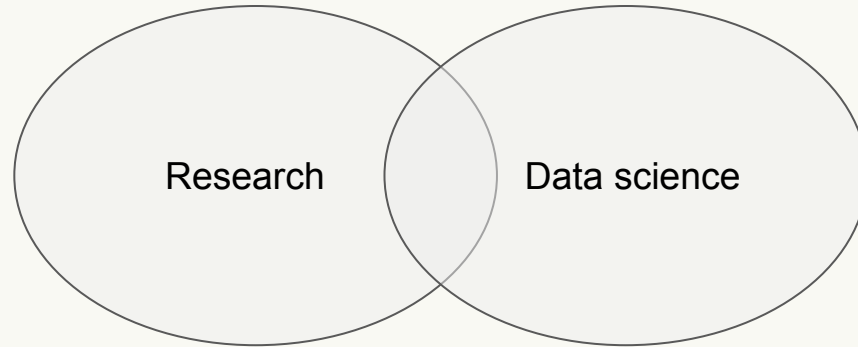# Nice to meet you

- Data scientist in Melbourne

# Nice to meet you

- Data scientist in Melbourne

- PhD in astrophysics at the University of Groningen

# How did I become a data scientist?

Research     Data science

# Why does a data scientist use Python?

Python can support the entire data science lifecycle!

Data wrangling
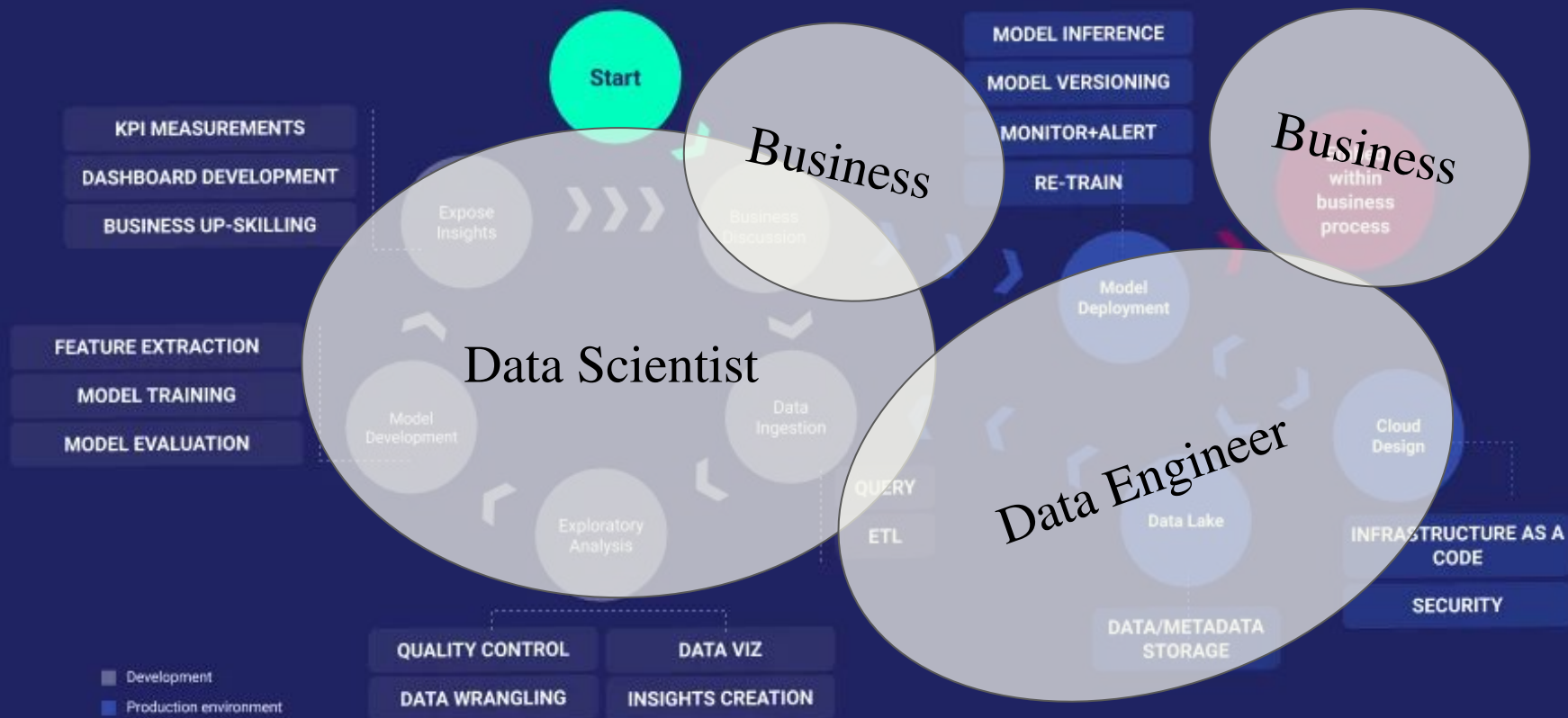
Statistical analysis

Visualisations

ML models

Deployment

Tools

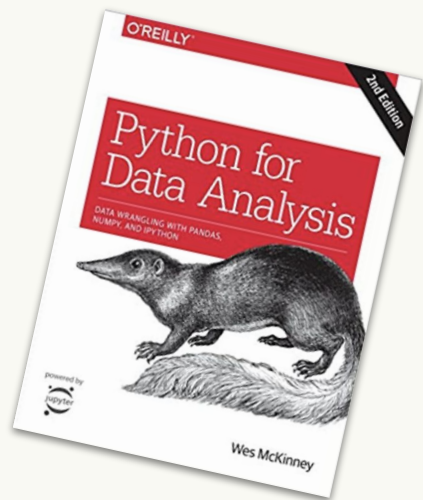# Data cleaning and feature extraction

**pandas**

- DataFrame object for data manipulation

- Reading and writing

- Label-based computations

- Missing data, reshaping

- Group by, merge, concat

- Time-series functionality

# Data cleaning and feature extraction



```
[8]:   # Import pandas
       import pandas as pd

       # Read in dataset
       df = pd.read_csv('kaggle_datasets/Churn_Modelling.csv')
       df.head()
```

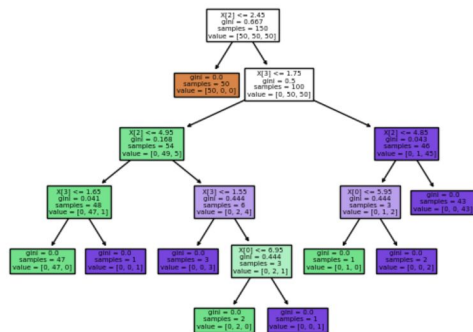| [8]: | | CreditScore | Geography | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | EstimatedSalary | Exited |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 619 | France | 42 | 2 | 0.00 | 1 | 1 | 1 | 101348.88 | 1 |
| | 1 | 608 | Spain | 41 | 1 | 83807.86 | 1 | 0 | 1 | 112542.58 | 0 |
| | 2 | 502 | France | 42 | 8 | 159660.80 | 3 | 1 | 0 | 113931.57 | 1 |
| | 3 | 699 | France | 39 | 1 | 0.00 | 2 | 0 | 0 | 93826.63 | 0 |
| | 4 | 850 | Spain | 43 | 2 | 125510.82 | 1 | 1 | 1 | 79084.10 | 0 |

# Machine learning



- Simple and efficient tools for predictive data analysis

- Built on NumPy, SciPy, and matplotlib

- Open source, commercially usable

# Machine learning



```
>>> from sklearn.datasets import load_iris
>>> from sklearn import tree
>>> X, y = load_iris(return_X_y=True)
>>> clf = tree.DecisionTreeClassifier()
>>> clf = clf.fit(X, y)
```
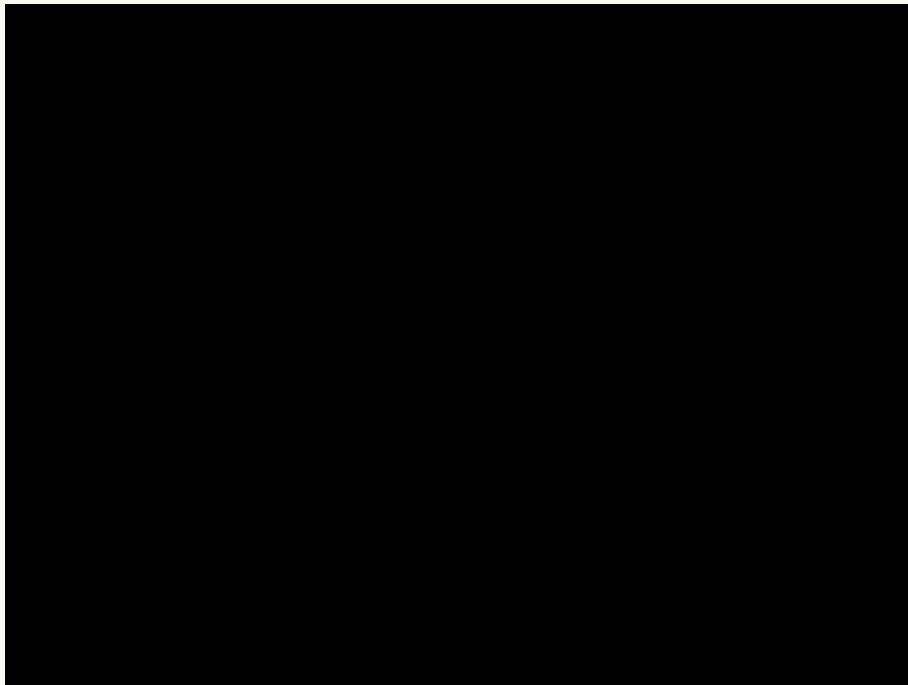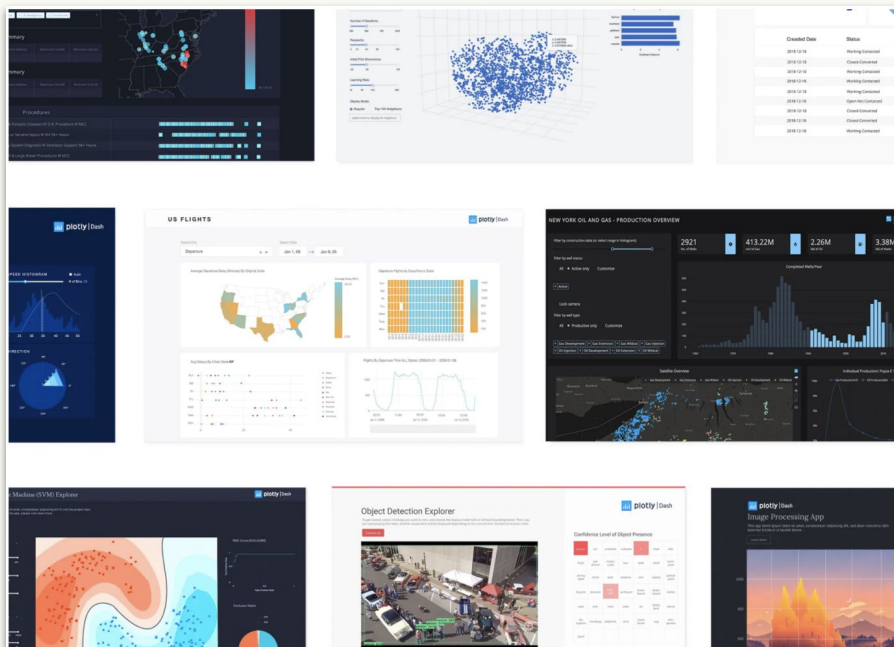
```
>>> tree.plot_tree(clf)
```

# Data visualisation

plotly

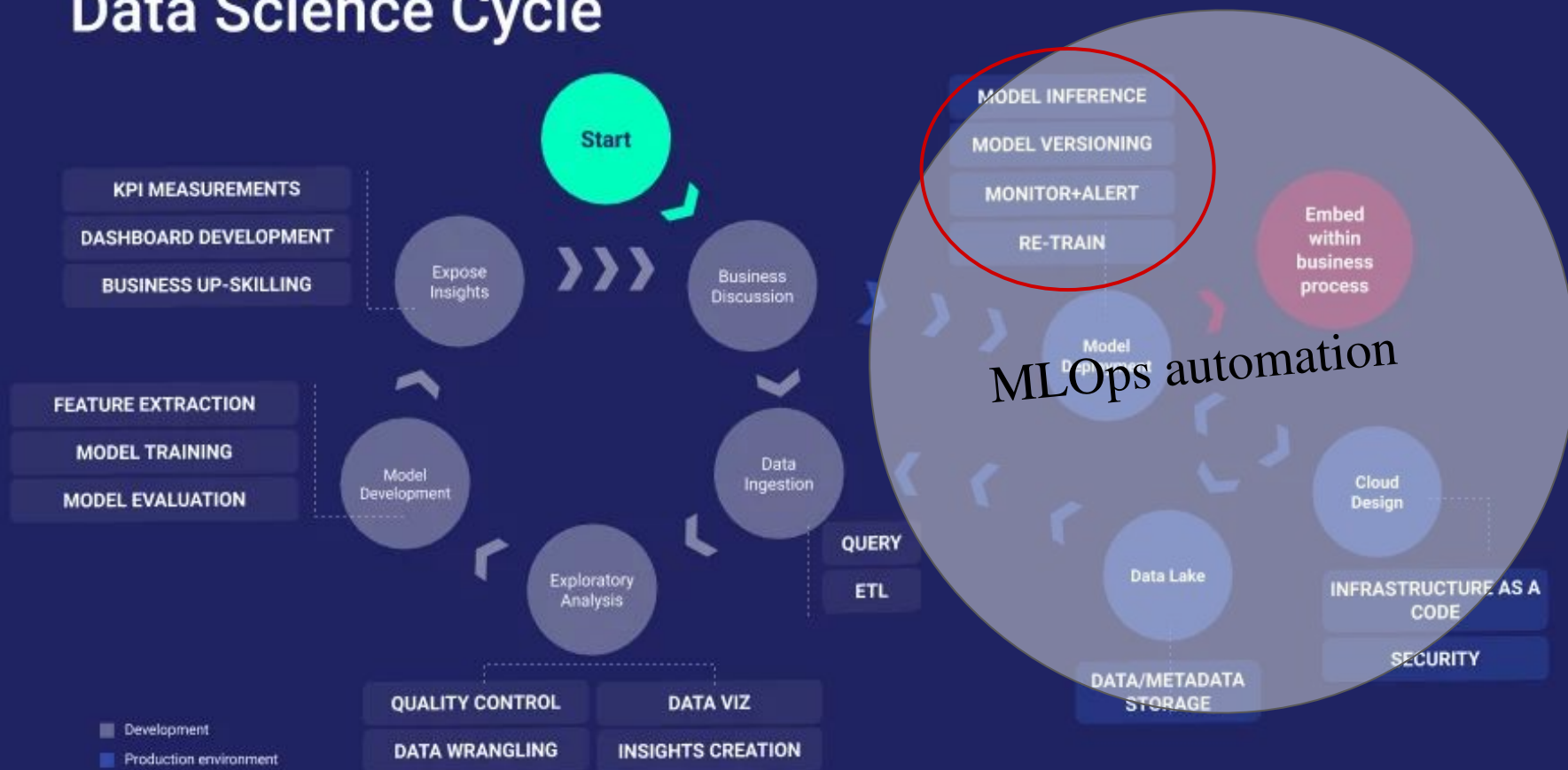- Interactive visualisations

# Web apps



- Build and deploy apps

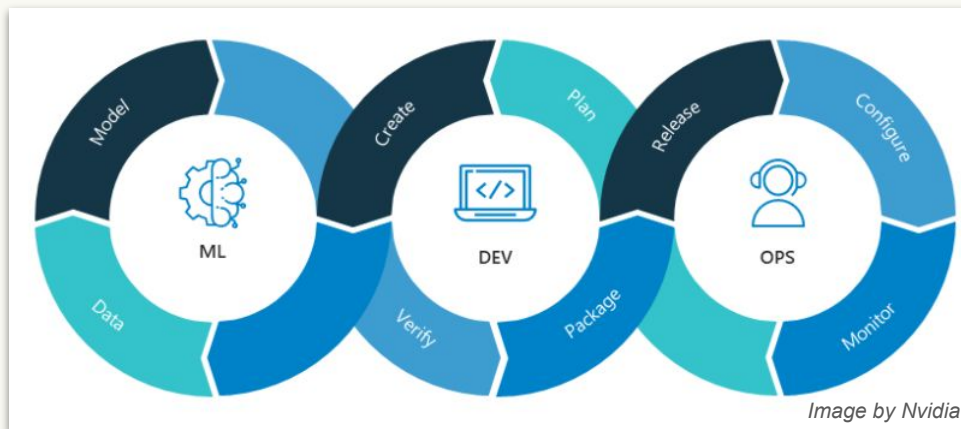# Towards machine learning automation

# MLOps



Image by Nvidia

- Trackable
- Reproducible
- Self-sustaining
- Automated

# MLOps in the cloud

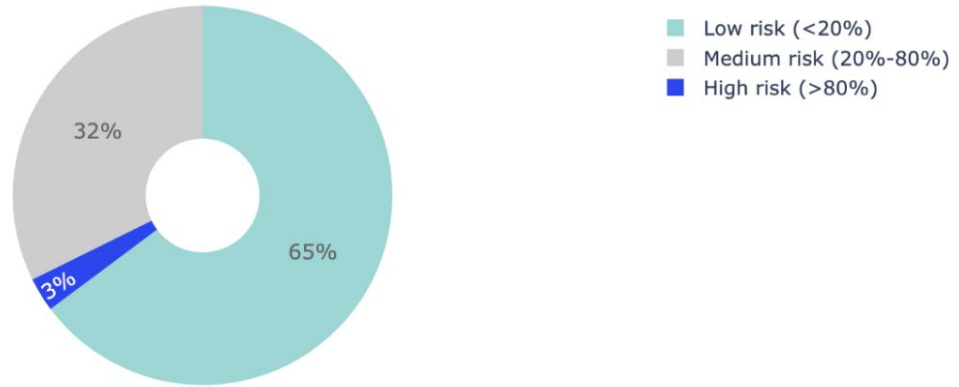# Use cases

# Churn Prediction

```
df = pd.read_csv('Churn_Modelling.csv')
df.head()
```

| CreditScore | Geography | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | EstimatedSalary | Exited |
|---|---|---|---|---|---|---|---|---|---|
| 619 | France | 42 | 2 | 0.00 | 1 | 1 | 1 | 101348.88 | 1 |
| 608 | Spain | 41 | 1 | 83807.86 | 1 | 0 | 1 | 112542.58 | 0 |
| 502 | France | 42 | 8 | 159660.80 | 3 | 1 | 0 | 113931.57 | 1 |
| 699 | France | 39 | 1 | 0.00 | 2 | 0 | 0 | 93826.63 | 0 |
| 850 | Spain | 43 | 2 | 125510.82 | 1 | 1 | 1 | 79084.10 | 0 |

# Churn Prediction

```python
# Train the logistic regression model
from sklearn.linear_model import LogisticRegression
model = LogisticRegression(random_state=0)
model.fit(X, y)
```

### Churn risk profiling



- Low risk (<20%)
- Medium risk (20%-80%)
- High risk (>80%)

32%

3%

65%

# Risk profiling for customer churn analysis

Katinka Gereb · Aug 1, 2020 · 6 min read



@katinka-gereb

# Clustering



*https://cs.stanford.edu/people/karpathy/cnnembed/*

# Customer 360

**Churn risk profiling**



- Low risk (<20%)
- Medium risk (20%-80%)
- High risk (>80%)

65%
32%
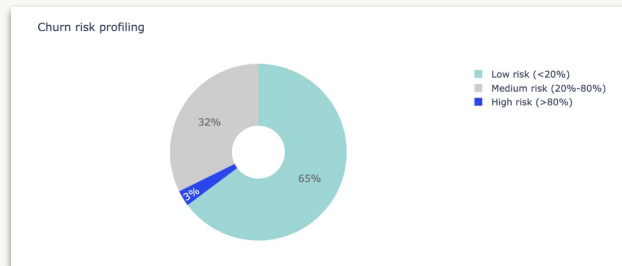3%



**Transactions:**
ID, date time, value, SKUs, transaction history)

**Customer data:**
Membership ID, industry, demographics, lifetime value, segmentation metrics, churn risk profile

**Loyalty:**
Discounts, rewards, redemption history

**Product data:**
SKUs, product categories, price, purchase frequency, volume forecasts

**Customer Insights**

**External data:**
- Government data
- Competitor data
- Ad-hoc, such as COVID-19 data

**Payment/credit:**
Date and time, credit rating and risk, debt history, fraud predictions

**Customer interactions:**
Call center data, chatbot interactions, emails, social media posts with # reference

**Marketing data:**
Click-through analysis, response analytics, campaigns, offers, cohort analysis
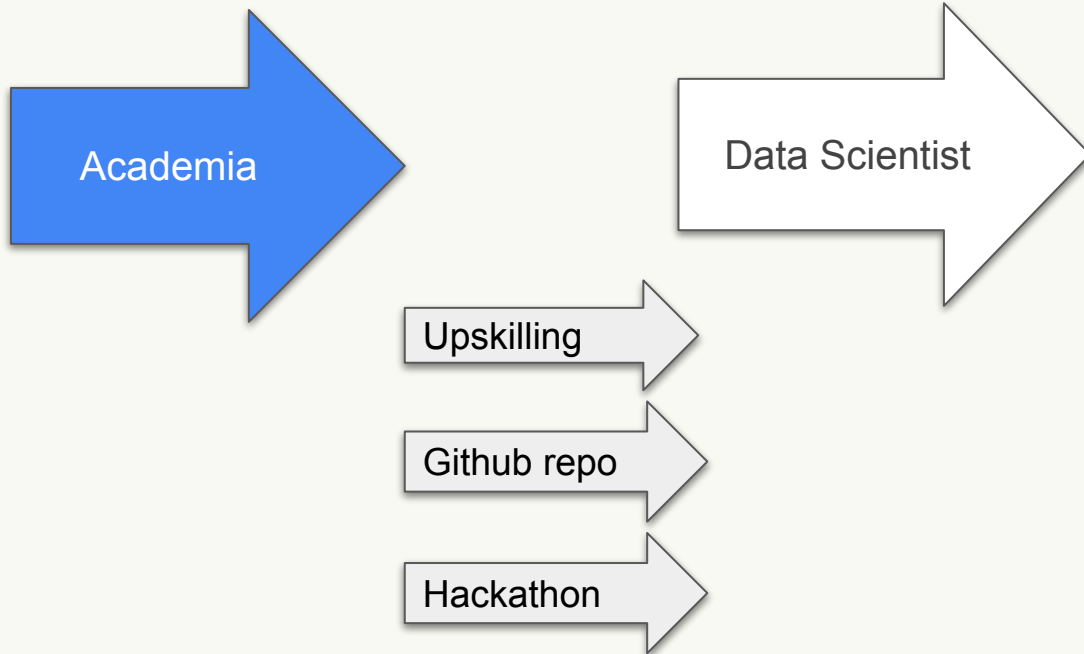
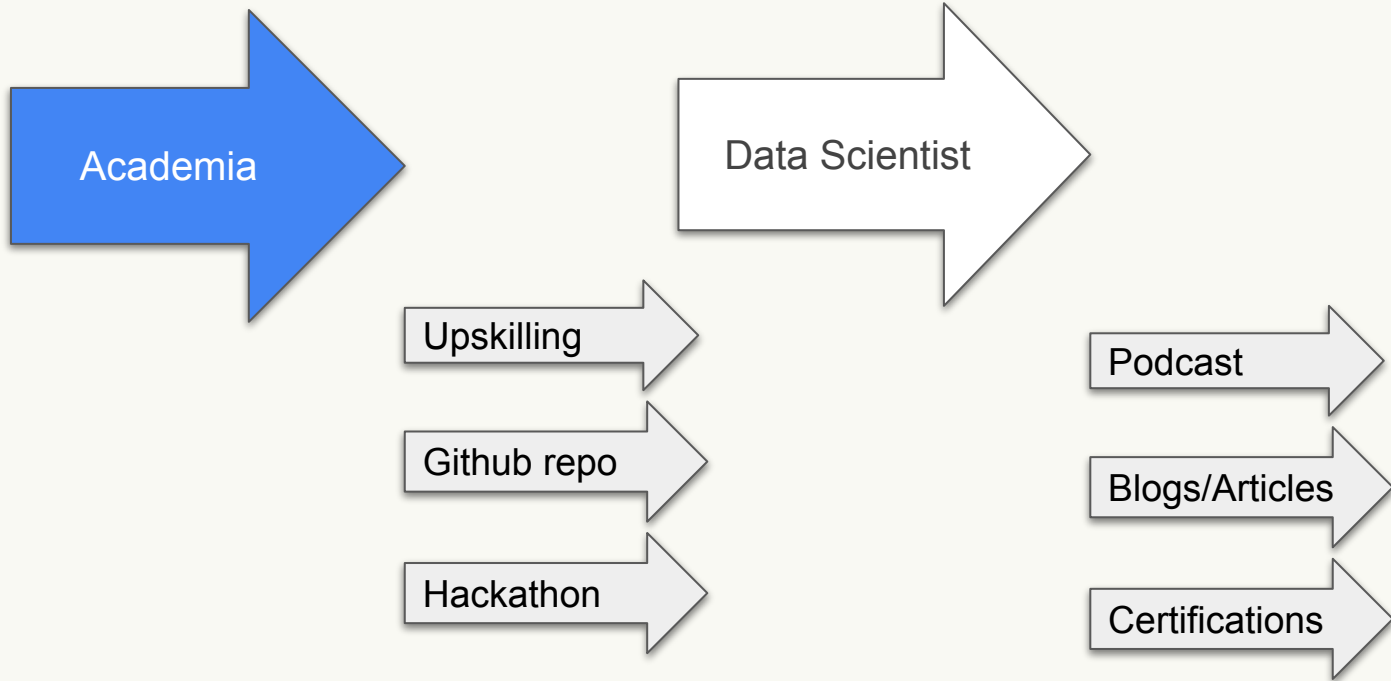# Advancing your career


Academia


Data Scientist

# Advancing your career

# Advancing your career

@katinka-gereb

# Thanks for listening!



Nova Tech Podcast

https://novatech.buzzsprout.com/