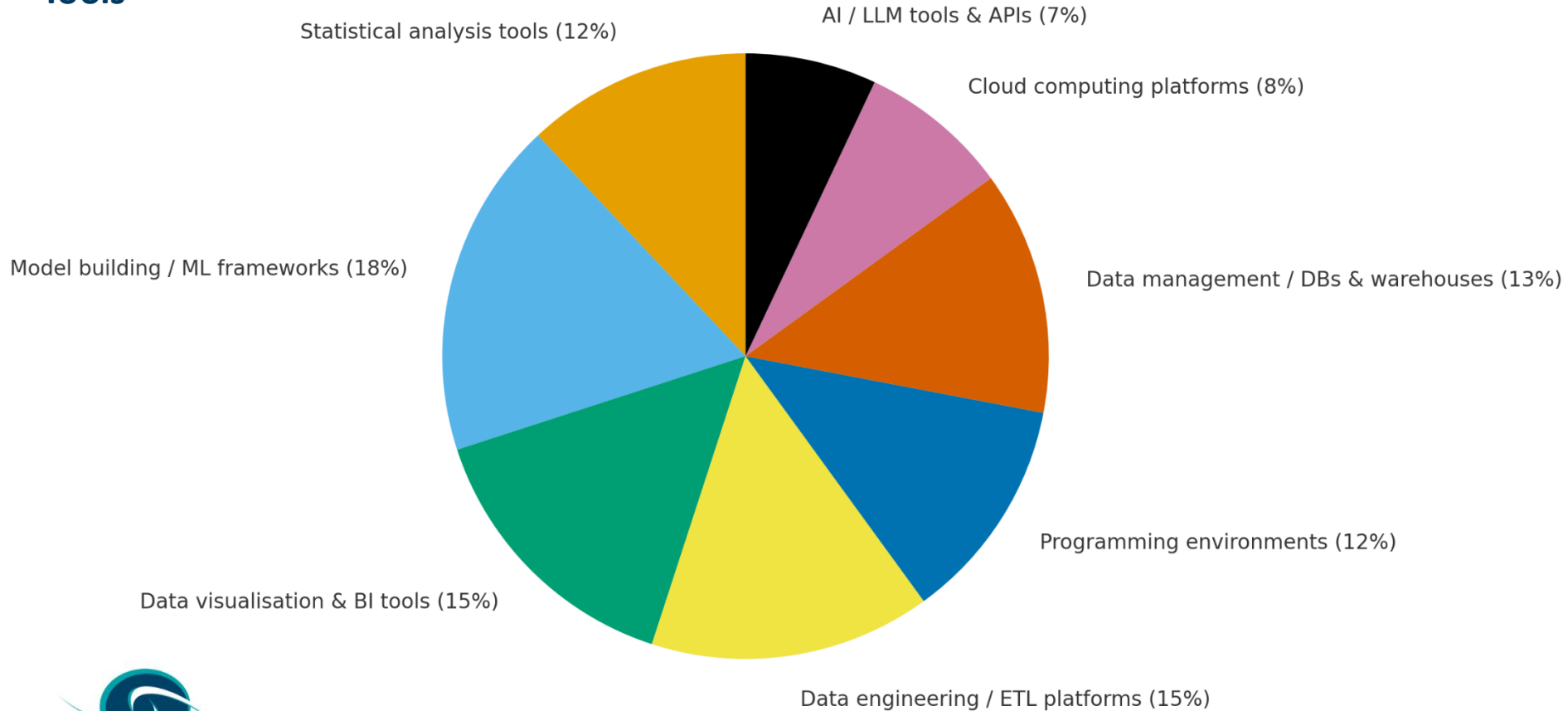- Data engineering / ETL platforms (15%)
- Data visualisation & BI tools (15%)
- Model building / ML frameworks (18%)
- Statistical analysis tools (12%)

# *What does Data Science actually produce?*

Although Data Science is sometimes about driving insight...

...it's the *new data* it creates which is really important.

# What do these new data represent?

- **Likelihoods**
    - recommend to a friend / complete a tv series / renew a subscription / click an offer / return to a store / make an insurance claim / choose a university / require maintenance / need a biopsy / make a complaint / fail a warranty / complete a course / return to hospital / fall into arrears / leave employment / defect to a competitor / commit fraud / show up for a flight / repay a loan / cause an accident / prevent infection / report a crime / vote for a party

# What do these new data represent?

- **Estimates & Forecasts**

  – Student scores / regional sales / time to completion / blood pressure readings / pollution levels / website hits / survival times / growth rates / museum visits / medical costs / fuel consumption / crop yields / traffic volumes / causality patients / monthly expenditures / pupil numbers / power consumption / maintenance jobs / supply interruptions / flooding events / passenger volumes / property prices / infection rates / tickets sold

# What do these new data represent?

- **Categories & Recommendations**

  – Customer segments / fault causes / medical diagnoses / tumour classes / replacement parts / treatment risk groups / preferred movie genres / political affiliations / fashion preferences / mobile phone plans / satisfaction levels / recommended crop types / product assortments / suggested drug regimes / targeted advert recommendations / content filters / document categories / customer sentiments / image classifications / speech-emotion classes

# Typical Data Science Applications

- Customised Offer Creation
- Subscriber Retention
- Drug Performance Prediction
- Patient Outcome/Readmission
- Predictive Maintenance/Issue Forecasting
- Fraud Detection
- Loyalty Modelling
- Path to Purchase

- Capacity planning & scheduling
- Anomaly Detection
- Association Analysis
- Forecasting
- Chatbots and virtual assistants
- Text Classification
- Optimisation Engines

**It's important to understand that depending on the circumstances, some of these applications may be driven by old statistical methods whilst others rely on cutting edge AI algorithms**

SMARTVISION
Europe

# So where does AI like Chat GPT fit into Data Science?

- **Code Generation and Debugging**: Writing code snippets and fixing errors

- **Feature Engineering**: Helping to prepare the data for modelling

- **Data Transformation**: Turning unstructured data like text, images and audio into structured data like numbers

- **Exploratory Data Analysis (EDA):** Carrying out basic data analysis and generating summary reports and visualisations

- **Documentation and Explanation:** making projects more transparent and understandable for other team members

# Building a Data Science Model

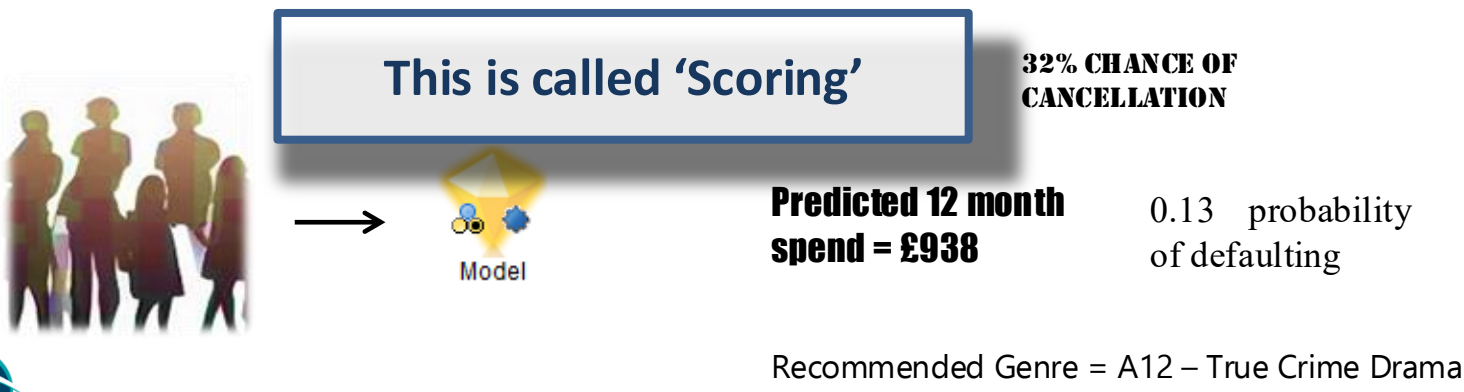# At the heart of a Data Science application is a model

- Typically uses historical data from many people/incidents/assets

- Age, Gender, Spending, Region, Tenure, Usage etc.

- With a known outcome/result

- Responded, upgraded, defaulted, recommended, cancelled, donated, failed, renewed etc.

- To create an accurate, usable model

Model

**This is called 'Training'**

# At the heart of a Data Science application is a model

- We can take new data from new individuals or incidents…

- Age, Gender, Spending, Region, Tenure, Usage etc.

- Using a model based on the same information…

- Generate likelihood scores, estimates and classifications

- In other words,…..predictions

**This is called 'Scoring'**

**32% CHANCE OF CANCELLATION**

→

Model

**Predicted 12 month spend = £938**

0.13 probability of defaulting

Recommended Genre = A12 – True Crime Drama

# At the heart of a Data Science application is a model

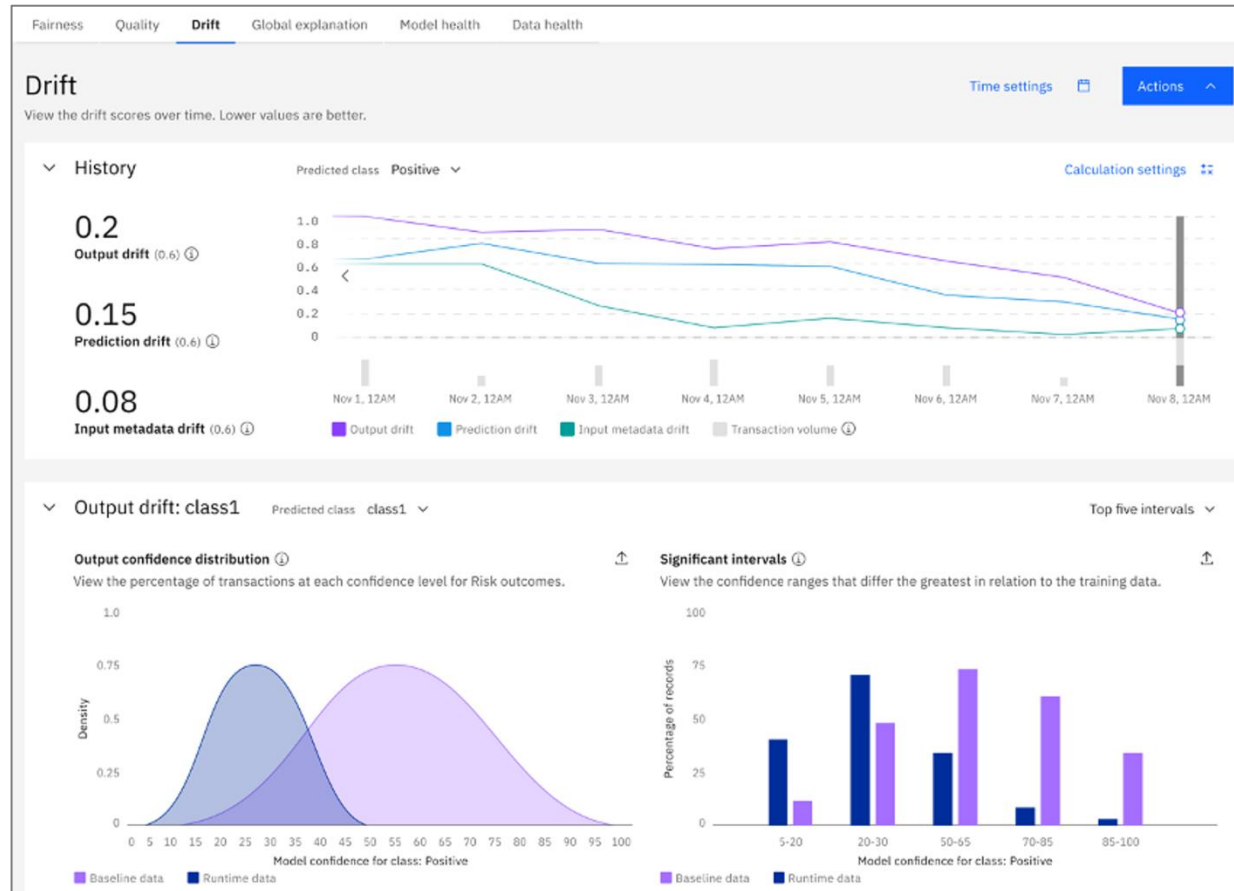- We can then send the model scores to different platforms to drive better outcomes



This is called 'Deployment'

Model

SMARTVISION
Europe

# However, a Model is *not* an Application...

# Until it is used in the real world to drive outcomes

$$\text{maximize } f(c_1 \ldots c_n) = \sum_{i=1}^{n} c_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i c_i (\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)) y_j c_j$$

$$= \sum_{i=1}^{n} c_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i c_i k(\mathbf{x}_i, \mathbf{x}_j) y_j c_j$$

$$\text{subject to } \sum_{i=1}^{n} c_i y_i = 0, \text{ and } 0 \leq c_i \leq \frac{1}{2n\lambda} \text{ for all } i.$$

$\neq$

Recommendations:
- Check pupil dilation
- Eye tracking function
- Assess nausea level
- Check reflex reaction

# And applications usually require Governance

# And applications usually require Governance

- **Prevents harmful or biased outcomes:** Rules for data access, model use, and fairness checks, reducing the risk of discrimination, reputational damage.

- **Keeps models accurate over time (drift monitoring):** Monitoring picks up when data or behaviour changes (model drift), so you can retrain, recalibrate, or retire models before performance falls off a cliff.

- **Supports compliance and auditing:** Documentation and audit trails for models (who built it, what data, which checks) to help satisfy regulators, internal audit, and legal – especially in finance, healthcare, HR, and credit.

- **Improves reproducibility and handover:** Good practice standards for version controls, security, and model lifecycle  ensure others can reproduce results, debug issues, and safely update or extend a model.

- **Builds trust with stakeholders:** Decision makers are reassured there's a formal process for reviewing, approving, and monitoring models in production, they're more willing to rely on data science for real decisions, not just "nice dashboards."

# *What are the real-world challenges with Data Science?*

If you build it,
*they will come*

# All the gear but no idea

- Even big companies make the mistake of thinking that Data Science/AI is all about having the right resources:

  - A new data science team

  - A cloud-based AI platform

  - Sophisticated data storage/process architecture

# Do they have a use for Data Science or AI?

- A regular complaint among newly-hired but highly-qualified Data Scientists and AI specialists is that they find their roles consist of fairly basic analytical tasks such as running SQL queries or building dashboards

- Some companies may use the term "data scientist" as a buzzword to attract talent, without a clear understanding of what the role entails

Hired as a Data Scientist, not doing Data Science work. - Reddit

2 Jun 2021 — Hired as a Data Scientist, not doing Data Science work. : r/datascience.

Reddit · r/datascience

Big problem with companies now is they hire data scientist for task ...

31 Aug 2022 — Big problem with companies now is they hire data scientist for task that don't require data...

Reddit

Current "Data Science" job is unfulfilling and demotivating. I want to ...

12 Dec 2021 — It feels awful. Lately, I don't even know if I want to be in data science anymore because this...
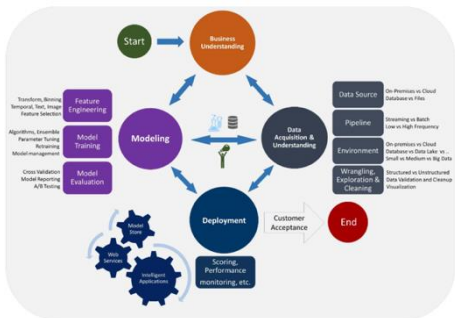
Reddit

**SMARTVISION**
Europe

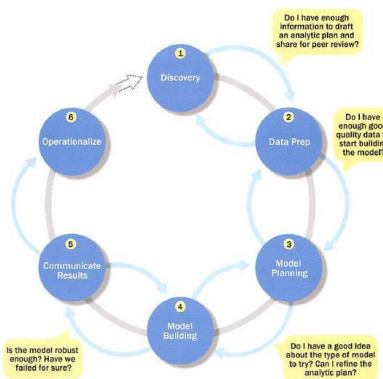# What are the biggest challenges in implementing Data Science?

- The time and effort taken to consolidate, blend and prepare data so it can used effectively

- Coordination and communication across business units

- Matching the capabilities of data science to the needs of the organisation i.e. creating valuable applications

- Measuring the value of the application

- Creating a feedback cycle to manage things operationally

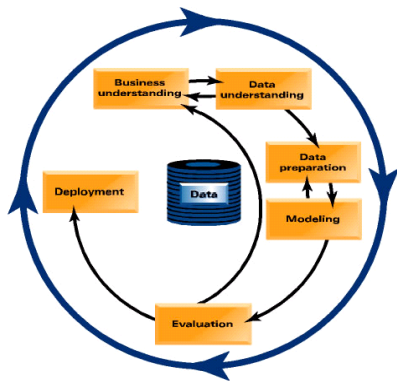# It's helpful to know there are several methodologies dedicated Data Science

- Microsoft's Team Data Science Process (TDSP)



- EMC's Data Analytics Lifecycle



IBM's Analytics Solution Unified Method (ASUM-DM)





- Cross-Industry Standard Process for Data Mining (CRISP-DM)

# And they illustrate that it's not all just building models



Business understanding → Data understanding → Data preparation → Modeling → Evaluation → Deployment

Data

Often Data Scientists may spend 50 % to 70% of their time just wrangling and preparing the data when working on a new project

Important questions for intrepid adventurers

A SELECT INTERNATIONAL COMPANY

- What does 'good' look like?

- What will you do differently?

- How will you know it worked?

- Does everyone agree or understand?

- What methodology will you use?

SMARTVISION
Europe

# Think of a methodology as a route map to successful deployment

- CRISP-DM: Cross-Industry Standard Process for Data Mining
- Each application can be developed and progressed through a series of key phases
- http://crisp-dm.eu/

# Download our e-book for free



THE INSIDERS' GUIDE TO
**PREDICTIVE
ANALYTICS**

**Jarlath Quinn**

Smart Vision Europe

The insider's guide to predictive analytics

£0.00

- | 1 | + | **Add to basket**

Category: books

A SELECT INTERNATIONAL COMPANY

# ChatSPSS

**CONVERSATIONS**

▼ How can I remove unwanted LOG output?

    How can I remove unwanted LOG output

    Whats the difference between ANOVA and ANCOVA

To remove unwanted LOG output in IBM SPSS Statistics, you can hide the log messages in the Viewer. Here's how to do it:

1. From the menus, choose:

   `Edit > Options...`

2. Click the **Viewer** tab.
3. Double-click the **Log** icon to hide log messages in the output.
4. Click **OK** to save your changes and close the dialog box.

This will hide the log output from the Viewer, helping to declutter your output display.

> Whats the difference between ANOVA and ANCOVA?

## Difference Between ANOVA and ANCOVA

ANOVA (Analysis of Variance) and ANCOVA (Analysis of Covariance) are both statistical methods used to compare the means of different groups. However, they differ in their approach and application:

- **ANOVA**: This method is used to compare the means of three or more groups to see if there are any statistically significant differences between them. It does not take into account any other variables that might influence the dependent variable.

- **ANCOVA**: This method is an extension of ANOVA that includes one or more continuous independent variables, known as covariates. ANCOVA adjusts the means of the groups based on the covariates, which can help reduce the error variance and increase the power of the test.

In summary, ANCOVA is used when you want to control for the effects of one or more continuous variables (covariates) while comparing the means of different groups, whereas ANOVA is used when you only want to compare the means without considering any other variables.

For more detailed information, please consult external sources or the official SPSS documentation.

**DOCUMENT REFERENCES**

IBM_SPSS_Statistics_Core_System_User_G...

IBM_SPSS_Statistics_Brief_Guide.pdf

IBM_SPSS_Statistics_Base.pdf

IBM_SPSS_Missing_Values.pdf

**Training**

Training notes

In person training

Online self-paced training courses

**Consulting**

SPSS Boost

One-to-one guided consultancy

**Support**

Technical support

Frequently asked questions

ChatSPSS – AI-Powered SPSS Tutor

# Meet ChatSPSS – Your AI-Powered SPSS Tutor

## Get instant, accurate answers to your SPSS and statistics questions – anytime you need them.

Smart Vision Europe is proud to introduce **ChatSPSS**, the first AI-enabled tutor designed specifically for SPSS users. Whether you're new to SPSS or an experienced analyst, ChatSPSS helps you get the very most out of your software, saving you time and boosting your confidence with data analysis.

## What can ChatSPSS do for you?

✅ Guide you through loading and preparing data in SPSS

✅ Show you how to summarise results and run statistical tests

✅ Explain outputs from regression models and other advanced techniques

✅ Answer almost any SPSS or statistics-related question, on demand

# Smart Vision Europe: Services and Expertise

We have decades of experience providing guidance, training and consultancy in the delivery of effective data science initiatives.

Contact us:

+44 (0)207 786 3568
info@sv-europe.com
Twitter: @sveurope
Follow us on Linked In
Sign up for our Newsletter

# Thank you