

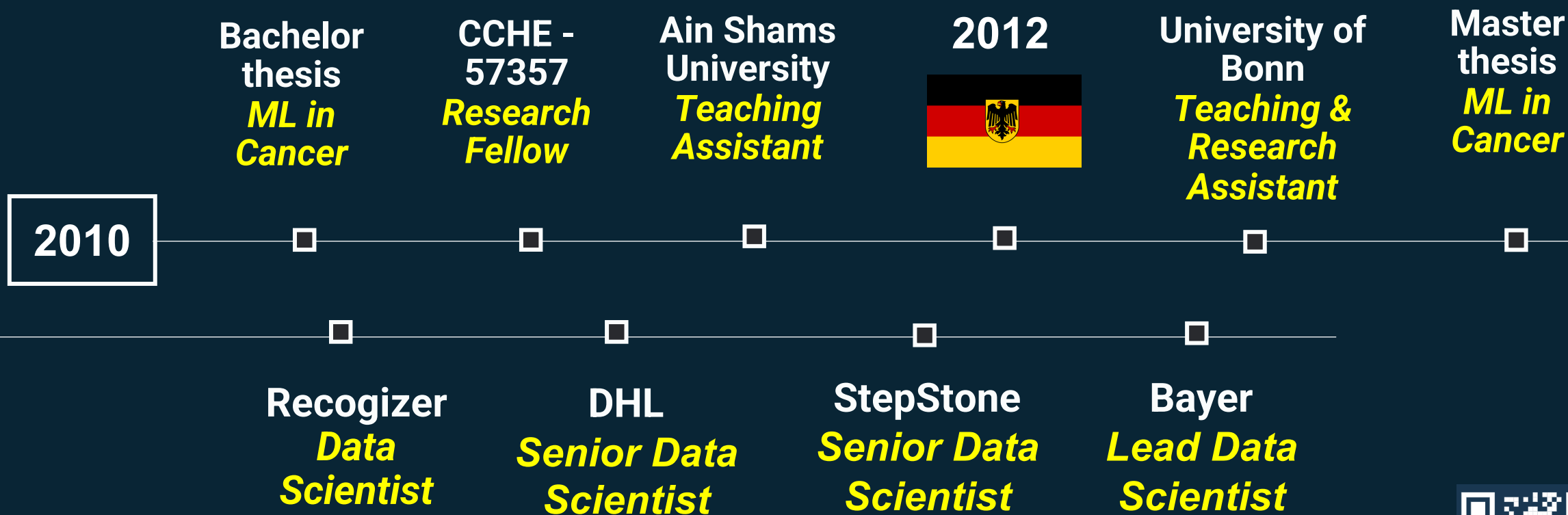


An End-to-End Data Science Project

Deena Gergis
Lead Data Scientist @ Bayer



Who am I





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

- Various levels of difficulties
- 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



Workshop overview:

Session 1 Preparation 10.04.2022

Start with the business problem, find data source, preprocess data & start the descriptive analytics pipeline

Session 2 Analytics 17.04.2022

Analyze and understand your data. Gain insights and prepare for the predictive modeling

Session 3 Machine learning x.05.2022

Build and evaluate prediction model(s), use Mlflow to keep track of the various experiments

Session 4 Production x.05.2022

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:

- **Form a team of 3 members**
- **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



& 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.

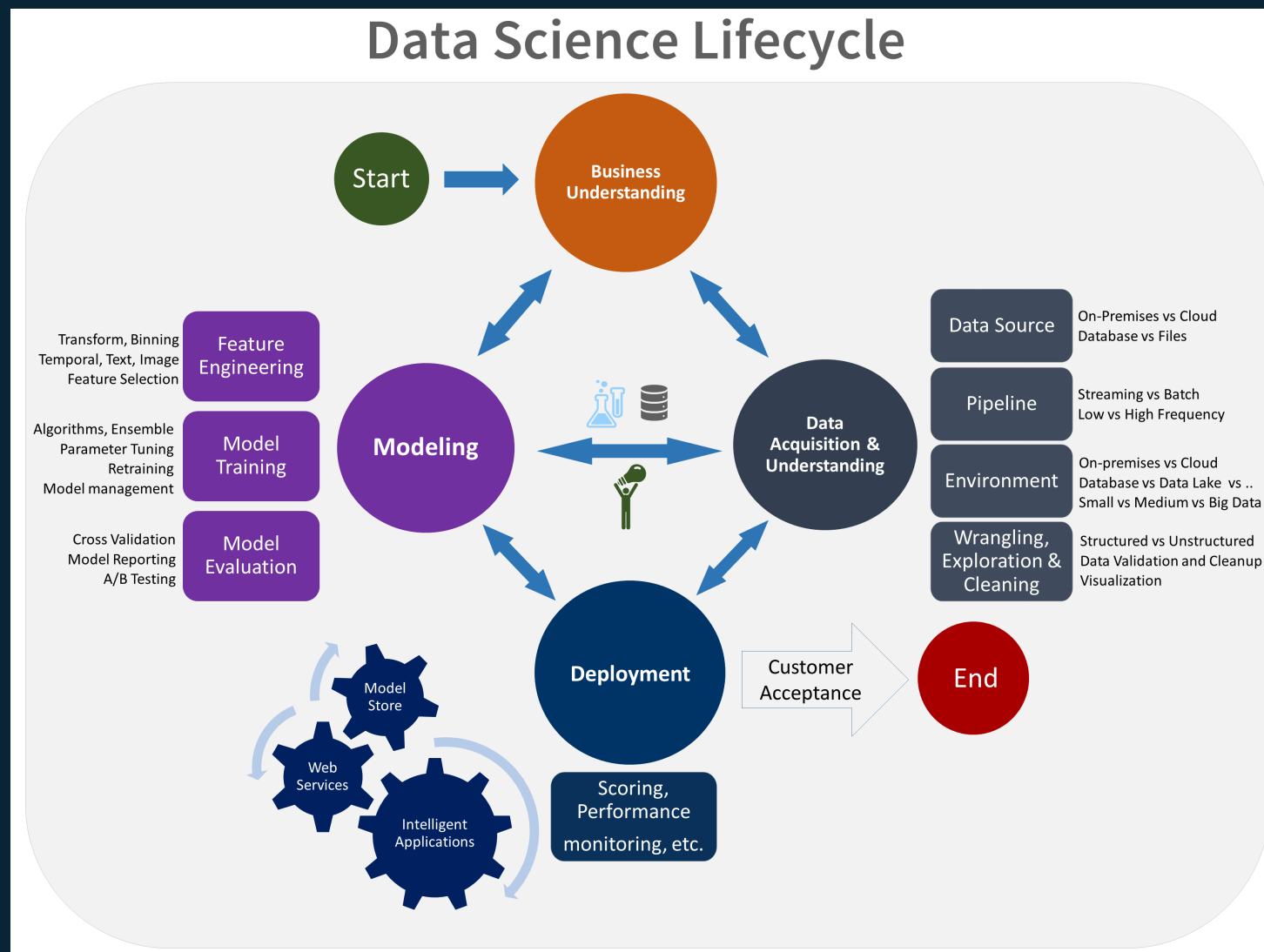
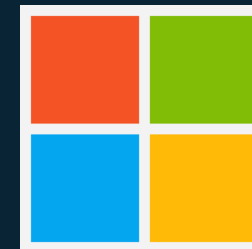
“Do I need to learn C++ to be a Data Scientist?” “Do DevOps and System admins use the same technologies?” “I really like JavaScript; can I use it in Data Analytics?” Those are some of the questions that our students ask.

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://docs.microsoft.com/de-de/azure/architecture/data-science-process/lifecycle>





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**



2. *Data*



It's your turn:
What is your Data Source?

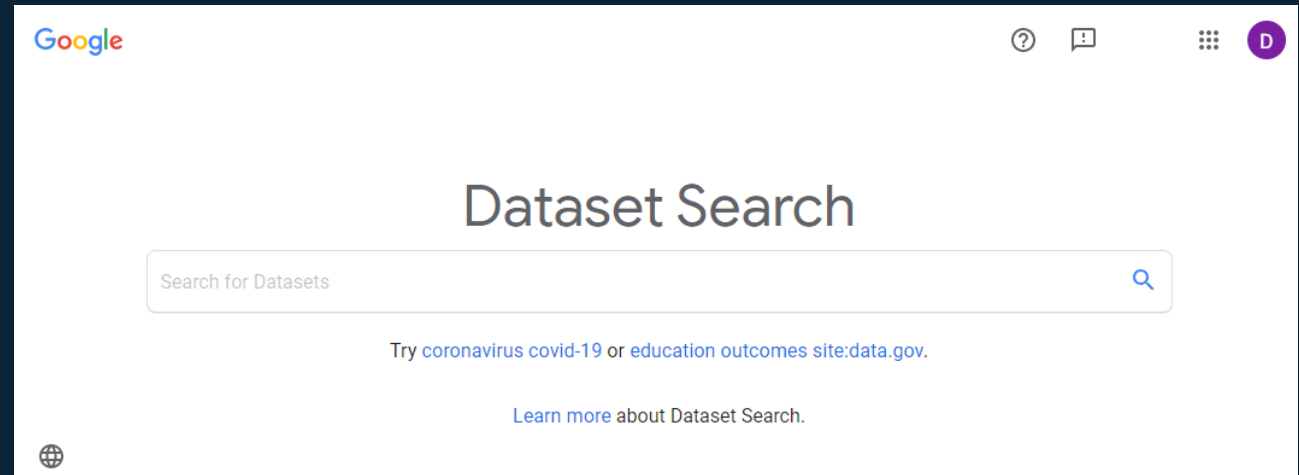
Our client doesn't have any internal data sources that could be used for this project.
Find the data source that you will use to build the solution



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



Data source

Chosen: Stack Overflow developers survey

<https://insights.stackoverflow.com/survey/2021>





3. *Foundations*



1. Legal and data privacy check

Global:

<https://www.privacyaffairs.com/gdpr-fines/>

GDPR Fines Tracker & Statistics

Total Number of GDPR Fines
907

Total Amount of GDPR Fines
€1,544,575,254

Largest Fine
€746,000,000

Amazon Europe Core S.a.r.l. on July 22, 2021 - Luxembourg

Smallest Fine
€28

Unknown on November 18, 2020 - Hungary


Most Recent GDPR Fines

*Only includes finalised cases

DATE	ORG	FINE
01/26/2022	Uppsala hospital board	€152,000
01/26/2022	Uppsala regional board	€28,500
01/21/2022	Property Owner Community	€1,200
01/20/2022	Kaufland România SCS	€3,000
01/18/2022	Garlex Solutions, S.L.	€15,000

TOP 5 BIGGEST GDPR FINES

*Only includes final & binding fines

 Amazon Europe Core S.a.r.l.	€746,000,000
 WhatsApp	€225,000,000
 Google LLC	€90,000,000
 Facebook Ireland Ltd.	€60,000,000
 Google Inc.	€50,000,000

Local:

<https://www.privacylaws.com/media/3263/egypt-data-protection-law-151-of-2020.pdf>

Data Protection Law

قانون

بإصدار قانون حماية البيانات الشخصية

In the name of the People

باسم الشعب

The President of the Republic

رئيس الجمهورية

The Parliament has resolved and issued the following Law:

قرر مجلس النواب القانون الآتي نصه، وقد أصدرناه:

Article (1)

المادة (1)

The provisions of this law and the accompanying law shall apply with regards to the protection of personal data of natural persons partly or fully processed electronically by any holder, controller or processor.

يعمل بأحكام هذا القانون والقانون المرافق في شأن حماية البيانات الشخصية المعالجة إلكترونياً جزئياً أو كلياً لدى أي حائز أو متحكم أو معالج لها، وذلك بالنسبة للأشخاص الطبيعيين.

Article (2)

المادة (2)

The provisions of this law shall apply to any person that commits any of the violations stipulated in the accompanying law, if the offender is an Egyptian national inside or outside the Arab Republic of Egypt, or a non-Egyptian residing within the Arab Republic of Egypt, or a non-Egyptian outside the Arab Republic of Egypt provided that the act is punishable in any form in the country where it occurred, and the data subject of the crime belongs to Egyptian nationals or non-Egyptians residing within the Arab Republic of Egypt.

تسري أحكام هذا القانون والقانون المرافق له على كل من ارتكب إحدى الجرائم المنصوص عليها في القانون المرافق متى كان الجاني من المصريين داخل الجمهورية أو خارجها، أو كان من غير المصريين المقيمين داخل الجمهورية، أو كان غير المصريين خارج الجمهورية إذا كان الفعل معاقباً عليه في الدولة التي وقع فيها تحت أي وصف قانوني، وكانت البيانات محل الجريمة لمصريين أو أجانب مقيمين داخل الجمهورية.

Article (3)

المادة (3)

The provisions of the accompanying law do not apply to the following:

لا تسري أحكام القانون المرافق على ما يأتي:

1. Personal data of third parties retained by natural persons and processed for personal use.

1. البيانات الشخصية التي يحتفظ بها الأشخاص الطبيعيين للغرض ويتم معالجتها للاستخدام الشخصي.



2. How to structure your project

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

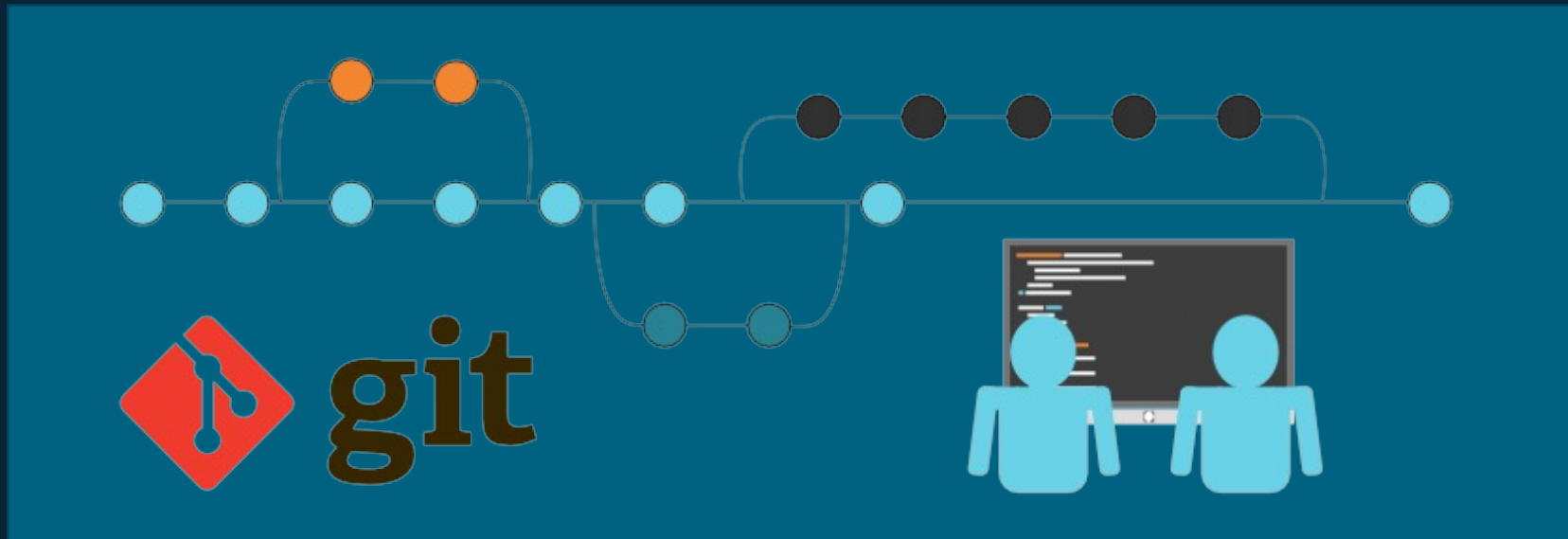
Directory structure

```
|— 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
|   └─ 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/>





4. *Preprocessing*



Preprocessing at first glance

...

1. String values in years need to be replaced



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

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

...



Jira - Kanban

Jira Software

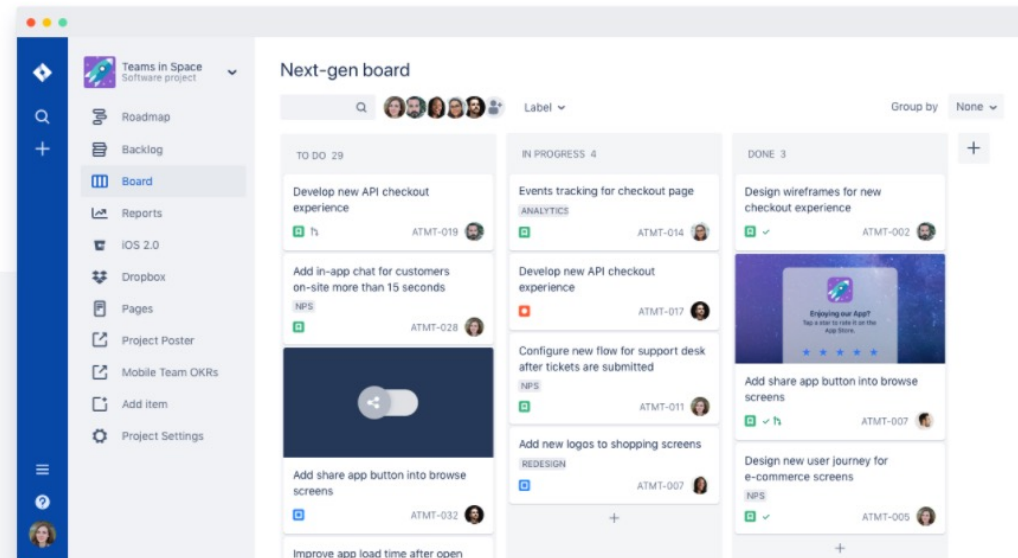
Features

Product guide

Pricing

Enterprise

Get it free



A Jira scrum board for every team

Although Jira scrum boards are ideal for highly technical teams who practice agile methodologies, teams of all types can take advantage of the key concepts of scrum and use the Jira scrum board to facilitate smooth project management. Here are a few ideas.



5. Descriptive Analytics



“Asking the right question is half of the answer”



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



Levels of descriptive analytics

- 1. Stats or summary tables*
- 2. Visualizations*
- 3. Unsupervised learning
(e.g. clustering)*



Wrap Up

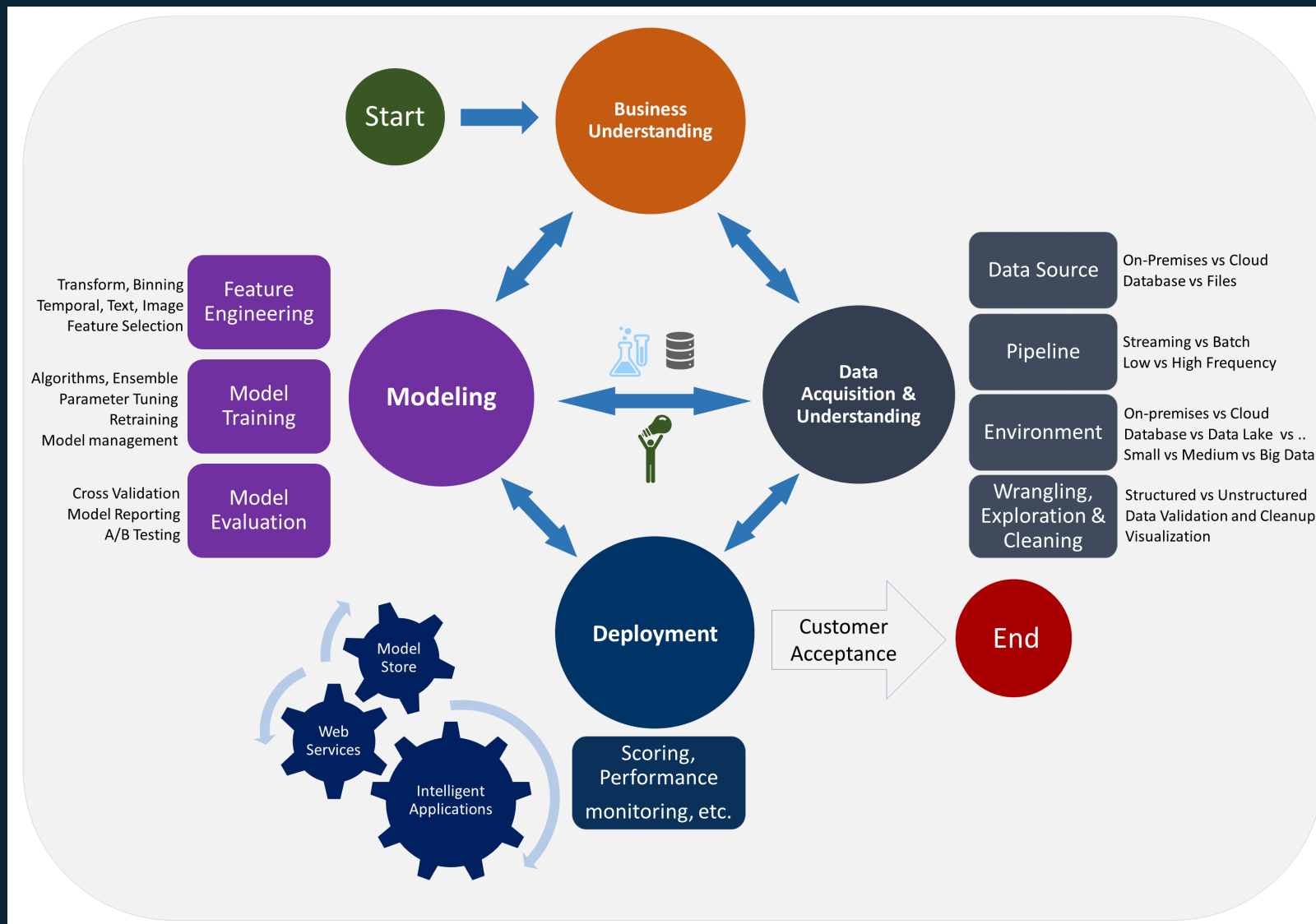
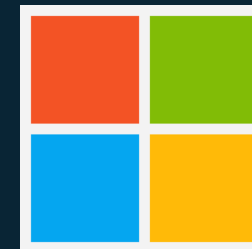


Wrap-up: **Today you have learned about**

- *Build a business case*
- *Find suitable data sources*
- *Verify legal rights*
- *Track your project via Git*
- *Explore and preprocess data*
- *Collaborate with your team using Kanban*



Wrap-up:





Till next time:

- Form your team and create your Kanban board
- Create your project directory and track in a new GitHub repo
- Preprocess your raw data and export it to a pickle file
- Complete your descriptive analytics part – understand your data and get insights to be used in the modelling



Questions?