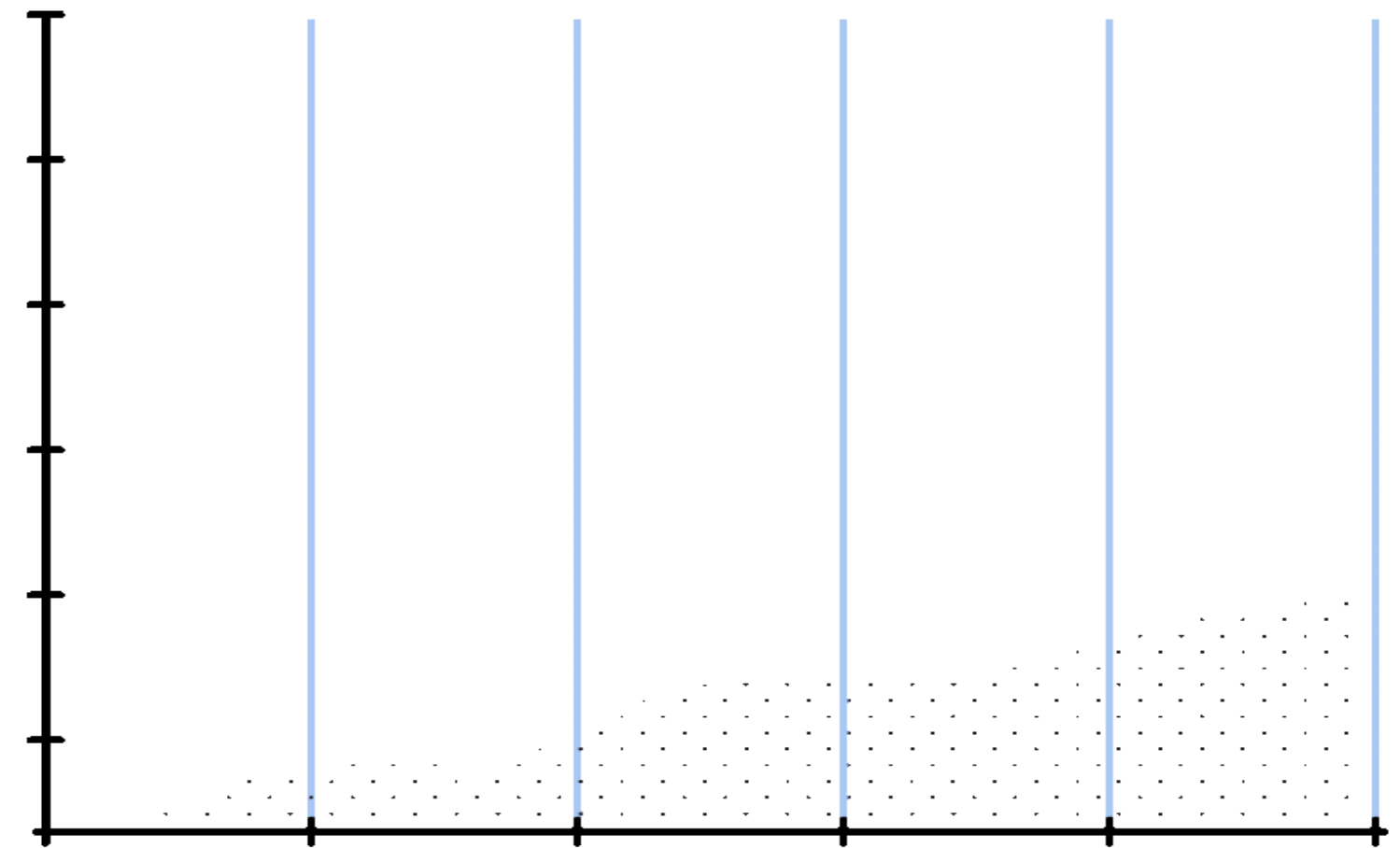# Data Science University

**Analyzing the trends in Data Science Industry**

Institutional Goal :

"Empowering the next generation of Analytics leaders through education and guidance."

BY THYNK ANALYTICS

# OUR TEAM

**Anagh Shrivastava**
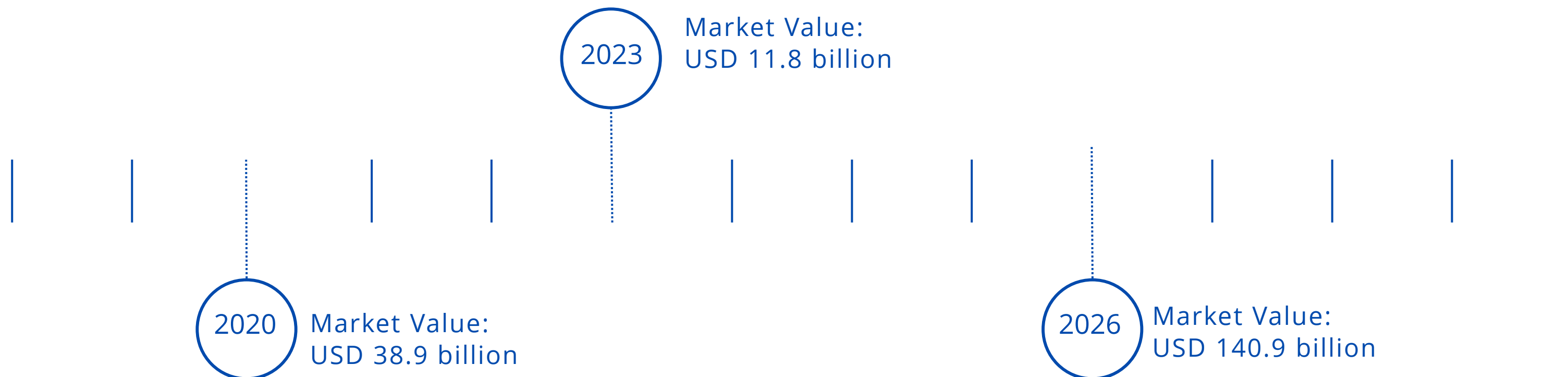
**Khyati Desai**

**Kritika Dhingra**

Parth Aloni

**Sara Dhoot**

# EXECUTIVE SUMMARY

- Examining trends in the analytics profession to provide insights/recommendations to an educational institution on survey responses from Kaggle Survey 2022 .

- It recorded over 25,000 data science practitioner responses which includes students and professionals.

- Conducted data cleaning/preparation and statistical analysis

- Used the kNN-supervised machine learning algorithm to build a tool to predict salaries based on user input.

- The Analysis provided insights into the significant factors that drive the salaries of data science professionals, and what kind of technologies and tools are prevalent and emerging in the field.

- Based on predictive modeling and insights the team is providing the following recommendations to the institute -
  - Focus on Higher Education
  - Build industry-specific knowledge
  - Keep up with the most used and emerging technologies and tools in Analytics, ML, and AI fields.

# DATA SCIENCE

**The global data science market size was valued at USD 38.9 billion in 2020 and is expected to reach USD 140.9 billion by 2026, growing at a compound annual growth rate (CAGR) of 24.2% from 2021 to 2026***

2023    Market Value:
        USD 11.8 billion

2020    Market Value:
        USD 38.9 billion

2026    Market Value:
        USD 140.9 billion

# PROBLEM STATEMENT

To effectively prepare and counsel students interested in analytic professions, but they lack understanding of the latest trends.
Also, they need a tool to provide students with a projected salary based on their qualifications.
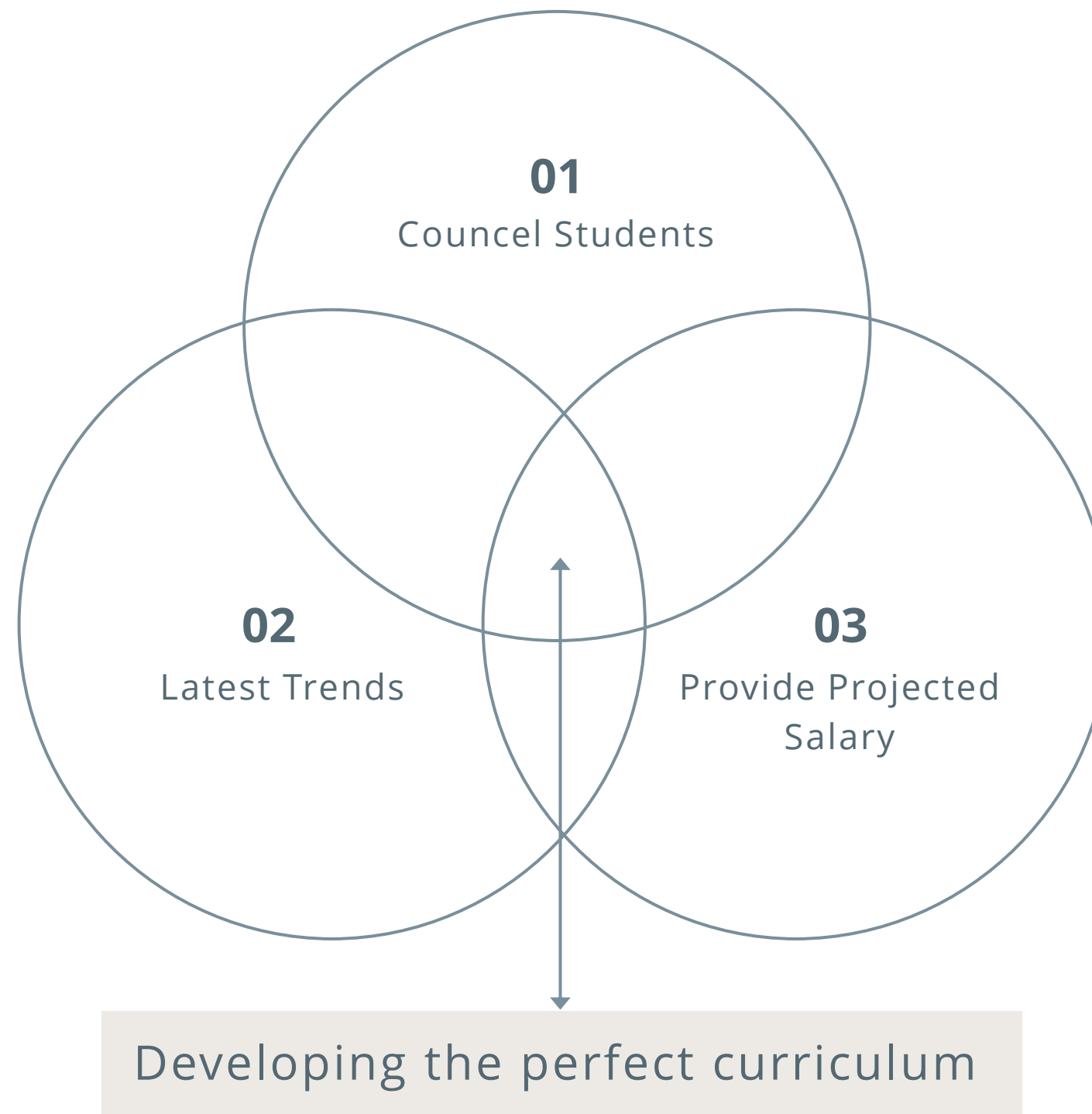
**Goal**:
To solve this and to examine the trends and create a tool that will help students prepare for the fast-growing field of data analytics.

## 75%

By the end of 2024, 75% of enterprises will shift from piloting to operationalizing AI, driving a 5X increase in streaming data and analytics infrastructures.

# PROBLEM ANALYSIS

**01**

Councel Students

**02**

Latest Trends

**03**

Provide Projected
Salary

Developing the perfect curriculum

# DATA ANALYSIS

To ensure that the dataset is representative for the analysis, the following criteria will be applied to include only responses from employed professionals:

- The respondents must **not currently be students**.
- The respondents must **currently be employed**.
- The respondents must have provided information on the **industry in which they are currently employed**, or their most recent employer if retired.
- The respondents must have answered the question regarding the **industry of their current employer or contract.**

The resulting dataset will include **approximately 37.9%** of the total responses and will be used for the analysis.

**1**

Data Cleaning

**2**

Data Manipulation

**3**

Data Analysis

# SIGNIFICANT FACTORS DRIVING DATA SCIENCE PROFESSIONALS' SALARIES?



- Country

- Education

- Experience

- Current Role

- Industry/ Segment

- Size of Organization

- Existing ML practices in Org.

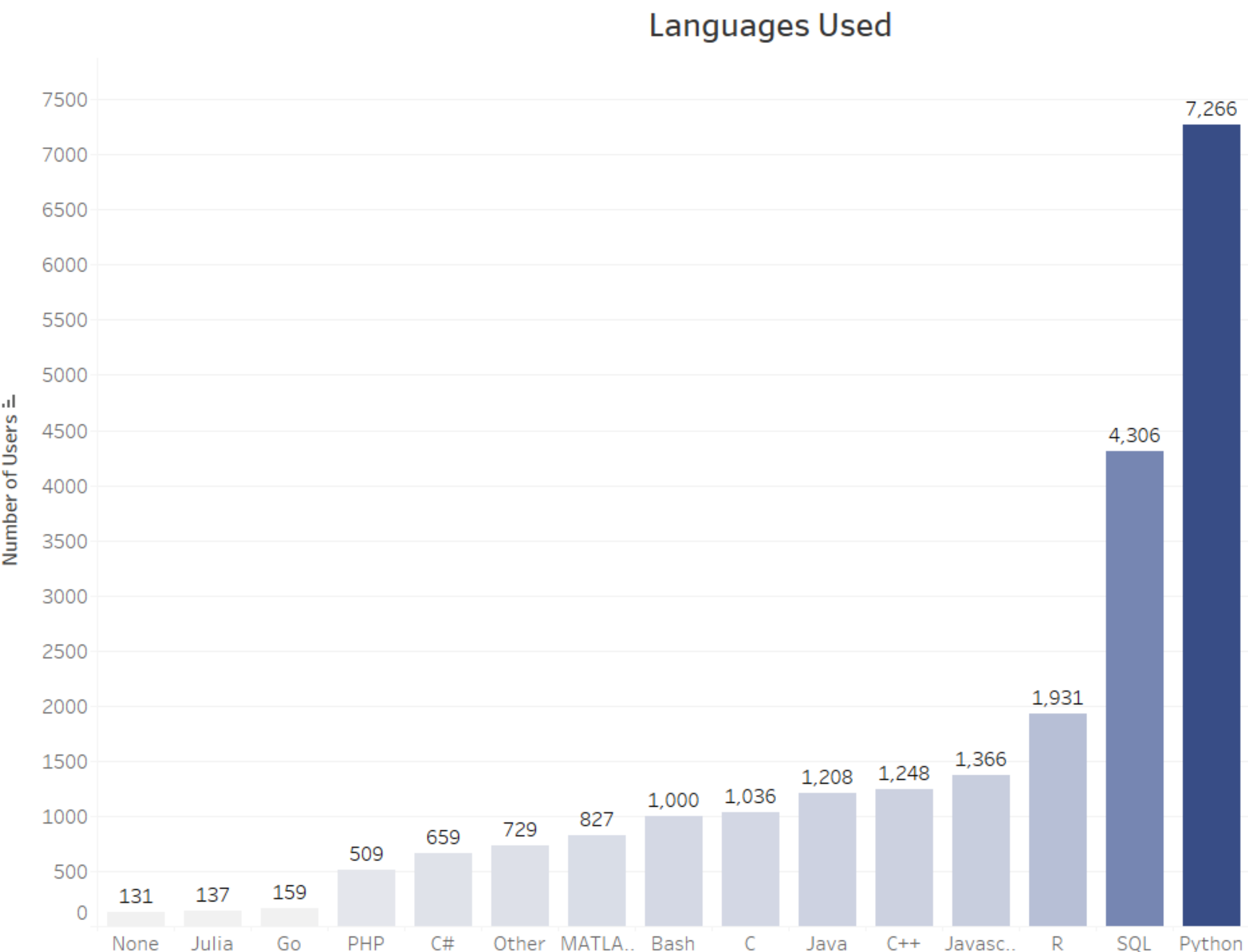- Investment in Learning ML, Cloud, AI

Data Science University

THE ROLES AND RESPONSIBILITES OF A DATA SCIENCE PROFESSIONAL:

## Important parts of the role:

1 Analyze and understand data to influence product or business decisions

2 Build prototypes to explore applying machine learning to new areas

3 Build and/or run a machine learning service that operationally improves my product or workflows

4 Experimentation and iteration to improve existing ML models
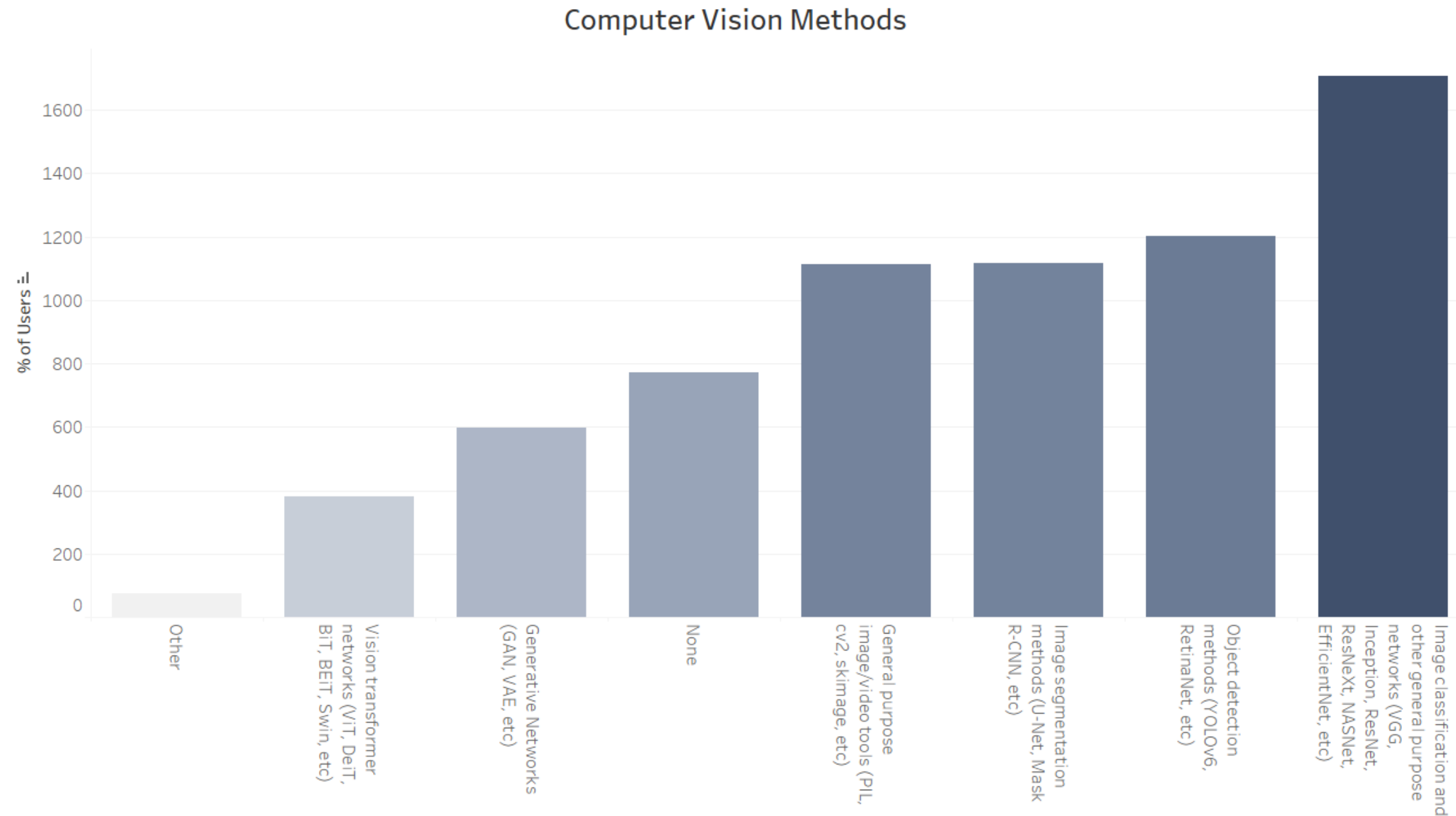
# PREVELANT TOOLS AND TECHNIQUES

Languages Used

**Most Used Languages in the field of Data Science:**

- Python
- SQL
- R

**Why to use these languages:**
Python, R, and SQL are popular choices in the field of data science due to their versatility, functionality, and ability to work together to analyze and manipulate large datasets.

# MOST USED COMPUTER VISION METHODS
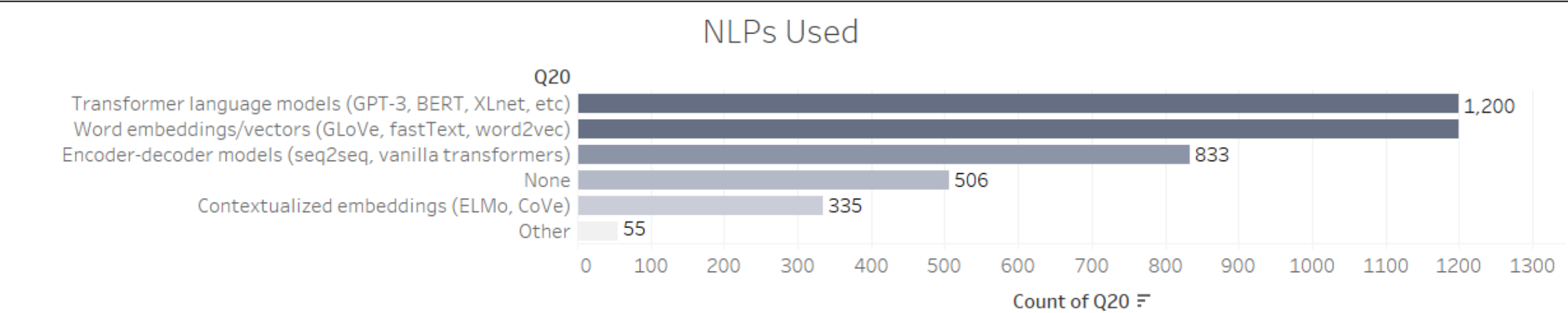


Computer Vision Methods

## Why to use computer vision methods?

Computer vision techniques can generate new data from existing datasets, improving accuracy of machine learning models with limited data.
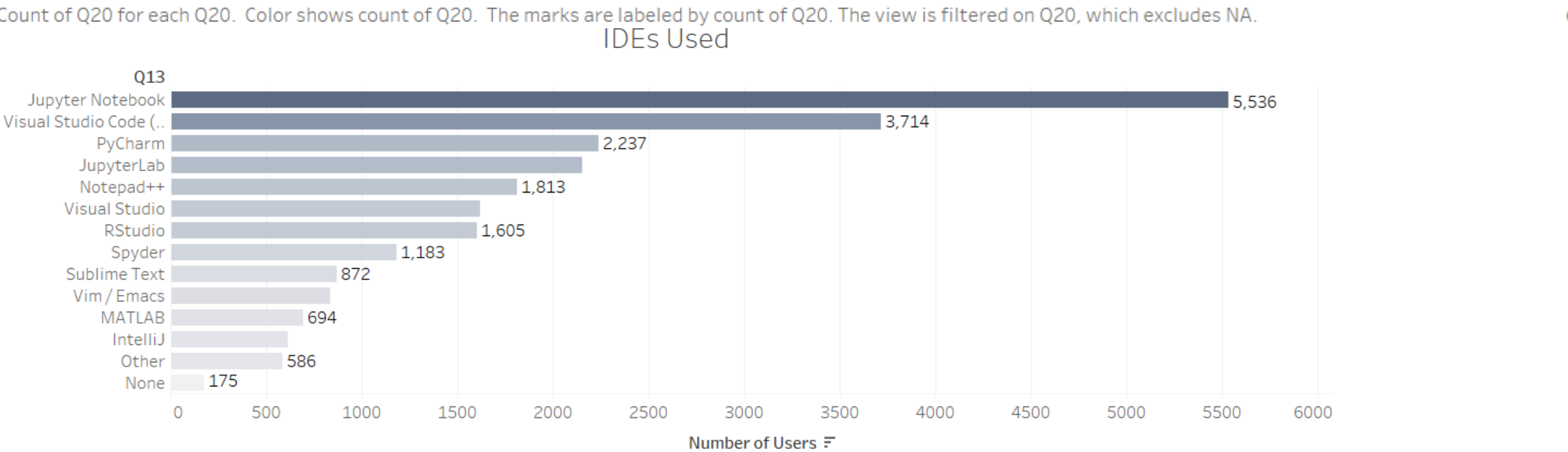
Most used:

Image classification and other general purpose networks (VGG, Inception, ResNet, ResNeXt, NASNet, EfficientNet, etc)
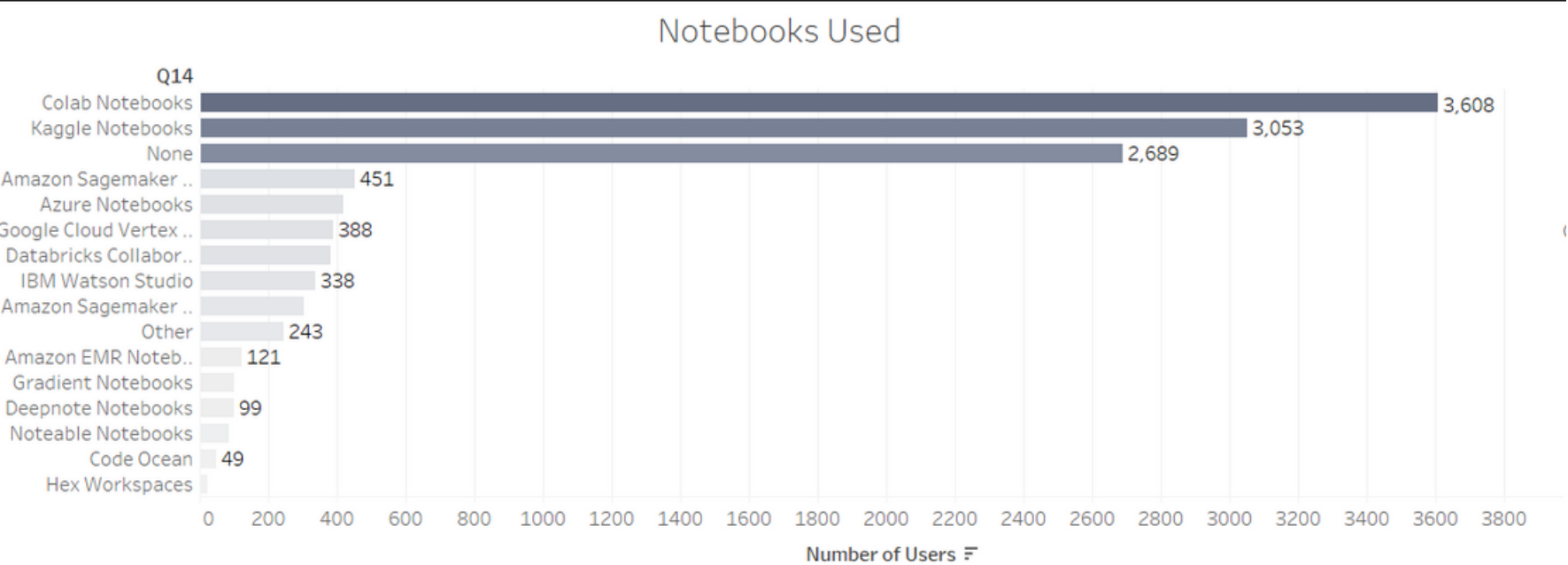
# THE MOST POPULAR USED TOOLS BY DATA SCIENCE PROFESSIONALS:

## NLPs Used

Q20

| NLP Tool | Count |
|---|---|
| Transformer language models (GPT-3, BERT, XLnet, etc) | 1,200 |
| Word embeddings/vectors (GLoVe, fastText, word2vec) | 1,200 |
| Encoder-decoder models (seq2seq, vanilla transformers) | 833 |
| None | 506 |
| Contextualized embeddings (ELMo, CoVe) | 335 |
| Other | 55 |

Count of Q20

Count of Q20 for each Q20.  Color shows count of Q20.  The marks are labeled by count of Q20. The view is filtered on Q20, which excludes NA.

## IDEs Used

Q13

| IDE | Number of Users |
|---|---|
| Jupyter Notebook | 5,536 |
| Visual Studio Code (.. | 3,714 |
| PyCharm | 2,237 |
| JupyterLab | |
| Notepad++ | 1,813 |
| Visual Studio | |
| RStudio | 1,605 |
| Spyder | 1,183 |
| Sublime Text | 872 |
| Vim / Emacs | |
| MATLAB | 694 |
| IntelliJ | |
| Other | 586 |
| None | 175 |

Number of Users

## Notebooks Used

Q14

| Notebook | Number of Users |
|---|---|
| Colab Notebooks | 3,608 |
| Kaggle Notebooks | 3,053 |
| None | 2,689 |
| Amazon Sagemaker .. | 451 |
| Azure Notebooks | |
| Google Cloud Vertex .. | 388 |
| Databricks Collabor.. | |
| IBM Watson Studio | 338 |
| Amazon Sagemaker .. | |
| Other | 243 |
| Amazon EMR Noteb.. | 121 |
| Gradient Notebooks | |
| Deepnote Notebooks | 99 |
| Noteable Notebooks | |
| Code Ocean | 49 |
| Hex Workspaces | |

Number of Users

- NLP: Transformer language models (GPT-3, BERT, XLnet, etc)

- IDE: Jupyter Notebook

- Notebook: Colab Notebooks

THYNK
ANALYTICS

# MACHINE LEARNING ALGORITHMS



## Why to use ML Algorithms?

### Identify:

- Patterns and trends in large datasets
- Automate processes
- accurate predictions

### Advantages:

- Businesses make data-driven decisions
- Improve efficiency
- Competitive Edge

**Most used:**
- Linear or Logistic Regression
- Decision Trees or Random Forests
- Gradient Boosting Machines

**THYNK
ANALYTICS**

# RETURN ON INVESTMENT

**How much a person has invested in their education and are they getting return on their investment:**
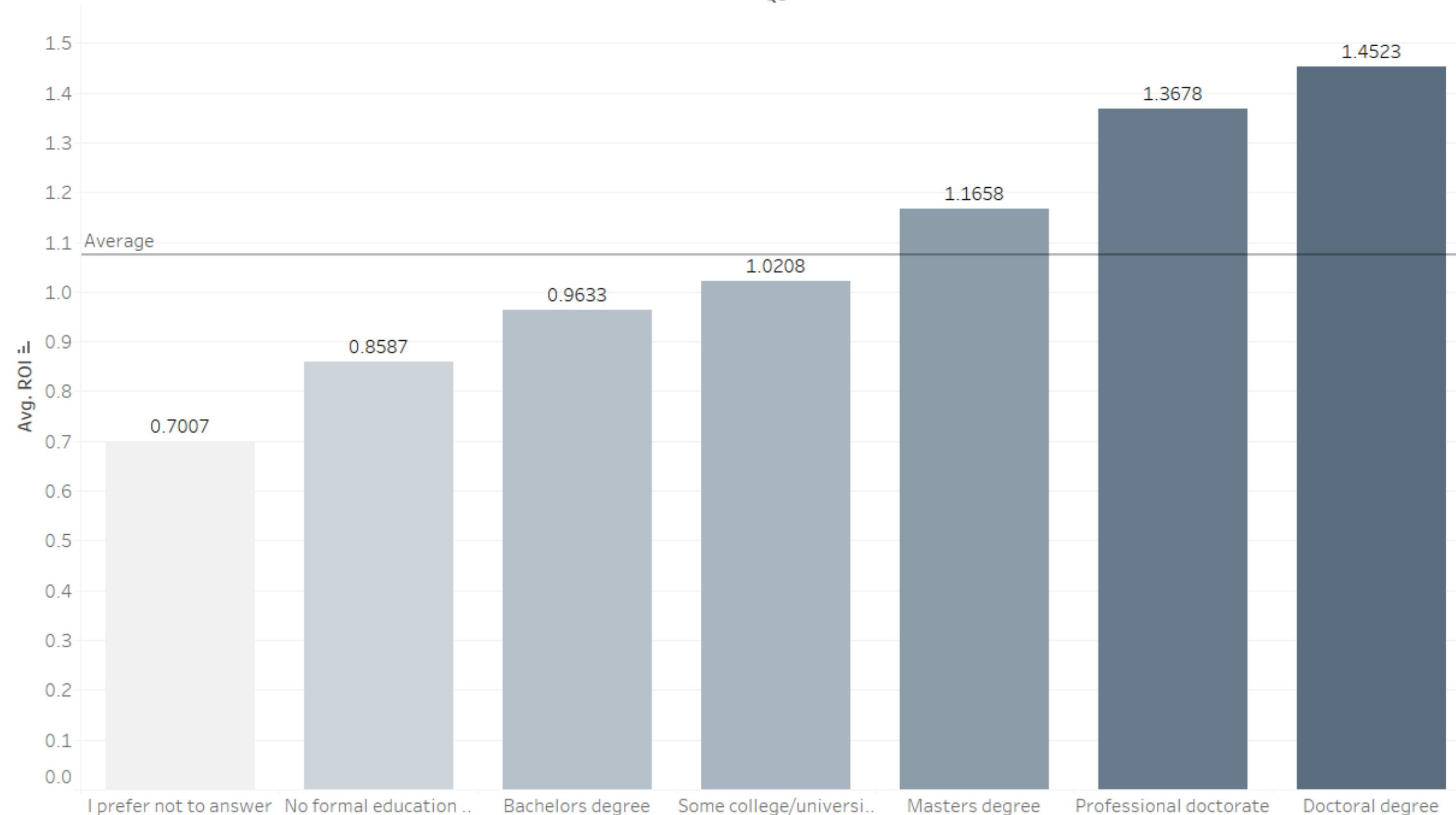
Return on Investment (ROI)

**ROI = Salary/Investment**

**THYNK
ANALYTICS**

## Degree v/s Return on Investment

Q8



**The ROI**

The professionals who have pursued :
- **Doctoral Degree**
- **Professional Degree**
- **Master's Degree**

Have the return on investment higher than the average return on investment in the market.

# PREDICTION MODEL PERFORMANCE

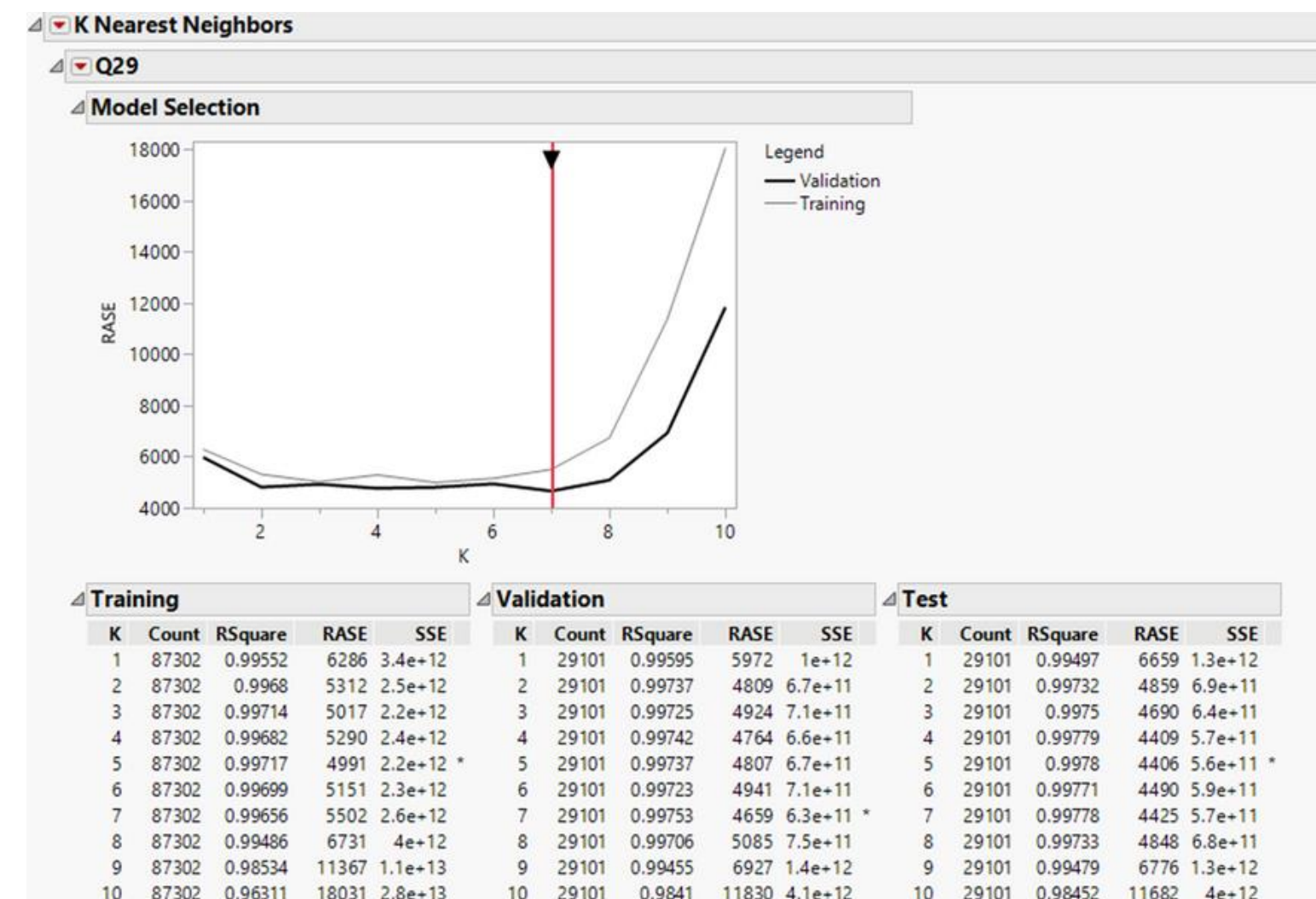The model is based on kNN (K-Nearest Neighbor) Algorithm of supervised Learning

Best Model fit with tuning parameter: **k = 7**

Model Performance is analyzed with statistical measures :
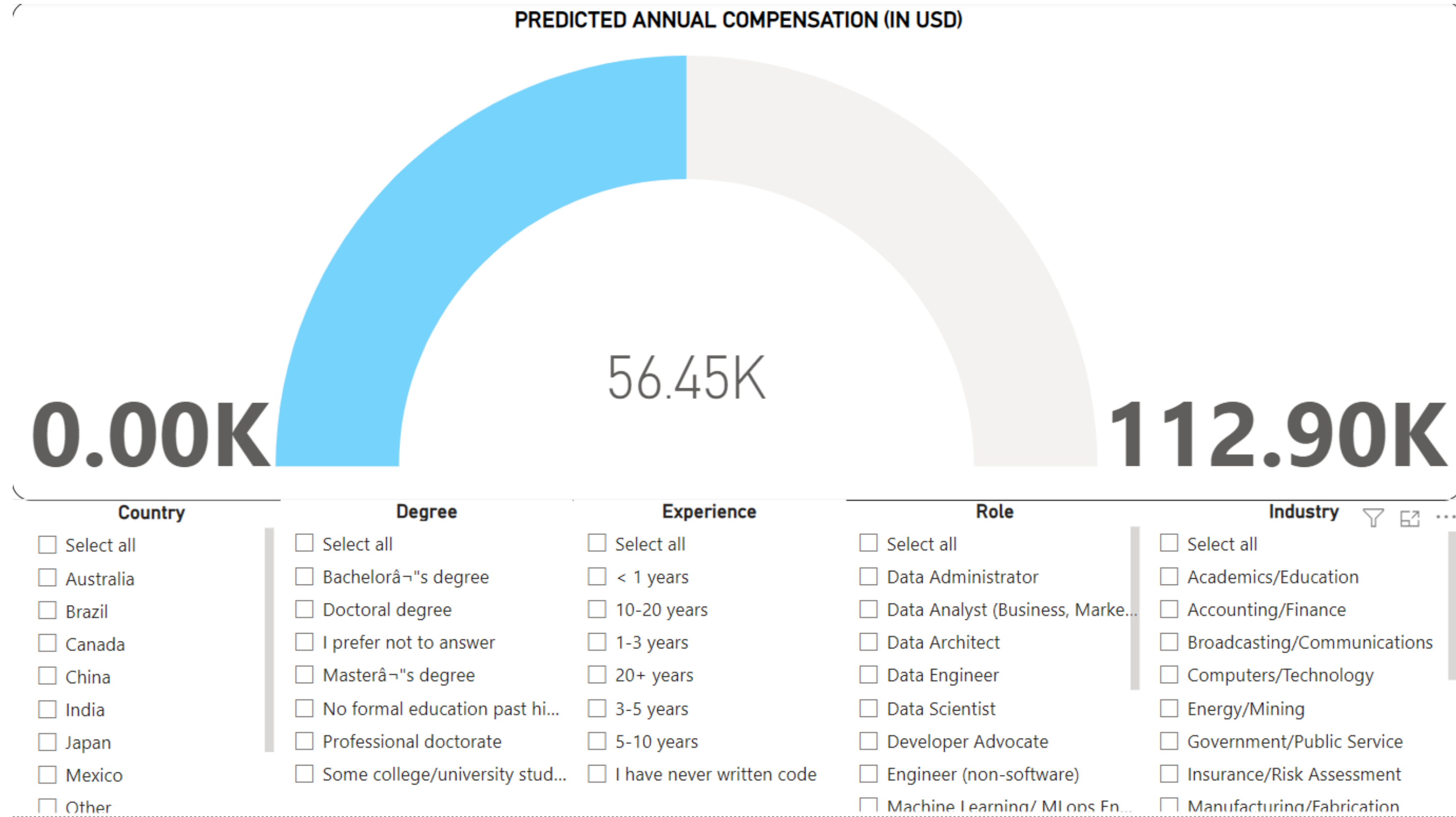
**r-square** = 0.99 (on the test set)

**RASE** = 4425

** Data only for current professionals

# Income Prediction Tool



PREDICTED ANNUAL COMPENSATION (IN USD)

56.45K

0.00K                                          112.90K

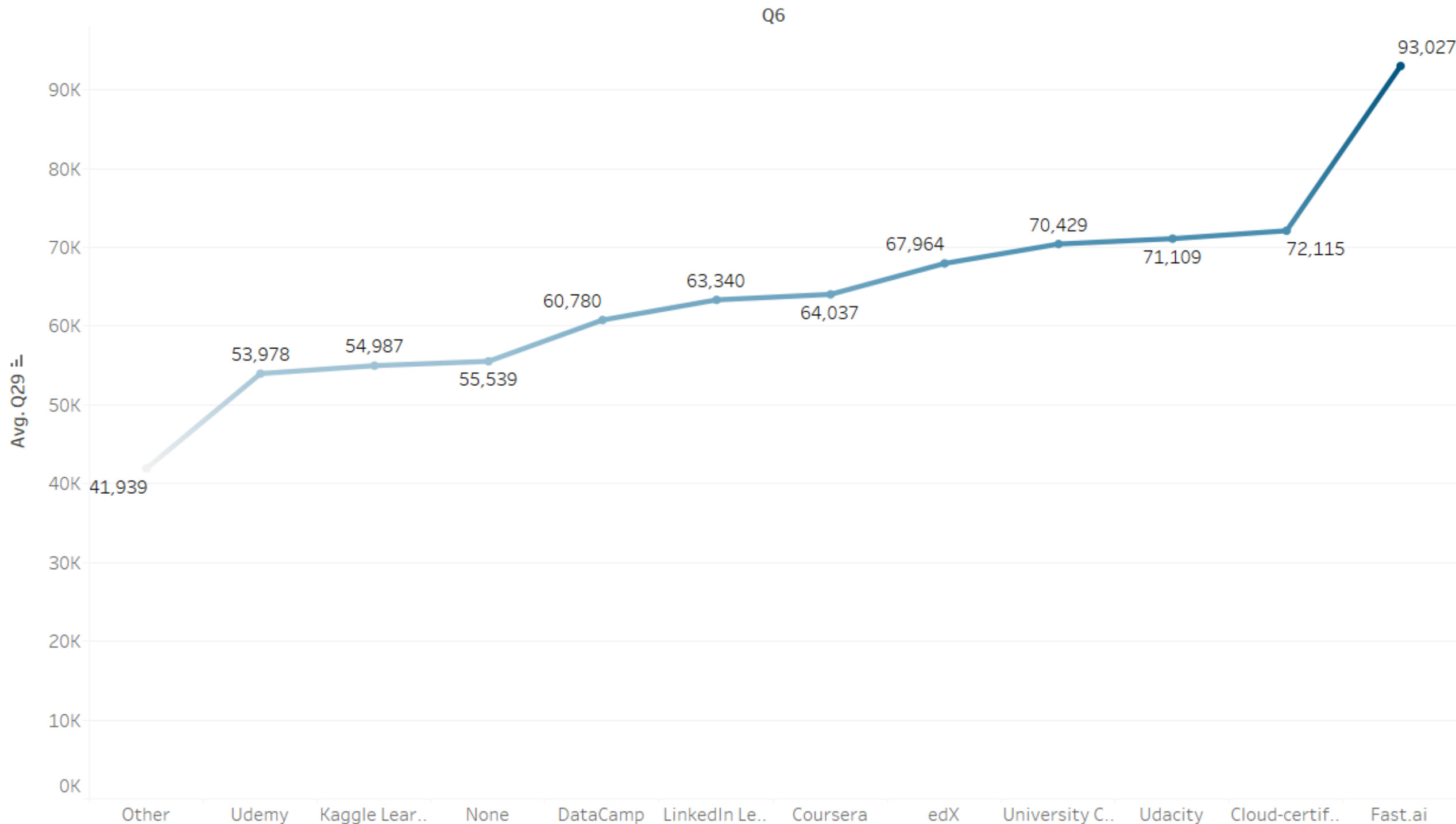| Country | Degree | Experience | Role | Industry |
|---------|--------|------------|------|----------|
| ☐ Select all | ☐ Select all | ☐ Select all | ☐ Select all | ☐ Select all |
| ☐ Australia | ☐ Bachelorâ¬"s degree | ☐ < 1 years | ☐ Data Administrator | ☐ Academics/Education |
| ☐ Brazil | ☐ Doctoral degree | ☐ 10-20 years | ☐ Data Analyst (Business, Marke... | ☐ Accounting/Finance |
| ☐ Canada | ☐ I prefer not to answer | ☐ 1-3 years | ☐ Data Architect | ☐ Broadcasting/Communications |
| ☐ China | ☐ Masterâ¬"s degree | ☐ 20+ years | ☐ Data Engineer | ☐ Computers/Technology |
| ☐ India | ☐ No formal education past hi... | ☐ 3-5 years | ☐ Data Scientist | ☐ Energy/Mining |
| ☐ Japan | ☐ Professional doctorate | ☐ 5-10 years | ☐ Developer Advocate | ☐ Government/Public Service |
| ☐ Mexico | ☐ Some college/university stud... | ☐ I have never written code | ☐ Engineer (non-software) | ☐ Insurance/Risk Assessment |
| ☐ Other | | | ☐ Machine Learning/ MLops En... | ☐ Manufacturing/Fabrication |

Link:
https://app.powerbi.com/links/IJ67FqhCAP?ctid=80f23f4a-91a4-4566-8db1-3bcabb21d1cb&pbi_source=linkShare

# EMERGING TOOLS TO INVEST TIME IN

## Tools and Techniques to Invest Time In

### Q6



The trend of average of Q29 for Q6. Color shows average of Q29. The marks are labeled by average of Q29. Details are shown for Q6. The view is filtered on Q6, which excludes NA and Null.

**The comparsion between the emerging online learning platforms using the average salaries in the data science industry:**

**Fast.ai:**
Eases building & training advanced ML models with pre-built models & tools for data preprocessing, visualization & interpretation. Built on PyTorch & TensorFlow.

# RECOMMENDATIONS:

**To design the curriculum holistically to integrate the use of the emerging technologies and tools.**

## Higher Education

Data science higher education offers valuable skills and knowledge with a high return on investment, making it a must for students seeking better job opportunities.

## Technologies

Focusing on the most used technologies in the industry, i.e Python, SQL, Jupyter, Collab, etc.
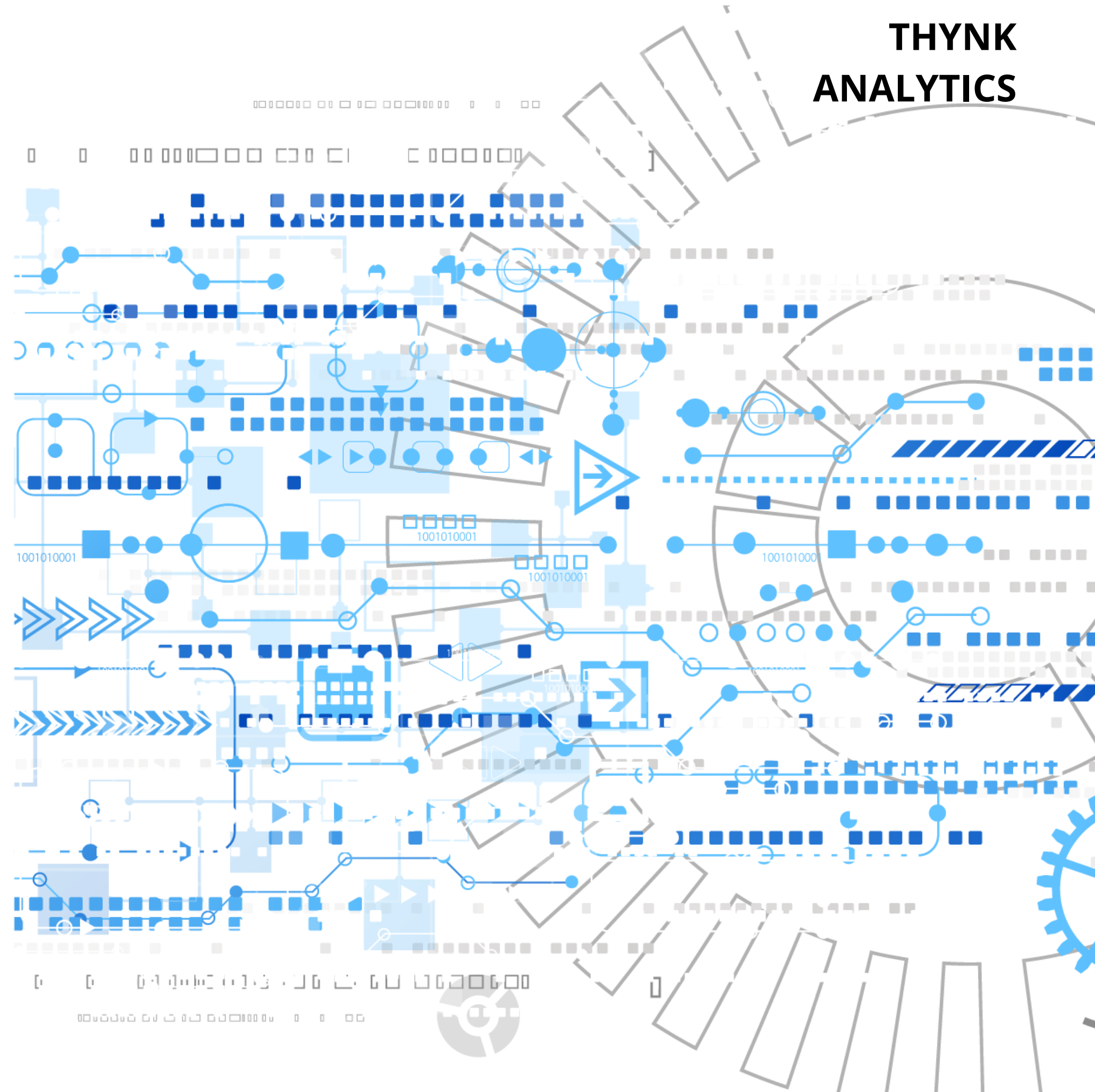
## Industry Specific Knowledge

. Students should be encouraged to focus on a major area of interest like education, computers/IT, Banking, etc. and build industry-specific knowledge.

# LIMITATIONS

- The dataset is biased as most of the responses are from India.

- The model has overfitted the data.

# Thank You