

Deloitte.

Deloitte Data Science
Academy

Pavel Milička



Professional
Growth Platform
by Deloitte

MAKING AN
IMPACT THAT
MATTERS
since 1845

Other events organized



26.-27.10. 2023	SAFe® 5 Product Owner / Product Manager	workshop
1.11.2023	Data Storytelling	workshop
9-10.11.2023	Leading SAFe	workshop
16.11.2023	OKR (Objectives and Key Results)	workshop
22.11.2023	Datová gramotnost pro začátečníky	workshop
23.-24.11.2023	SAFe® 5 Scrum Master (SSM)	workshop
30.11. 2023	Data-driven HR	workshop

More info at akce.deloitte.cz

Expectations?

- What are your expectations?
- What do you want to learn, why did you choose this course?
- What is your background, how familiar are you with SQL, Programming, Data Science?
- What tools do you use on daily basis?



Goal

Move from zero level to intermediate level in data science.

Learn basics of Python

Use Python and its data science tools

Understand standard machine learning tasks

Have an overview of methods and use cases related to AI

Develop models for Clustering, Classification and Regression

Be able to evaluate models



Lecturers' introduction



Pavel Milička

CE AI Lead

- Master degree in Artificial Intelligence from Czech Technical University, co-author of 3 research papers – Bio-inspired hexapod control, Bilevel optimization.
- Model for predicting performance based on CVs.
- Expedition optimization for major steel manufacturer.
- Technical owner of Deloitte Dynamic Pricing solution.
- Leading project for Debt Collection Optimization for foreign Health Insurer.
- Collaborating on automated open source intelligence system.
- Creator of internal Task Mining solution from activity logs.
- Invoice OCR Solution.
- Leading team of approx. 30 AI Specialists across CE.

Content

1	Introduction to Data Science, Data Science Tools, Python as a programming language, Python Basics, Git – code collaboration, Pure Python Data Exploration	HA1
2	Features and, Clustering, behavioral segmentation, Hands on training for segmentation with Python	HA2
3	Propensity models, single factor analysis, binning, binary classifiers, Hands on training for data exploration, classification and regression	HA3
4	Introduction to NLP and Computer Vision, Selected topics	HA4

Format

Overview

4

Days

During which we go through Data Science, Python, NLP, and much more

6

Hours

Every lecture has 6 hours.



Calls and emails outside the room, PC for hands-on exercises only.

4

Home assignments

You need all 4 home assignments accepted to qualify for a final certificate.

Major theme

During the course, we will work on various Data Science problems in the hands-ons or your homework. For the homework it is possible to come up with your own topics and assignments if it will have reasonable complexity. The homeworks topics are:

- 1) Python – coding exercise (scraping)
- 2) Clustering/Segmentation
- 3) Classification/Regression
- 4) Natural Language Processing / Computer Vision

You will learn how to work it end to end from data crunching, calculating predictors, target definition and feature engineering to modeling and final visualizations.

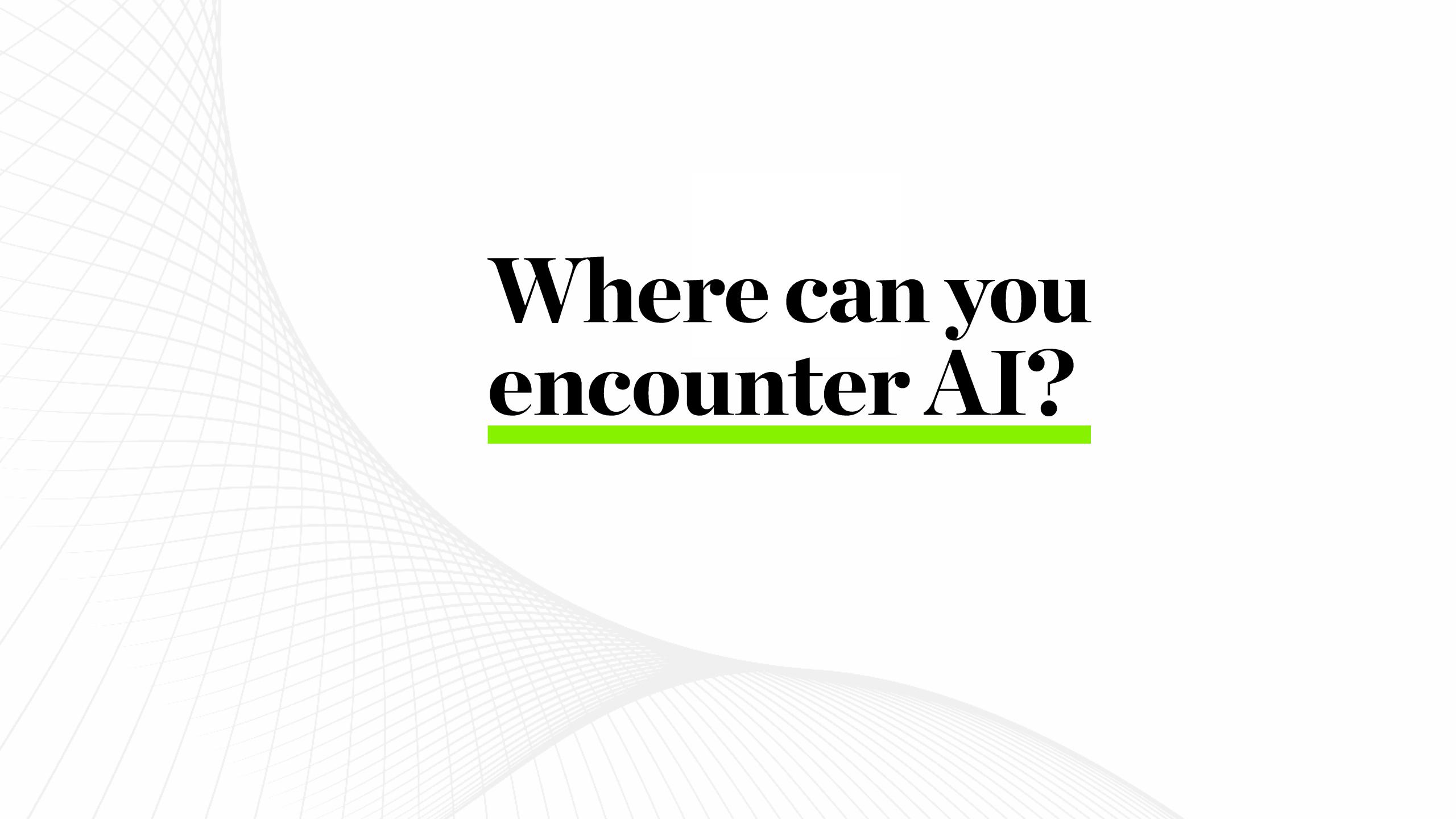
It is possible to come up with your own problems if the scale of the problem is comparable.

Today's program

1	Introduction to Data Science, Data Science Tools, Keboola introduction, Python basics, Git – Code collaboration, Pure Python data exploration	HA1
2	Features and dataset preparation, Clustering, Behavioral segmentation, Hands on training for data exploration and clustering with Python	HA2
3	Regression and Classification, Hands on training for data exploration, classification and regression	HA3
4	Introduction to NLP and Computer Vision, AutoML	HA4

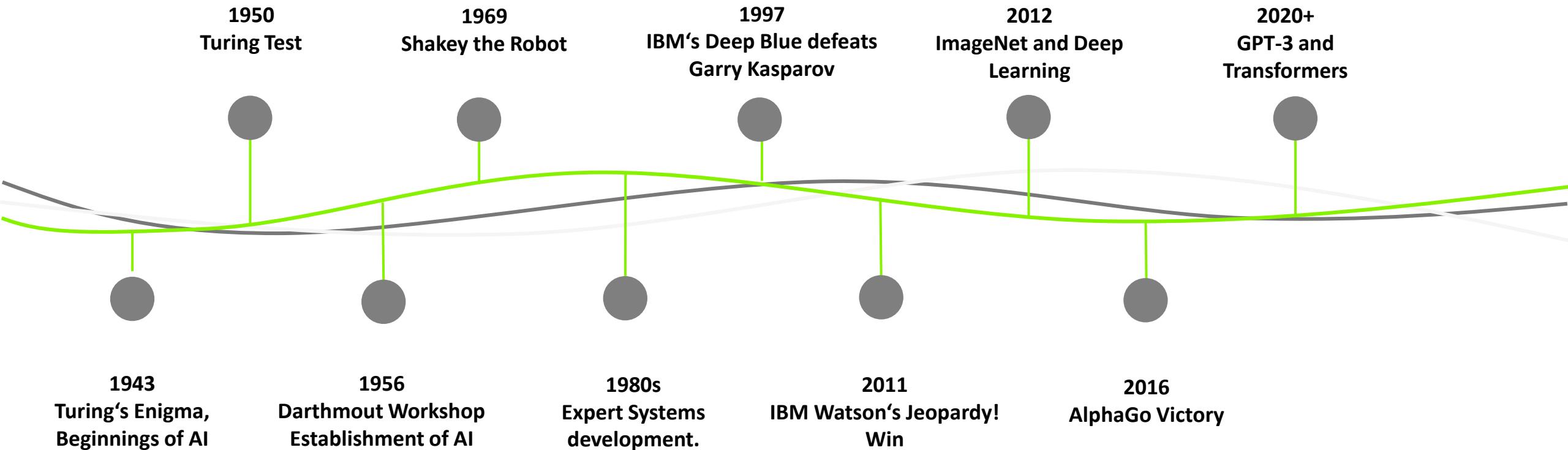
AI / ML / DS





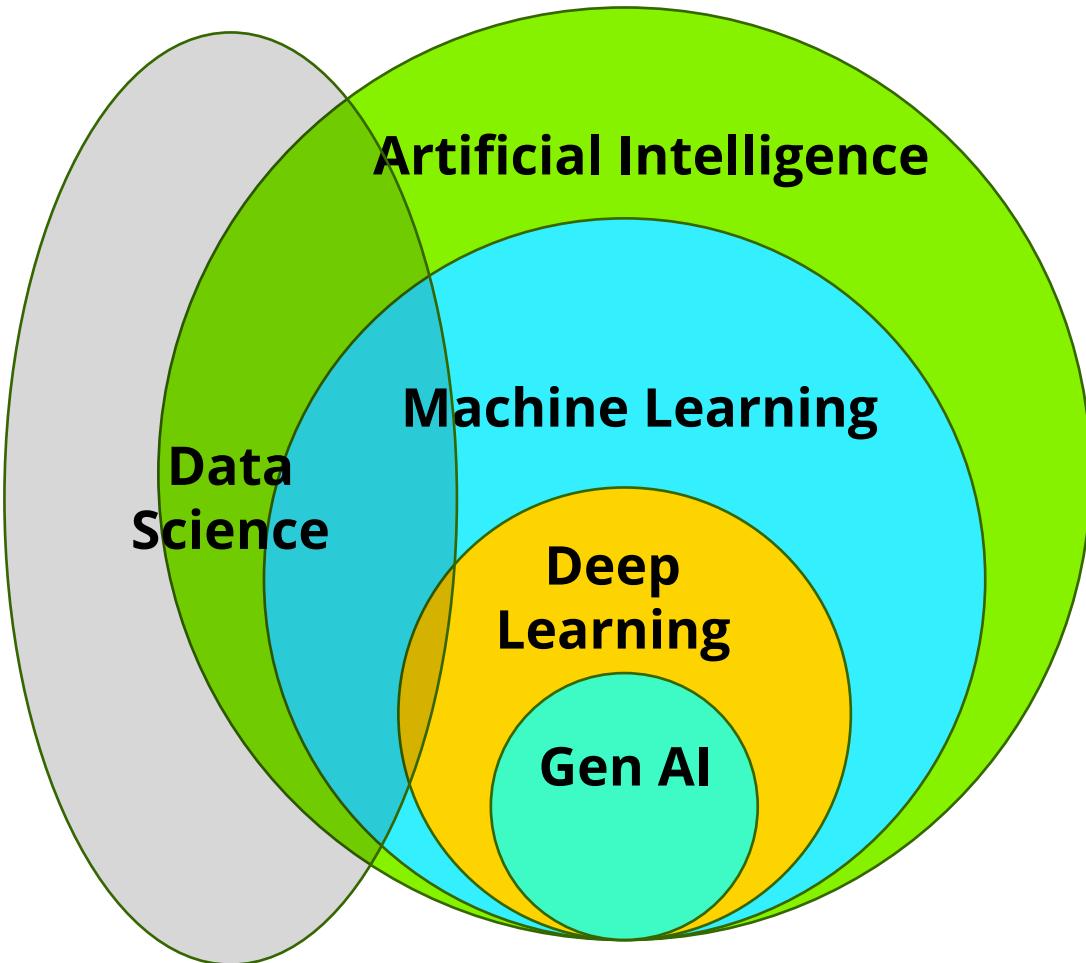
Where can you
encounter AI?

We are living in the era of AI Revolution.



The amount of development in the field is unprecedented in the last years.

What is AI



Artificial Intelligence is broadly defined as the ability of machines to mimic human behavior.

Machine Learning algorithms leverage statistical techniques to automatically detect patterns and make predictions or decisions based on historical data that they are trained on.

Deep learning is defined as a machine learning technique that uses deep neural networks which are able to teach computer to learn from the data that is inspired by humans.

Generative AI focuses on creating models capable of generating new content that resemble existing data.

<https://synoptek.com/insights/it-blogs/data-insights/ai-ml-dl-and-generative-ai-face-off-a-comparative-analysis/#:~:text=Generative%20AI%2C%20a%20branch%20of,might%20be%20created%20by%20humans.>

What is Machine Learning

Unsupervised Learning

„I have data and want to group them together based on their characteristics.“

- Client Segmentation
- Recommender Systems
- Anomaly Detection
- Customer Lifetime Value



Supervised Learning

„I have historical data and want to predict events/values based on them, i.e., find the dependencies and utilize them in my decision making.“

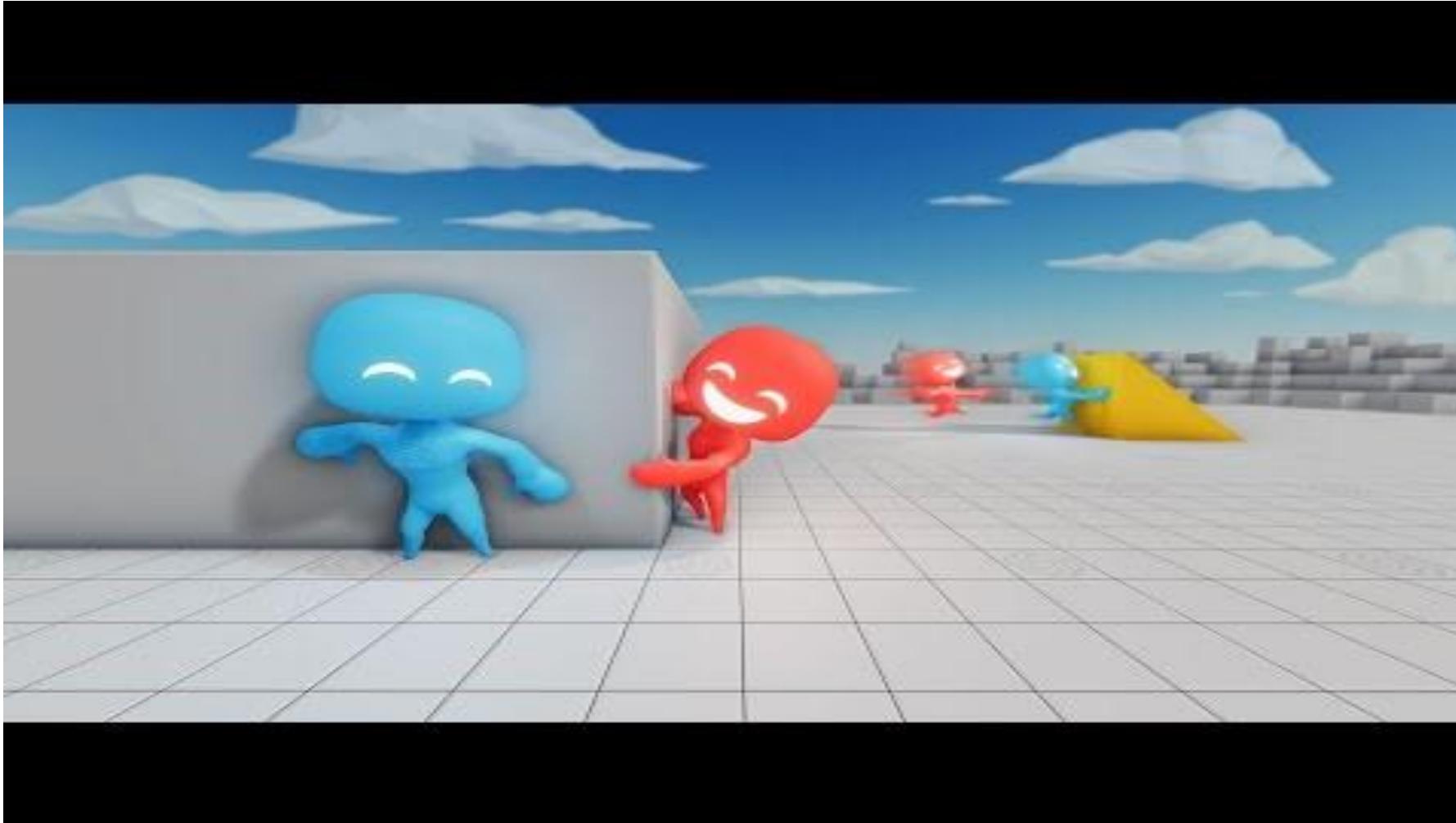
- Churn prediction
- Campaign targeting
- Demand forecasting

Reinforcement Learning

„I have an agent and representation of a world and I want my agent to learn how to achieve his goals by interacting with the goal.“

- Route planning
- Self-driving cars
- Gaming

Reinforcement learning

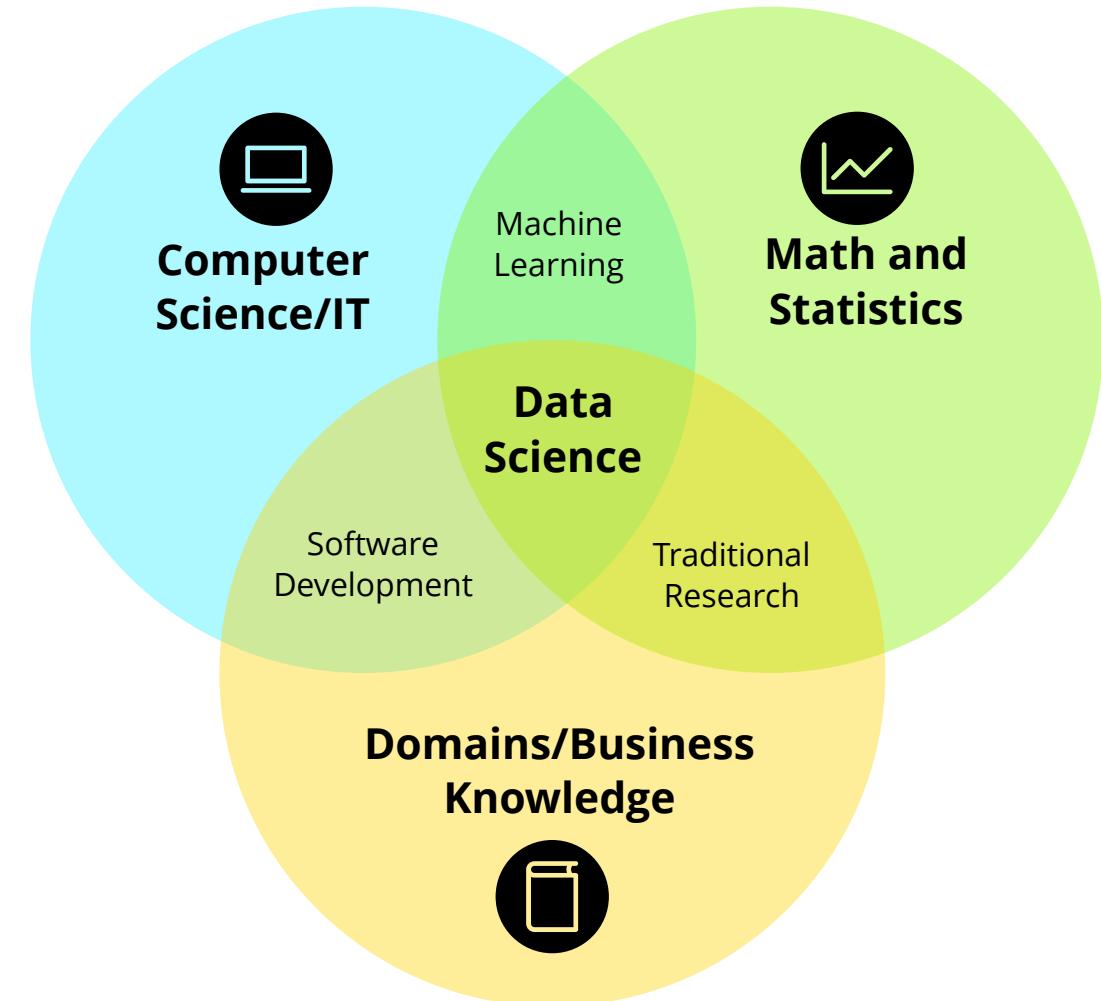


What is Data Science

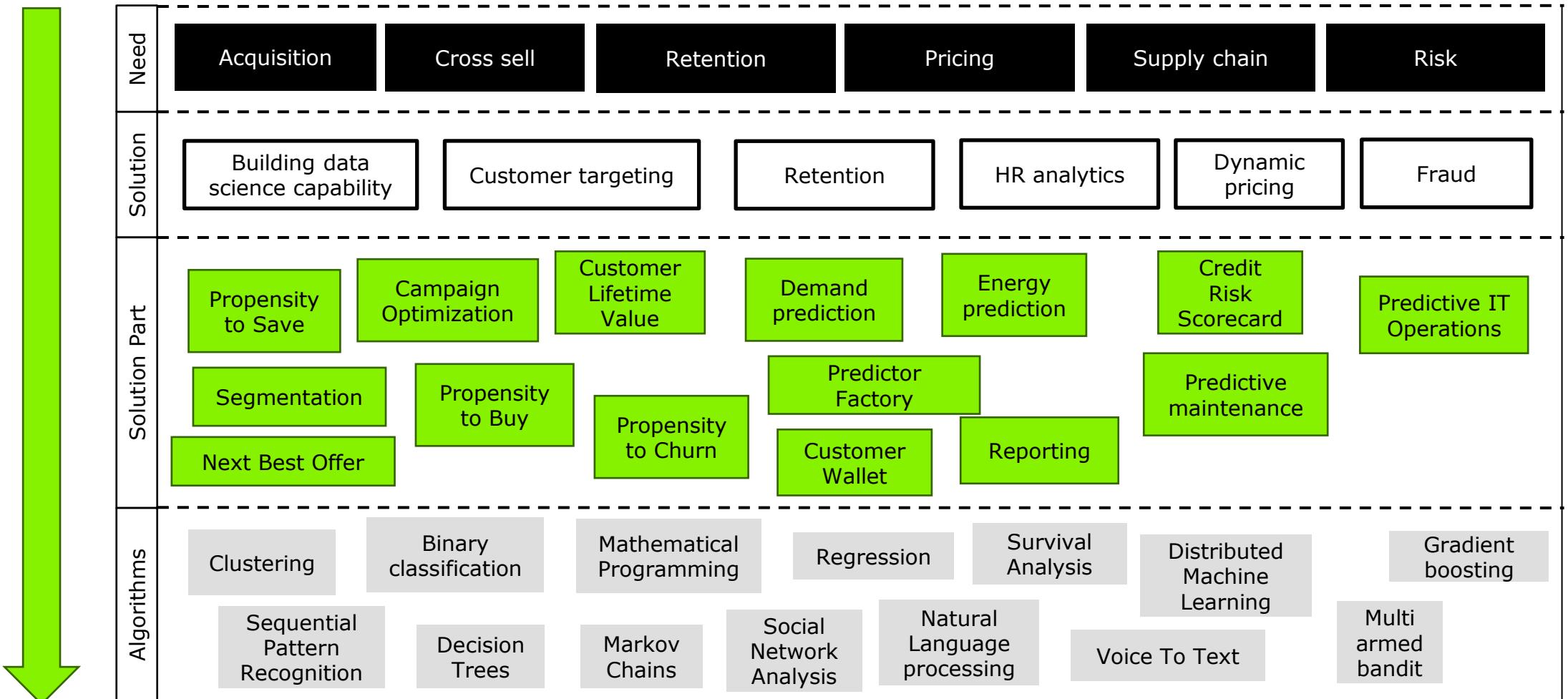
Data science is the extensive use of **data**, statistical and quantitative **analysis**, exploratory, predictive models, and fact based management to drive decisions and **actions**.

Data scientists come with flavours.

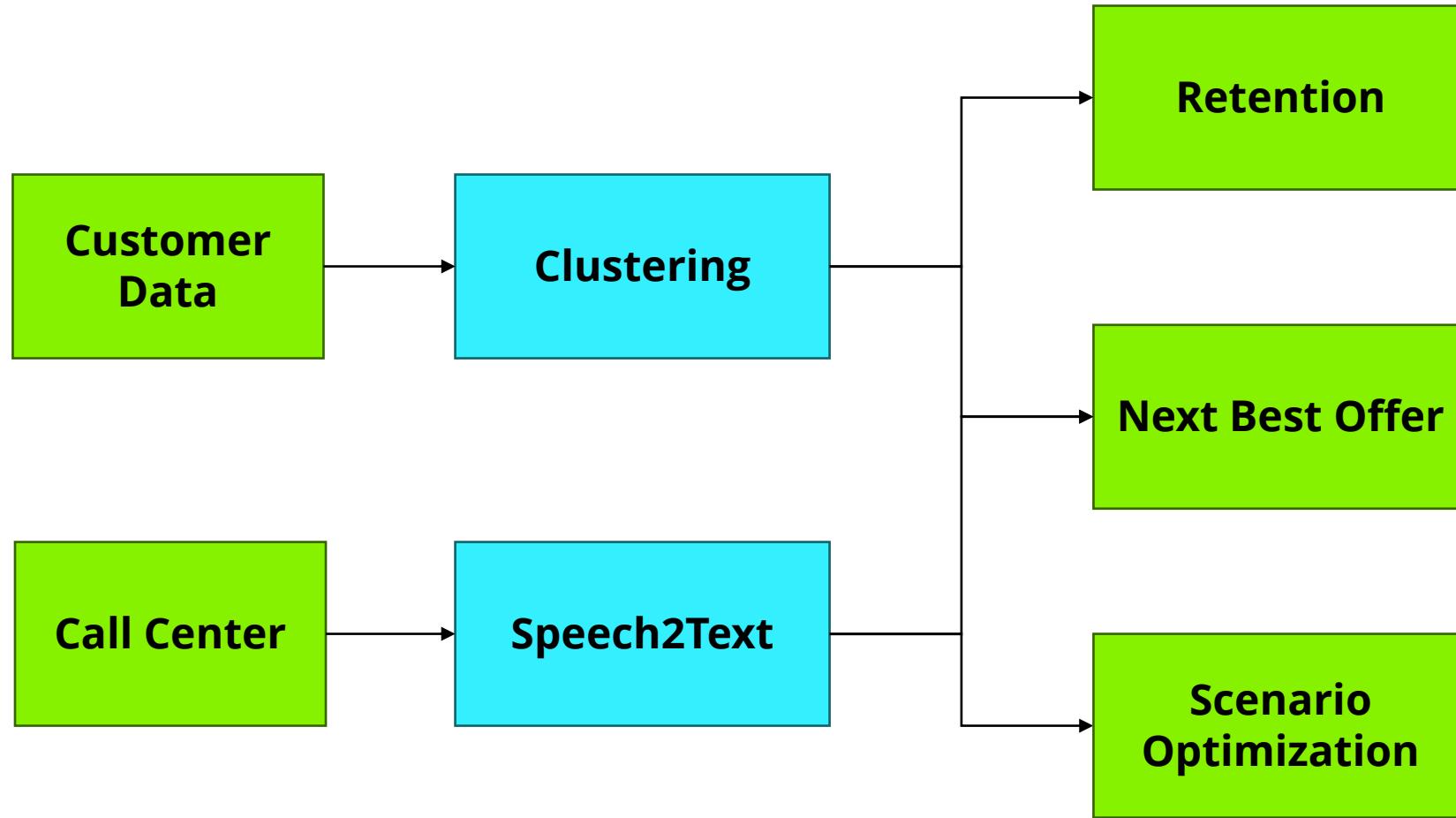
- Data Scientist
- ML Engineer
- Business Analyst
- „Prompt Engineer“



How to think about use cases

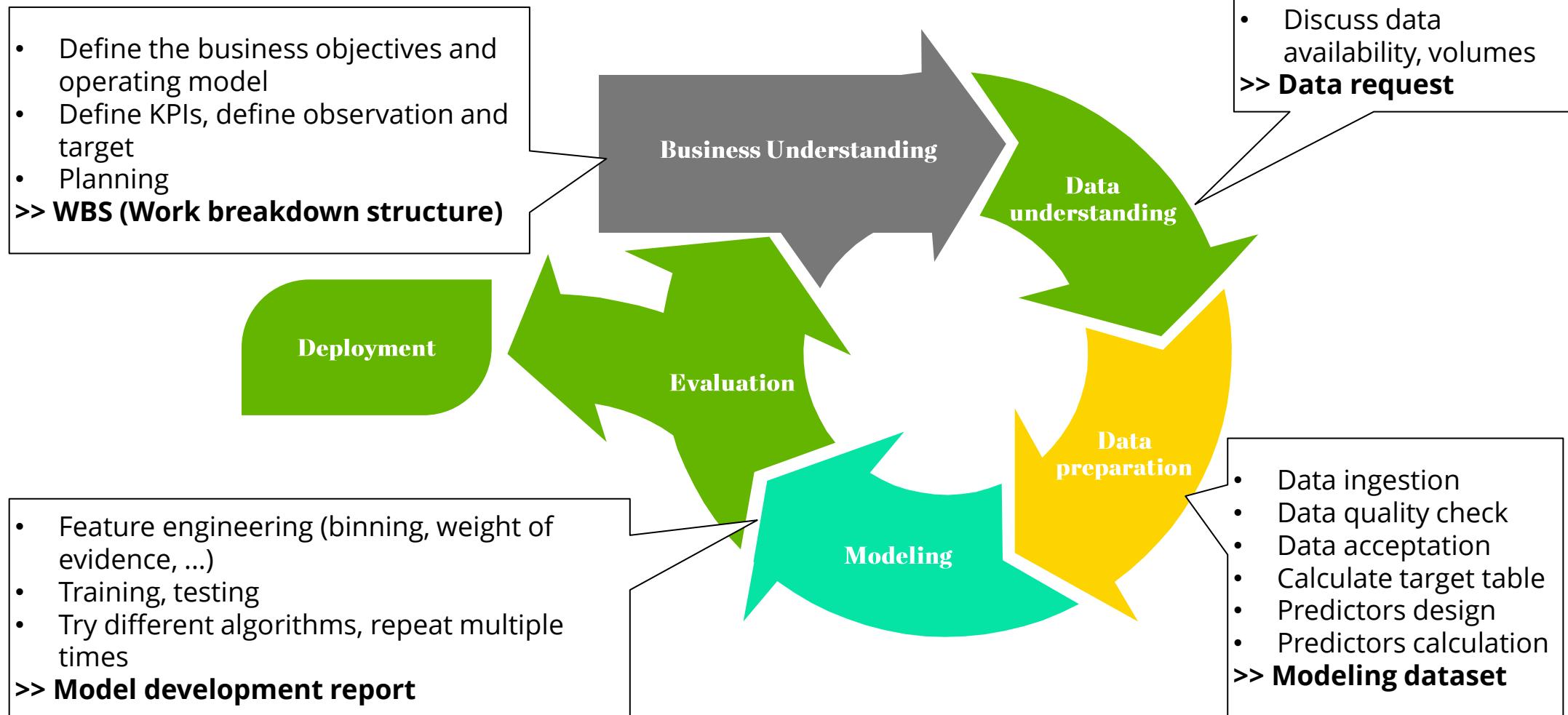


How to think about use cases – Example



CRISP-DM Framework

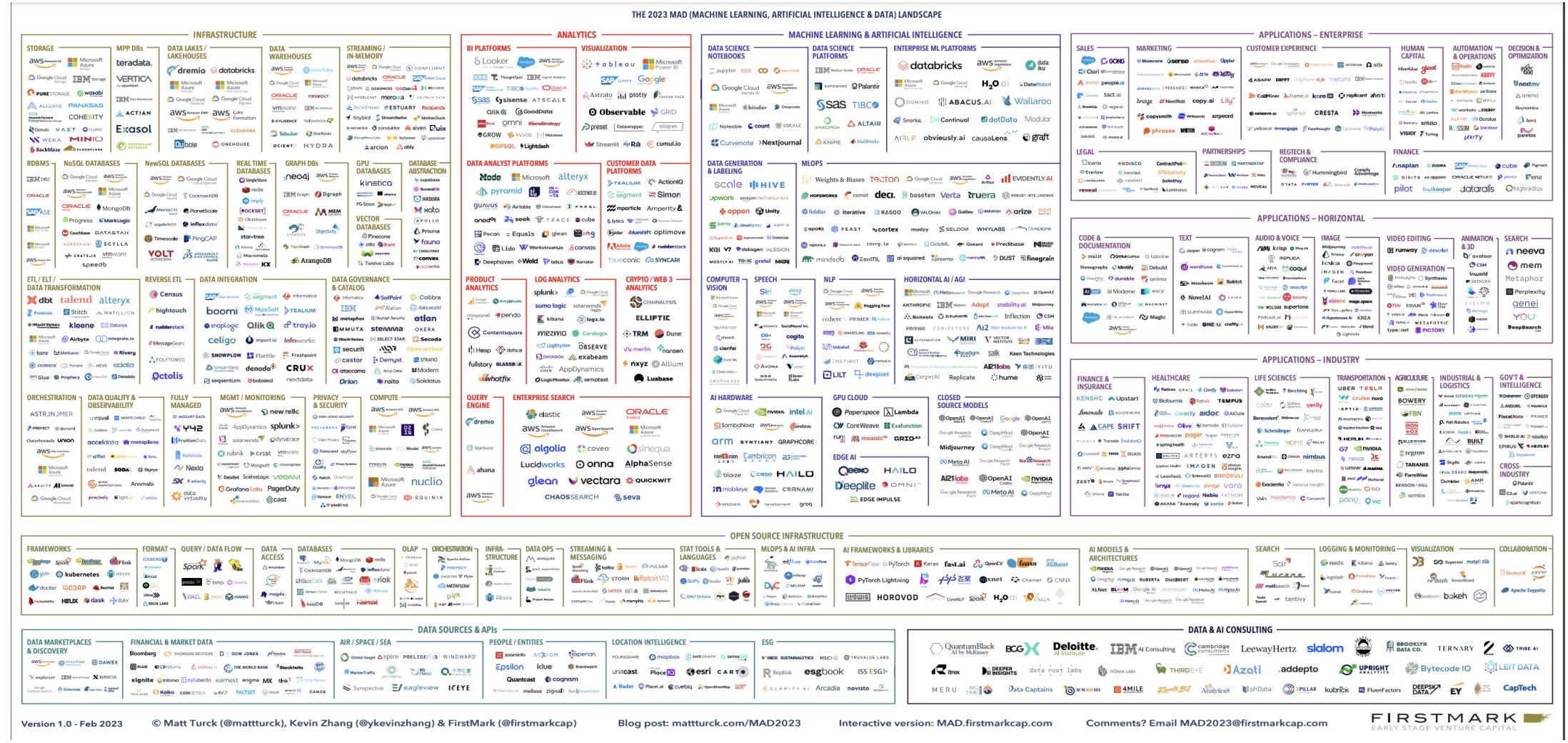
Applicable to all Data Science jobs



Tools



AI Tools Landscape



AI Tools Landscape – The Essentials

Languages and Utils

- Python
- PySpark
- SQL
- Git

Vizualize

- Python
 - Matplotlib
 - Seaborn
 - Plotly
 - Streamlit
- PowerBI
- Tableau
- Qlik

Transform

- Python
 - Numpy
 - Pandas
- SQL

Model

- Python
 - Scikit-learn
 - Tensorflow
 - PyTorch
 - PyCaret
 - ExplainerDashboard
- (OpenAI or similar GenAI tool).

Build (ML Engineering)

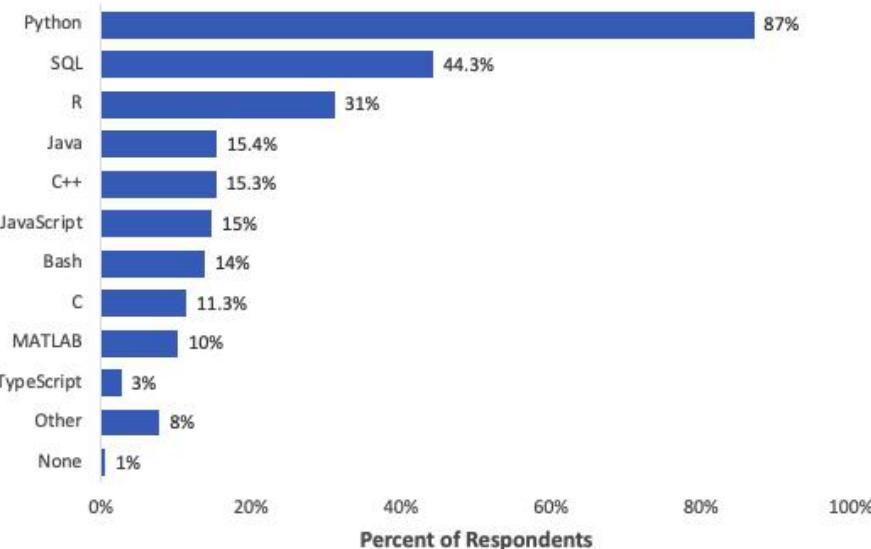
- Docker
- Kubernetes
- AWS/GCP/Azure
- Databricks
- Terraform
- API
- Python
 - Dagster

Other

- ElasticSearch
- R
- Julia
- Airflow

Most popular tools

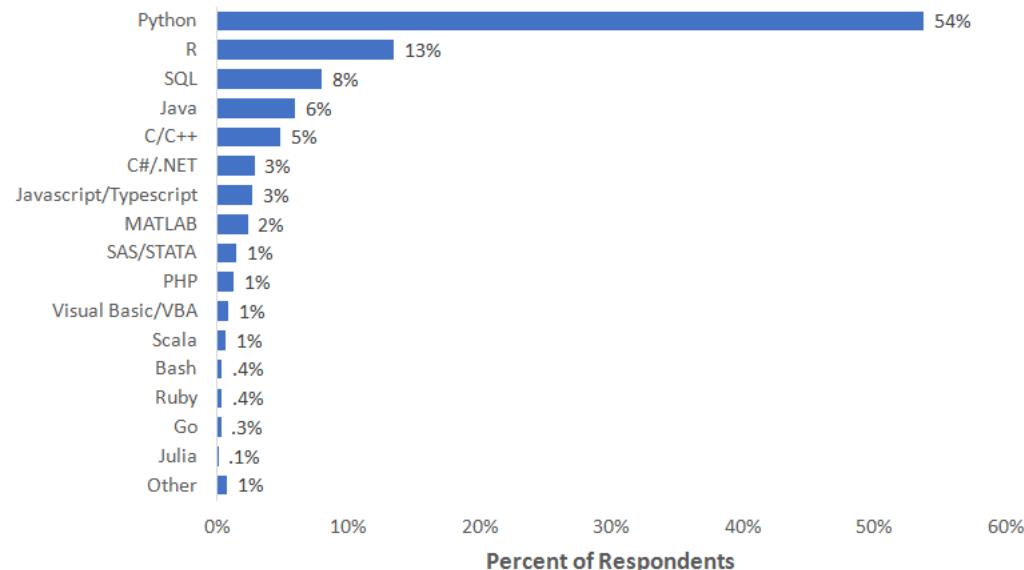
What programming languages do you use on a regular basis?



Note: Data are from the 2019 Kaggle ML and Data Science Survey. You can learn more about the study here: <https://www.kaggle.com/c/kaggle-survey-2019>.

A total of 19717 respondents completed the survey; the percentages in the graph are based on a total of 14762 respondents who provided an answer to this question.

What specific programming language do you use most often?

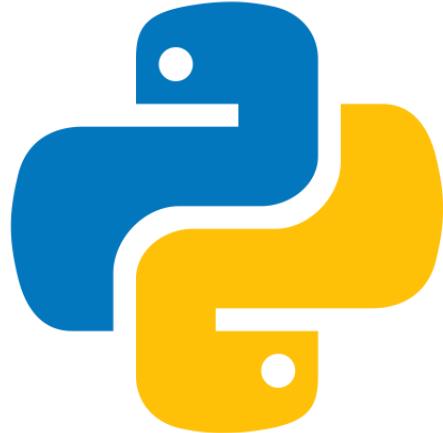


Note: Data are from the 2018 Kaggle ML and Data Science Survey. You can learn more about the study here: <http://www.kaggle.com/kaggle/kaggle-survey-2018>.

A total of 23859 respondents completed the survey; the percentages in the graph are based on a total of 15222 respondents who provided an answer to this question.

Python

"Python is powerful... and fast; plays well with others; runs everywhere; is friendly & easy to learn; is Open." **python.org**

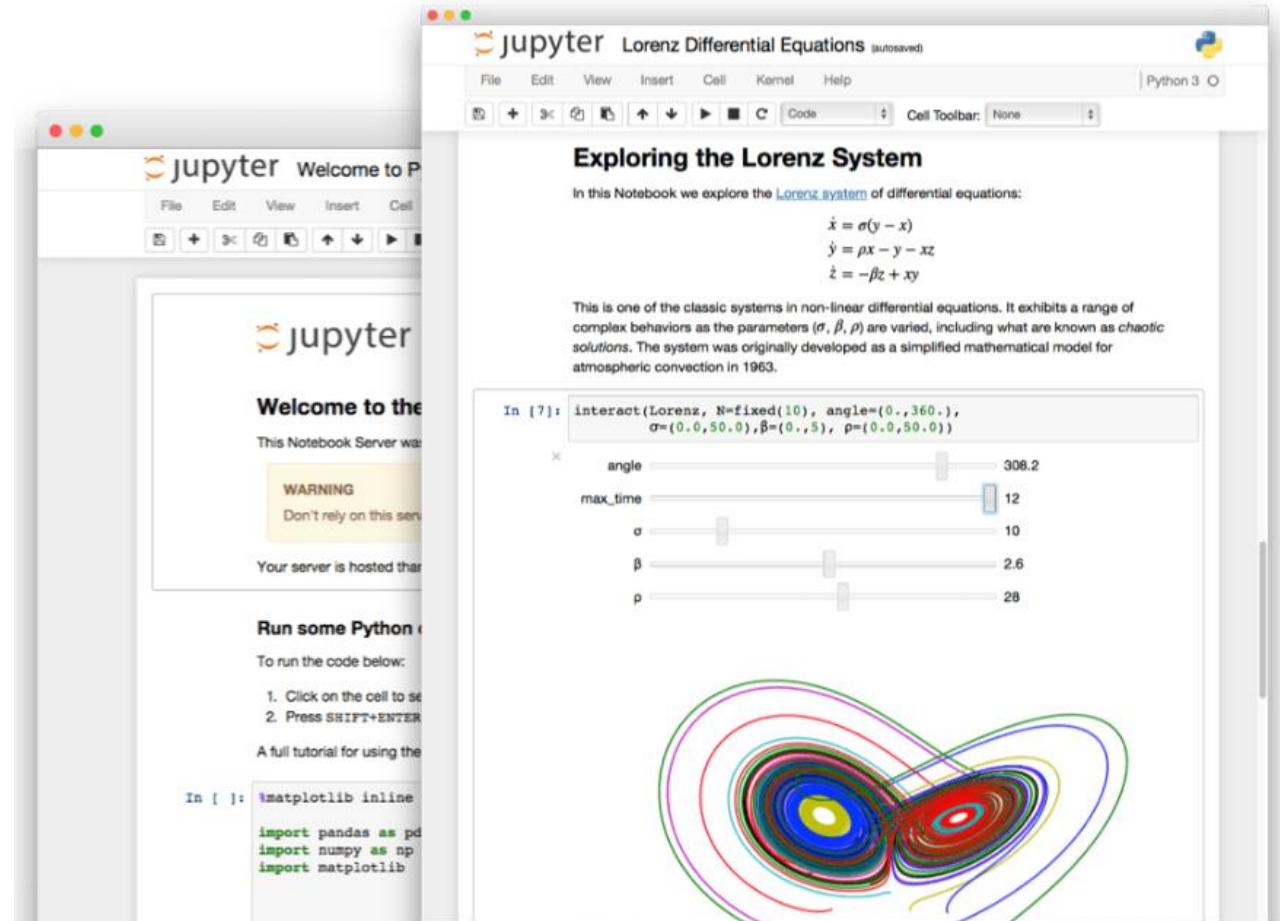


- High level programming language with strong support for data science. Created in 1991. Philosophy of code readability. Open source.
- Dynamic typing, automatic memory management. Procedural, object oriented, functional.
- IDEs (Spyder, PyCharm, VSCode, ...)
- Large community

IDEs - Jupyter notebooks / Jupyter lab

The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text.

<https://jupyter.org/>



<https://jupyter.org/>

IDEs – Spyder

- Scientific Python IDE
- Similar to RStudio
- Loved by R & Python bilingual data scientists
- Part of Anaconda distribution
- Version 4 released in 2020 bringing in many new features

The screenshot shows the Spyder IDE interface. The code editor displays a Python script named `interpolation.py` with the following content:

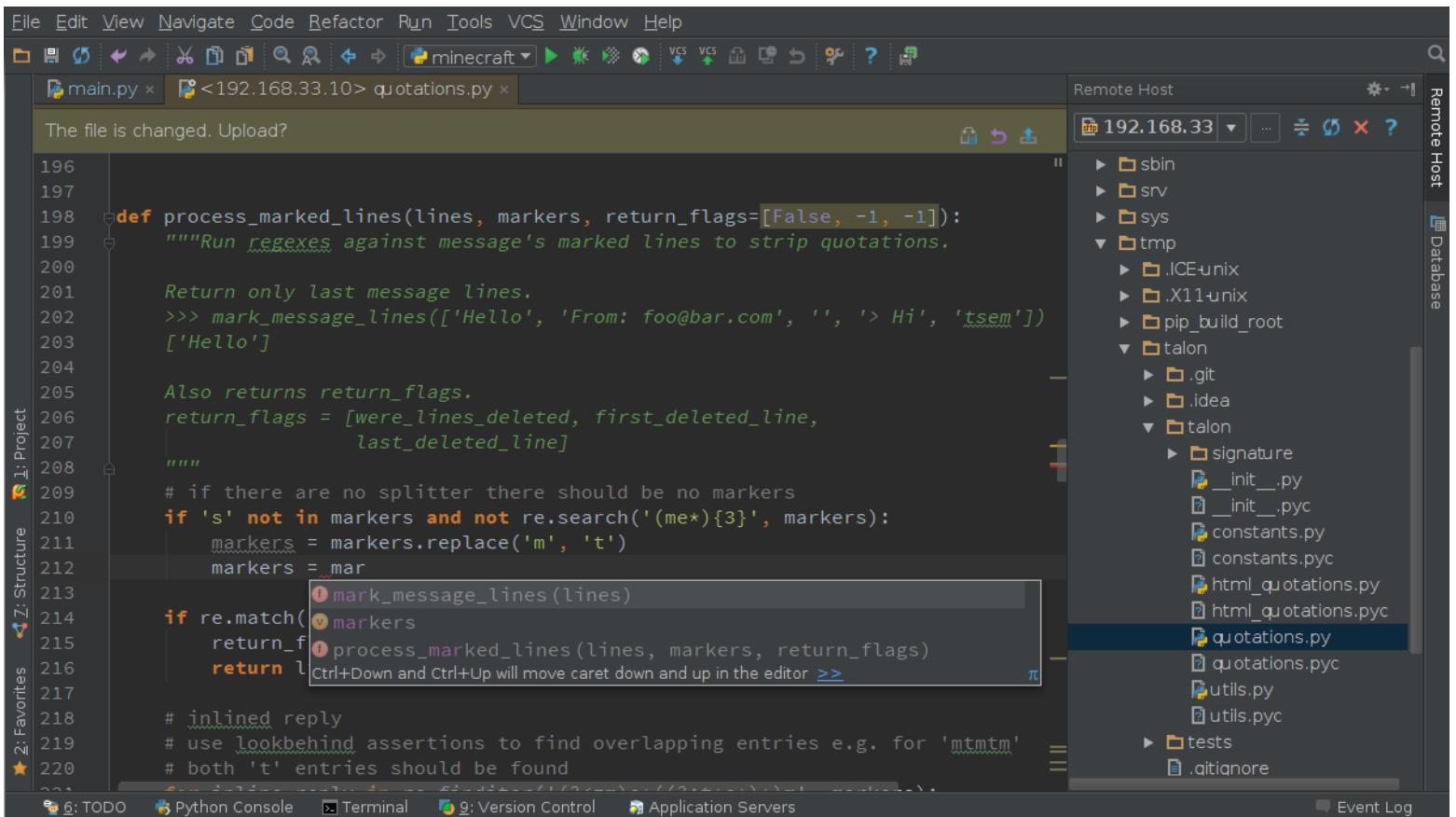
```
6 import pylab
7 from numpy import cos, linspace, pi, sin, random
8 from scipy.interpolate import splprep, splev
9
10 # %% Generate data for analysis
11 t = linspace(0, 1.75 * 2 * pi, 100)
12
13 x = sin(t)
14 y = cos(t)
15 z = t
16
17 x += random.normal(scale=0.1, size=x.shape)
18 y += random.normal(scale=0.1, size=y.shape)
19 z += random.normal(scale=0.1, size=z.shape)
20
21 # Add noise
22
23 # %% Perform calculations
24
25 # Spline parameters
26 smoothness = 3.0 # Smoothness parameter
27 k_param = 2 # Spline order
28 nests = -1 # Estimate of number of knots needed (-1 = maximal.)
29
30 # Find the knot points
31 knot_points, u = splprep([x, y, z], s=smoothness, k=k_param, nests=nests)
32
33 # Evaluate spline, including interpolated points
34 xnew, ynew, znew = splev(linspace(0, 1, 400), knot_points)
35
36 # %% Plot results
37
38 # TODO: Rewrite to avoid code smell
39 pylab.subplot(2, 2, 1)
40 data, = pylab.plot(x, y, 'bo', label='Data with X-Y Cross Section')
41 fit, = pylab.plot(xnew, ynew, 'r-', label='Fit with X-Y Cross Section')
42 pylab.legend()
43 pylab.xlabel('x')
44 pylab.ylabel('y')
45
46 pylab.subplot(2, 2, 2)
47 data, = pylab.plot(x, z, 'bo', label='Data with X-Z Cross Section')
48 fit, = pylab.plot(xnew, znew, 'r-', label='Fit with X-Z Cross Section')
49 pylab.legend()
50 pylab.xlabel('x')
```

The Variable explorer on the right shows various variables and their values, such as `array_int8`, `array_uint32`, `bars`, `df`, `filename`, `list_test`, `rows`, `r`, `radii`, `region`, `rgb`, `series`, and `test_none`.

The bottom right corner shows a 3D surface plot generated from the data.

IDEs – PyCharm

- Rich IDE made by JetBrains (Czech company)
 - Free and Professional edition
 - Mainly focused for software developers (but it has a scientific mode similar to Spyder)



IDEs – VSCode

- Open Source
- Versatile, options are limitless as you can download extensions for remote development, notebook display, etc.
- Git integration
- Add-ons – cloud, Data Wrangler, Jupyter notebooks, Code Completion

The screenshot displays the Visual Studio Code (VSCode) interface. The top bar shows the title 'concat.ipynb' and the status bar indicates 'el_paso_ocr (Python 3.11.5)'. The left sidebar has 'EXPLORER' selected, showing a message 'NO FOLDER OPENED' and buttons for 'Open Folder' and 'Clone Repository'. The main area contains two code editors. The top editor shows Python code for concatenating CSV files:import pandas as pd
import numpy as np
import sys
import os

The bottom editor shows more Python code for reading CSV files from a folder:FOLDER = f"C:\\Users\\pmilicka\\Deloitte (O365D)\\Yanev, Dimitar - El Paso - Invoice data capture\\intermediate"
files = os.listdir(FOLDER)

Below the code editors is a terminal window titled 'TERMINAL' which shows the command 'cmd' and the path 'C:\\Users\\pmilicka'. The status bar at the bottom right shows 'Cell 3 of 8'.

Python core Data Science packages

- Python can be augmented by using packages
- On 2020-07-30, the PyPI package repository features 253,532 packages
- Packages for everything – Web development (Django, Flask, ...), Natural Language Processing (gensim, spacy, nltk, ...), Web scraping (scrapy, selenium, BeautifulSoup...) etc.



Anaconda and Python environments

- Installing your own python environment
 - <https://www.anaconda.com/>
- Conda environments
 - `conda env list`
 - `conda create -n=<environment_name> python=3.8`
 - `conda activate <environment_name>`
- Package management
 - `pip install <package_name>`
 - `pip freeze`
 - `(conda install <package_name>)`

Python Basics



Python Basics hands-on

1. Open the D1_Python_101.ipynb notebook.

Version control systems – Git



Which file is the current one?

- 📄 Data_quality_control_script.sql
- 📄 Data_quality_control_script_Jan_new.sql
- 📄 Data_quality_control_script_latest.sql
- 📄 Data_quality_control_script_latest_adjusted_KP.sql
- 📄 Data_quality_control_script_old.sql
- 📄 Data_quality_control_script_v2_new.sql
- 📄 Data_quality_control_script_version1.sql
- 📄 Data_quality_control_script_version1_old.sql

Version control systems – Introduction

A version control system, or VCS, tracks the history of changes as people and teams collaborate on projects together. As the project evolves, teams can run tests, fix bugs, and contribute new code with the confidence that any version can be recovered at any time. Developers can review project history to find out:

- Which changes were made?
- Who made the changes?
- When were the changes made?
- Why were changes needed?

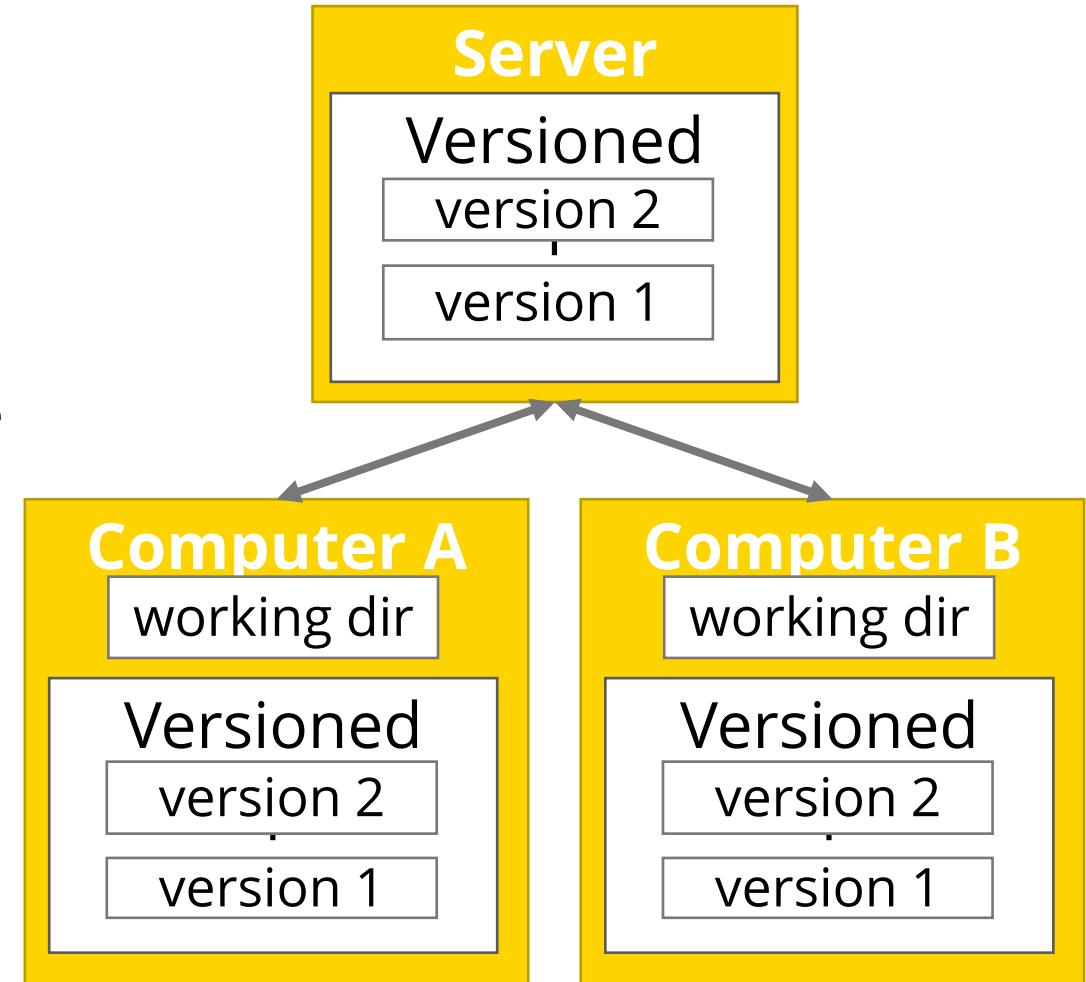
Version control systems – Benefits

- Protects source code from both catastrophe and the casual degradation of human error and unintended consequences.
- Allows work parallelization
- Helps organize changes and repairs in the code and document why the changes were made
- No need for file suffixing, all changes are backed by commit messages stating what was done and why.

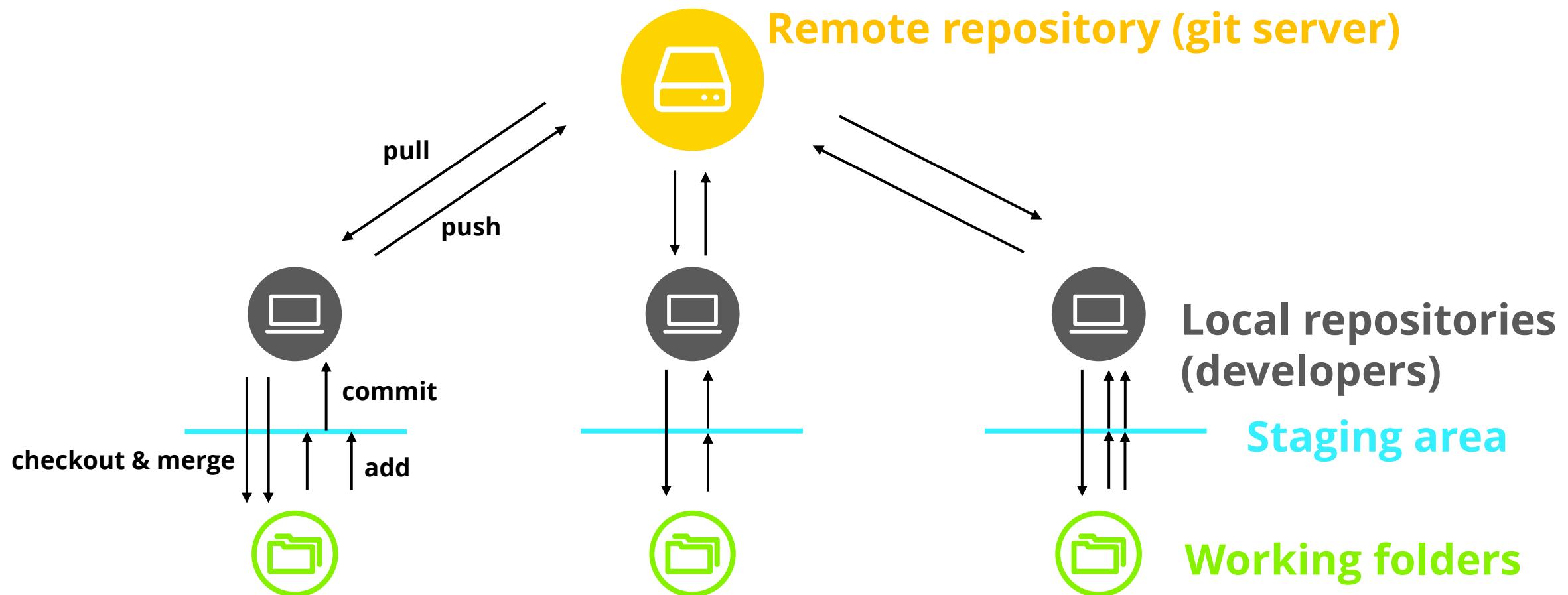
Most common version control systems are **git** and **SVN**.

Git – Introduction

- Created by author of Linux
- Repository on a server – interact with the remote by **clone**, **pull** and **push**
- **Local repository** is a complete copy of the remote server
- Many of the operations are done locally:
 - change branches
 - commit changes to the local repository
 - versioning
- Pushing the local repository makes the versioning visible for the central repository

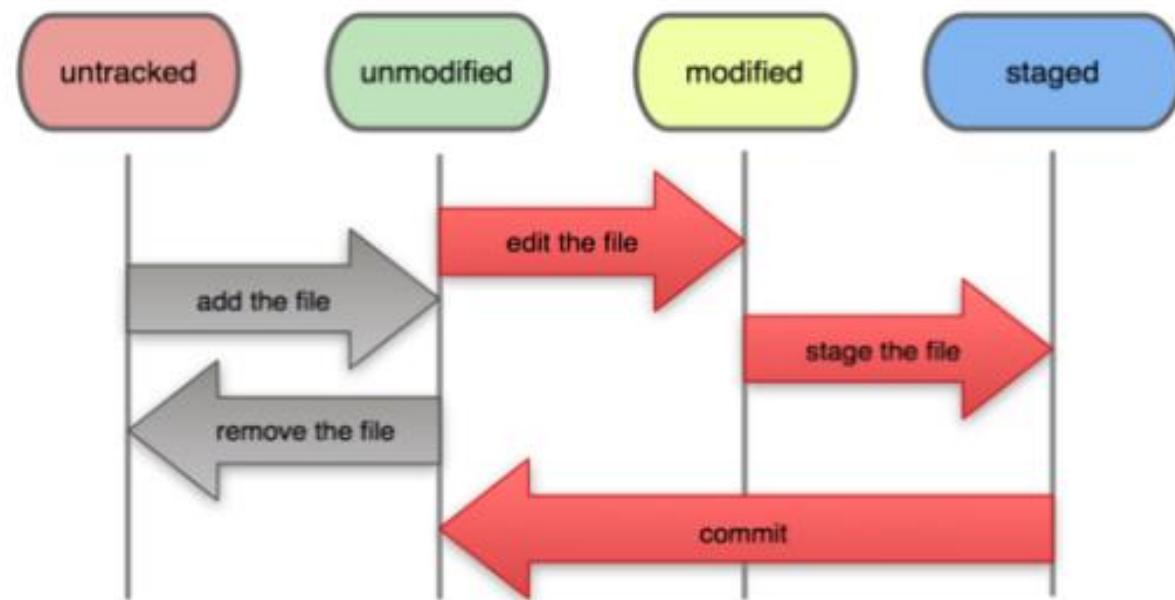


Git Schema



Local operations

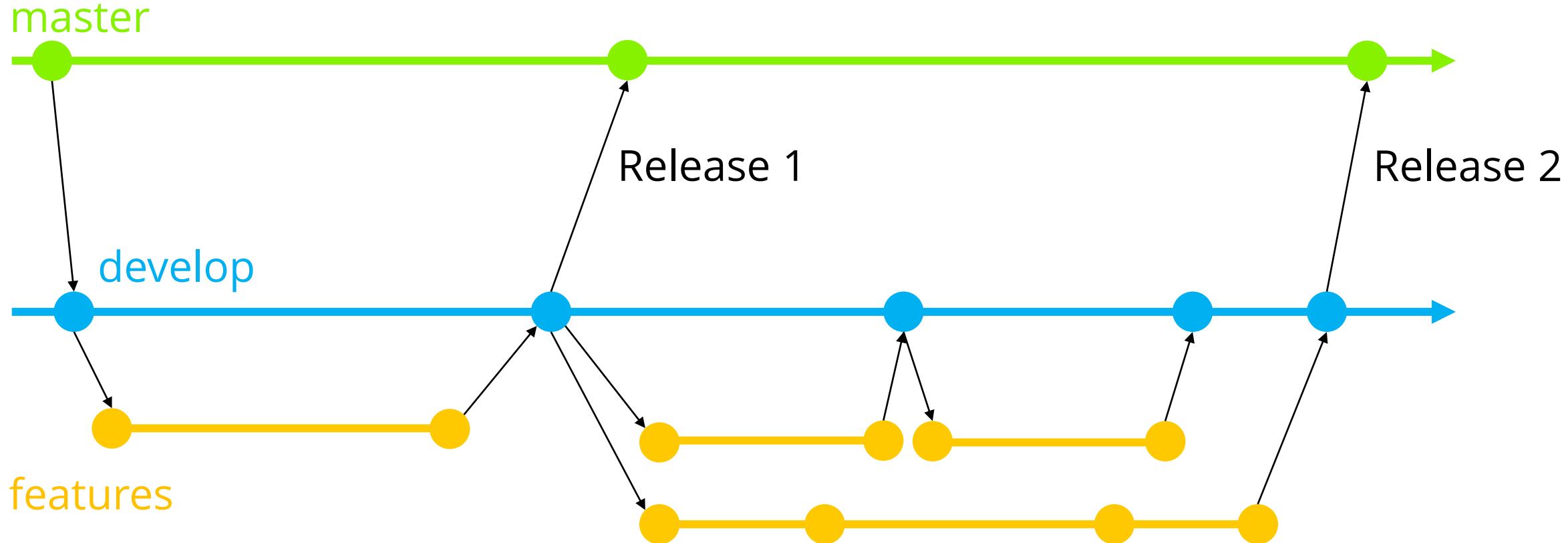
- **Add files** (git add) to tell git to version them
- **Stage files** create a snapshots of the files in he staging area
- **Commit** takes the files in the staging area and stores the snapshot permanently to the git repository



Branching

- Commits are associated with a branch.
- Default git branch is master.
- Branches are created from different branch (basically a copy)
- Branches allow to cluster commits belonging together (e.g. developers trying to implement specific feature vs others fixing a bug)
- Branches can be merged together when the new feature is ready and quality is assured to mess up the production branch.
- When two developers independently make conflicting changes, it creates a merge conflicts showing the conflicting parts with possibilities of resolution (use mine/use theirs/...)

Branching



Branching – Merge conflict

- File that where the same part was modified in two branches will result in **merge conflict**.
- If git is unable to resolve the conflict it adds <<<>> sections to indicate where the problem occurred, so the developer can decide which code is still relevant and which is faulty.

```
<<<<< HEAD:index.html
<div id="footer">todo: message here</div>
=====
<div id="footer">
    thanks for visiting our site
</div>
>>>>> SpecialBranch:index.html
```

branch 1's version

branch 2's version

Remote operations

- **Clone** takes a repository from a server and makes a local repository at given destination

```
git clone git://git.kernel.org/pub/scm/.../linux.git my-linux
```

- **Push** takes the local changes and writes them to the server.

```
git pull origin master
```

- **Pull** gets most recent changes from the repository.

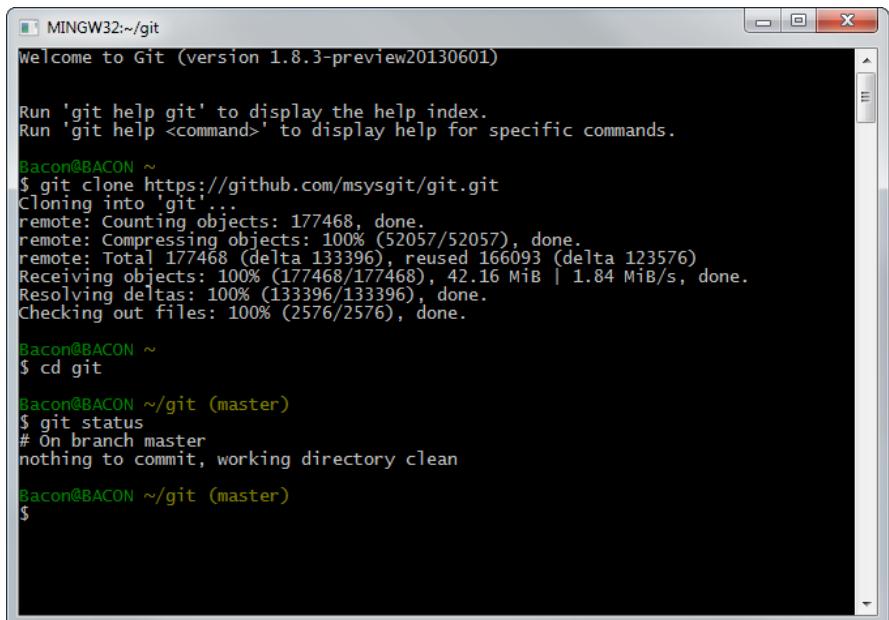
```
git pull origin master
```

Git – Workflow

- 1. Create a branch:** Topic branches created from the canonical deployment branch (usually master) allow teams to contribute to many parallel efforts. Short-lived topic branches, in particular, keep teams focused and results in quick ships.
- 2. Add commits:** Snapshots of development efforts within a branch create safe, revertible points in the project's history.
- 3. Open a pull request:** Pull requests publicize a project's ongoing efforts and set the tone for a transparent development process.
- 4. Discuss and review code:** Teams participate in code reviews by commenting, testing, and reviewing open pull requests. Code review is at the core of an open and participatory culture.
- 5. Merge:** Upon clicking merge, GitHub automatically performs the equivalent of a local 'git merge' operation. GitHub also keeps the entire branch development history on the merged pull request.

Git Bash vs Git GUI

- Less intuitive
- More technical approach
- Better control



```
MINGW32:~/git
Welcome to Git (version 1.8.3-preview20130601)

Run 'git help git' to display the help index.
Run 'git help <command>' to display help for specific commands.

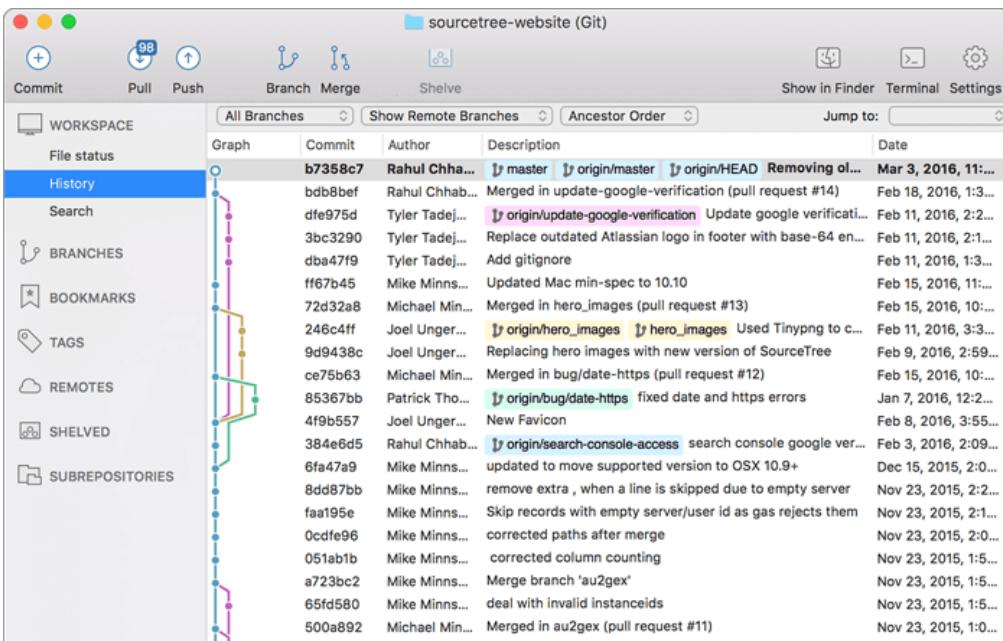
Bacon@BACON ~
$ git clone https://github.com/msysgit/git.git
Cloning into 'git'...
remote: Counting objects: 177468, done.
remote: Compressing objects: 100% (52057/52057), done.
remote: Total 177468 (delta 133396), reused 166093 (delta 123576)
Receiving objects: 100% (177468/177468), 42.16 MiB | 1.84 MiB/s, done.
Resolving deltas: 100% (133396/133396), done.
Checking out files: 100% (2576/2576), done.

Bacon@BACON ~
$ cd git

Bacon@BACON ~/git (master)
$ git status
# On branch master
nothing to commit, working directory clean

Bacon@BACON ~/git (master)
$
```

- Easy to start
- No need to know commands
- Blackbox



Git with GUI demo



- Git server repository.
- Ticketing system.
- Tickets (Issues) can be connected with branches – the issues can be discussed in the tickets.
- Assigning tickets to people. Issues labeling.
- When the developing is complete a merge request is made and the code is assigned to a person non-active in the issue fix. The person reviews the code and merge the changes into the production branch.
- Milestones.
- CI/CD
- ...

Git with GUI demo

1. Create a gitlab repository
2. Clone the repository
3. Add files
4. Push
5. Issue with branch and merge request
6. Switch to branch, add a file and change some of the already present
7. Stage and commit
8. Switch between branches (how working directory behaves)
9. Push changes
10. Merge and close the issue/branch

Git – Essential Commands

GIT Command	Description
git init	Create empty Git repo in specified directory. Run with no arguments to initialize the current directory as a git repository
git clone <repo>	Clone repo located at <repo> onto local machine. Original repo can be located on the local filesystem or on a remote machine via HTTP or SSH.
git branch	List all of the branches in your repo. Add a <branch> argument to create a new branch with the name <branch>.
git checkout -b <branch>	Create and check out a new branch named <branch>. Drop the -b flag to checkout an existing branch.
git add <directory>	Stage all changes in <directory> for the next commit. Replace <directory> with a <file> to change a specific file.
git commit -m "<message>"	Commit the staged snapshot, but instead of launching a text editor, use <message> as the commit message.
git status	List which files are staged, unstaged, and untracked.
git pull <remote>	Fetch the specified remote's copy of current branch and immediately merge it into the local copy.
git push <remote> <branch>	Push the branch to <remote>, along with necessary commits and objects. Creates named branch in the remote repo if it doesn't exist.

GIT Cheat Sheet: <https://www.atlassian.com/dam/jcr:8132028b-024f-4b6b-953e-e68fcce0c5fa/atlassian-git-cheatsheet.pdf>

GIT Official Documentation: <https://git-scm.com/docs/git-pull>

Git in CMD demo

1. git status, git branch in the repository
2. modify a file, add a file
3. git status
4. git commit +msg
5. git push
6. show the results on

Tasks

- Go to <https://learngitbranching.js.org/> and complete the introduction sequence

Git – Takeaway message

- Git is good for you, even though it may seem that it only adds work, in the end having track of all the changes and being able to develop parallelly outweighs the drawbacks by far.
- Helps with Quality Assurance workflow.
- Git is a versioning system vs GitLab/GitHub are cloud repository storage with further functions.
 - You do not have to be on cloud to use git.
- You do not need to work in team to benefit from git – you can use it solely to version your work.
 - One branch, just committing changes in time.
- It may seem overwhelming at first, but it is fairly simple once you get used to it. ☺

References:

- <https://git-scm.com/>
- <https://courses.cs.washington.edu/courses/cse403/13au/lectures/git.ppt.pdf>
- <https://www.slideshare.net/naimlatifi/gitpresentation-140814102916phpapp01>

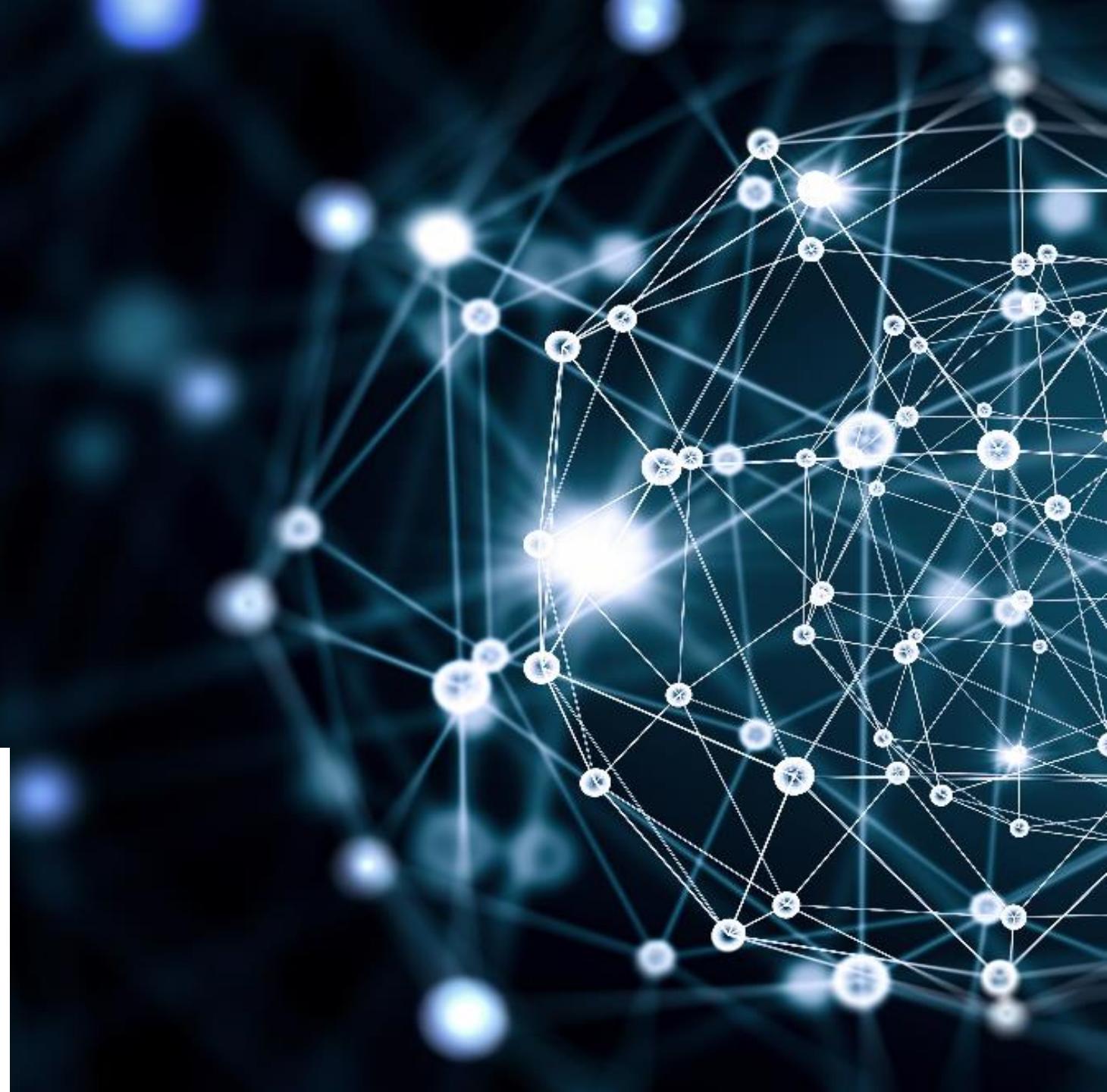
Pure Python Data Exploration



Pure Python Data Exploration

1. Upload the google playstore datasets
2. Open jupyter notebook D1_google_apps_empty.ipynb

Python Libraries



Python Libraries introduction

1. Open the D1_Python_libraries.ipynb notebook

Home Assignment HA1

Pick one of the following

1. Scrape a website with product you want to of your choosing and output a list of relevant information about the offers (e.g. <https://www.sbazar.cz/hledej/macbook?cena-dohodou=bez> , <https://www.mmreality.cz/nemovitosti/> , ...) (hint: BeautifulSoup)
2. For a list of stock symbols, get the last and future dividend payouts (e.g. <https://www.nasdaq.com/market-activity/stocks/t/dividend-history>) (hint: BeautifulSoup or find a free API)
3. Create a script (function) that states size of all files inside recursively and returns the path to the largest file. Raise error if given path is not a directory. Do this on local python installation in order to test it. (hint: os)
4. Any other ideas are welcome.

Save your script/notebook and send it to pmilicka@deloitteCE.com with subject **DSI_HA01_<surname>** for review.



Deloitte refers to one or more of Deloitte Touche Tohmatsu Limited ("DTTL"), its global network of member firms, and their related entities (collectively, the "Deloitte organization"). DTTL (also referred to as "Deloitte Global") and each of its member firms and related entities are legally separate and independent entities, which cannot obligate or bind each other in respect of third parties. DTTL and each DTTL member firm and related entity is liable only for its own acts and omissions, and not those of each other. DTTL does not provide services to clients. Please see www.deloitte.com/about to learn more.

Deloitte provides industry-leading audit and assurance, tax and legal, consulting, financial advisory, and risk advisory services to nearly 90% of the Fortune Global 500® and thousands of private companies. Our professionals deliver measurable and lasting results that help reinforce public trust in capital markets, enable clients to transform and thrive, and lead the way toward a stronger economy, a more equitable society and a sustainable world. Building on its 175-plus year history, Deloitte spans more than 150 countries and territories. Learn how Deloitte's more than 415,000 people worldwide make an impact that matters at www.deloitte.com.

This communication contains general information only, and none of Deloitte Touche Tohmatsu Limited ("DTTL"), its global network of member firms or their related entities (collectively, the "Deloitte organization") is, by means of this communication, rendering professional advice or services. Before making any decision or taking any action that may affect your finances or your business, you should consult a qualified professional adviser.

No representations, warranties or undertakings (express or implied) are given as to the accuracy or completeness of the information in this communication, and none of DTTL, its member firms, related entities, employees or agents shall be liable or responsible for any loss or damage whatsoever arising directly or indirectly in connection with any person relying on this communication. DTTL and each of its member firms, and their related entities, are legally separate and independent entities.