# *Workshop's goal*

The workshop will guide you through the process of completing an **end-to-end Data Science project**.

We will start with a **problem statement** and end with a **deployed product** that our client will be able to use.

We will utilize and connect various technologies, packages and programming paradigms to produce a functional product for our (fictional) client.

# *What to expect*

## *Not this*

- Course about the different technologies
- Deep development of any of the steps
- Information about specific markets or industries

## *But that*

- Simplified end-to-end life cycle of an AI solution development
- Connecting all the different tech and analytics pieces together
- Reflections on real commercial operations and projects & the associated best practices

#ACCELERATE

# *Workshop overview:*

## Session 1
### Prep & Analytics
## 12.09.2021

*Start with the business problem, set the foundation up, find data source, preprocess Start the descriptive analytics pipeline )*

## Session 2
### Machine learning
## 19.09.2021

*Implemented analytics pipeline, Build and evaluate prediction model(s), use MIflow to keep track of the various experiments*

## Session 3
### Deployment in Prod
## 26.09.2021

*Create prediction functions and production class, develop an API, create a dashboard that the user will access and call the API*

## *What you will do:*

- **During the sessions:** You will get tasks to be done

- **After the sessions:**
  - You will complete the whole covered phases
  - Dig deeper into the various technologies discussed

  *i.e.: No Spoon-feeding :-)*

#ACCELERATE

# & let's get started

# Problem statement

**Our** *(fictional)* **client is an IT educational institute. They have reached out to us has reach out with the following:**

"IT jobs and technologies keep evolving quickly. This makes our field to be one of the most interesting out there. But on the other hand, such fast development confuses our students. They do not know which skills they need to learn for which job.

"Do I need to learn C++ to be a Data Scientist?" "Do DevOps and System admins use the same technologies?" "I really like JavaScript; can I use it in Data Analytics?" Those are some of the questions that our students ask.
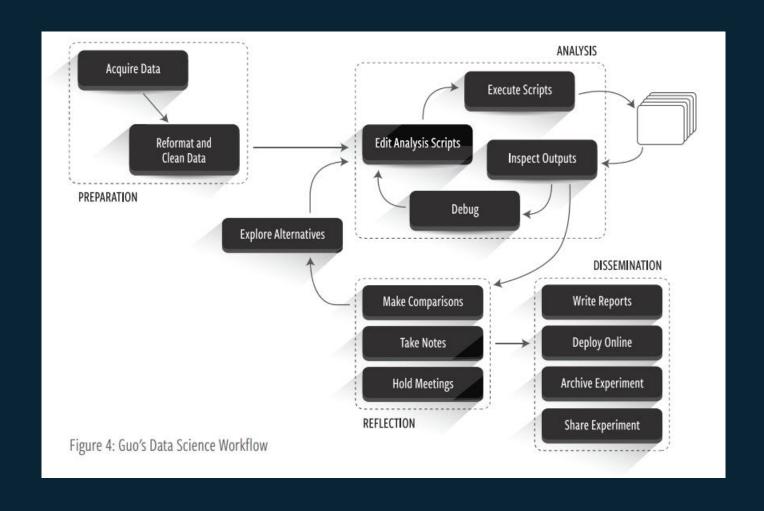
Could you please develop a data-driven solution for our students to answer such questions? They mostly want to understand the relationships between the jobs and the technologies.

# Data Science Workflow

*https://www.researchgate.net/figure/Guos-Data-Science-Workflow_fig2_319317441*



Figure 4: Guo's Data Science Workflow

# 1. Business Problem

# It's your turn:
## What is your Business case?

You are asking a commercial business to invest in a new project. You need to prove that your work will have a positive financial impact.
**How will you prove this? What are the KPIs that you will positively impact?**

# *Business case*

You are asking a commercial business to invest in a new project. You need to prove that your work will have a positive financial impact.
**How will you prove this? What are the KPIs that you will positively impact?**

1. **Higher enrollment rate due to the higher certainty**

2. **Decrease in drop-out rate**

3. **Time saved for the academic advisors**

*Learn more: https://www.youtube.com/watch?v=zQ5WqAz3myo*
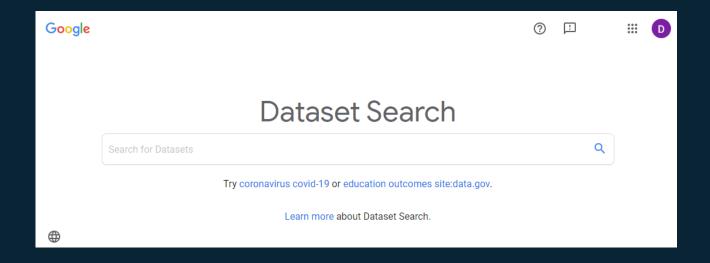
# 2. Data

# Data source

## Where to start?

https://datasetsearch.research.google.com/



## Be careful:

- Be thorough with the quality checks
- Make sure that your data will be updated on a regular base

#ACCELERATE

# *Data source*

## *Chosen: Stack Overflow developers survey*
https://insights.stackoverflow.com/survey/2020



2020
Developer
Survey

In February 2020 nearly 65,000 developers told us how they learn and level up, which tools they're using, and what they want.

Read the overview →    Methodology →

# 3. Foundations

# 1. Legal and data privacy check

## Global:
*https://www.privacyaffairs.com/gdpr-fines/*



## Local:
*https://www.privacylaws.com/media/3263/egypt-data-protection-law-151-of-2020.pdf*



#ACCELERATE

# 2. *How to structure your project*

https://drivendata.github.io/cookiecutter-data-science/

## Directory structure

```
├── LICENSE
├── Makefile           <- Makefile with commands like `make data` or `make train`
├── README.md          <- The top-level README for developers using this project.
├── data
│   ├── external       <- Data from third party sources.
│   ├── interim        <- Intermediate data that has been transformed.
│   ├── processed      <- The final, canonical data sets for modeling.
│   └── raw            <- The original, immutable data dump.
│
├── docs               <- A default Sphinx project; see sphinx-doc.org for details
│
├── models             <- Trained and serialized models, model predictions, or model summaries
│
├── notebooks          <- Jupyter notebooks. Naming convention is a number (for ordering),
│                         the creator's initials, and a short `-` delimited description, e.g.
│                         `1.0-jqp-initial-data-exploration`.
│
├── references         <- Data dictionaries, manuals, and all other explanatory materials.
│
├── reports            <- Generated analysis as HTML, PDF, LaTeX, etc.
│   └── figures        <- Generated graphics and figures to be used in reporting
│
├── requirements.txt   <- The requirements file for reproducing the analysis environment, e.g.
│                         generated with `pip freeze > requirements.txt`
│
```
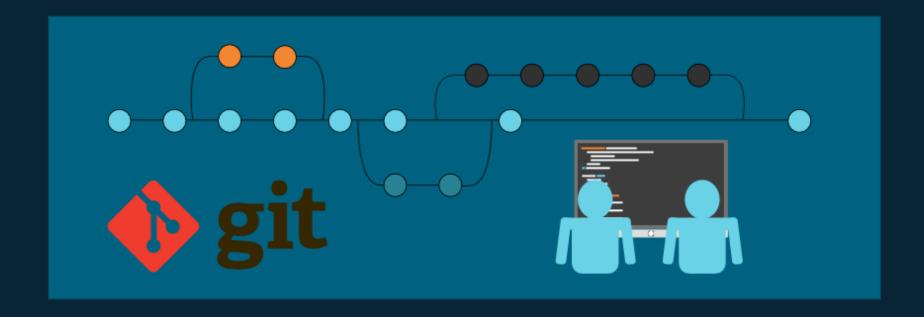
# 3. Your Git repo

https://developerhowto.com/2018/10/12/git-for-beginners/

# *It's your turn:*

- **Create your project directory using the cookie cutter**

- **Track your project in a new GitHub repo**

- **Download & save your data**

# 4. Preprocessing

# *Preprocessing at first glance*

...

1. **String values in years need to be replaced**

2. **Multiple values separated be `;` need to be splitted**

...

→

- *Prioritize task*
- *Create tickets in Jira*
- *Team members pick the tickets and solve them*

# *Jira - Kanban*

# *It's your turn :*

- **Preprocess your raw data and export it to a pickle file**

- **Push your work to your repo**

# 5. Descriptive Analytics

**It's your turn:**
**What are the descriptive questions that you will answer ?**

Think about what you want to do before you start doing it. Keep the original goal in mind

# My analytics question

**General:**
- Total number of answers
- Geographical distributions
- Missing answers

**Skills:**
- Frequency of each skill
- How are the skills correlated with each others

**Jobs:**
- Frequency of each job
- How are the jobs correlated with each others

**Relation:**
- How are the skills correlated with the jobs
- What is the specificity of each skill to a job

#ACCELERATE

# Levels of descriptive analytics

1. Stats or summary tables

2. Visualizations

3. Unsupervised learning (e.g. clustering)

# Assignment:

**Complete the setup, preprocessing & descriptive analytics phase**

# *Wrap Up*

# *Wrap-up:*
# *Today you have learned about*

- *Build a business case*
- *Find suitable data sources*
- *Verify legal rights*
- *Cookie-cutting your directory structure*
- *Track your project via Git*
- *Explore and preprocess data*
- *Collaborate with your team using Jira*
- *Planning framework for your descriptive analytics*

#ACCELERATE