

# finalproject

2024-12-03

## Contribution Statement

Student 1:

Student 2:

## Introduction

Understanding the global trends and dynamics of the data science and machine learning (DS/ML) community is critical in an era where these fields are driving innovation across industries. The Kaggle Survey 2020 dataset offers a comprehensive glimpse into this domain, capturing over 20,000 responses from individuals worldwide. With 350 columns detailing diverse aspects such as professional roles, tools, learning resources, and industry practices, the dataset presents a unique opportunity to explore the multifaceted landscape of DS/ML.

Key questions guiding this study include:

1. Demographic and Economic Correlation: How do factors like age, gender, and country of residence influence income levels among DS/ML professionals? Additionally, what role does the highest level of education play in shaping these trends?
2. Programming Language Adoption: Are there observable patterns in programming language preferences across varying experience levels? How do these preferences vary by industry?
3. Impact of Cloud Computing: How does the adoption of advanced machine learning tools vary across different cloud platforms (e.g., AWS, GCP, Azure)? Which platforms show the highest proportion of advanced ML users, and what insights can this provide about user preferences or specialization trends in cloud services?
4. Effectiveness of Educational Resources: How do respondents' experience levels and compensation distributions vary across different educational resources (e.g., MOOCs, certifications, books)? Can we identify which resources are associated with higher compensation levels or a greater proportion of experienced professionals?

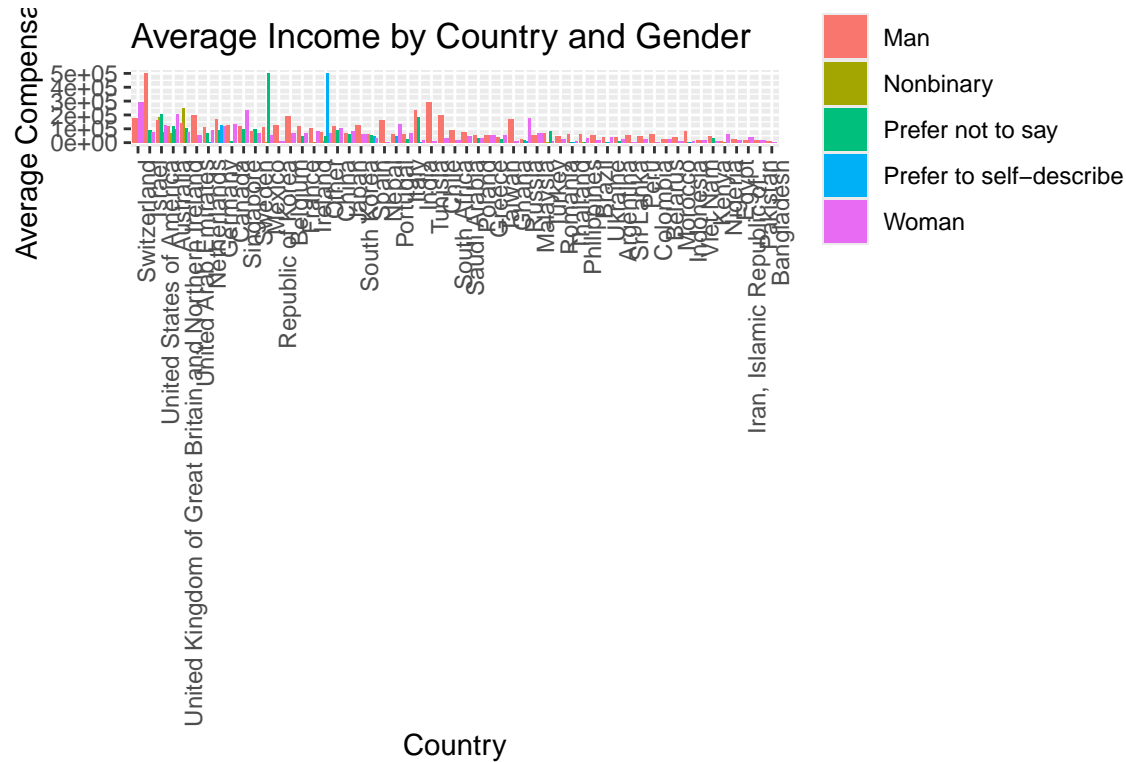
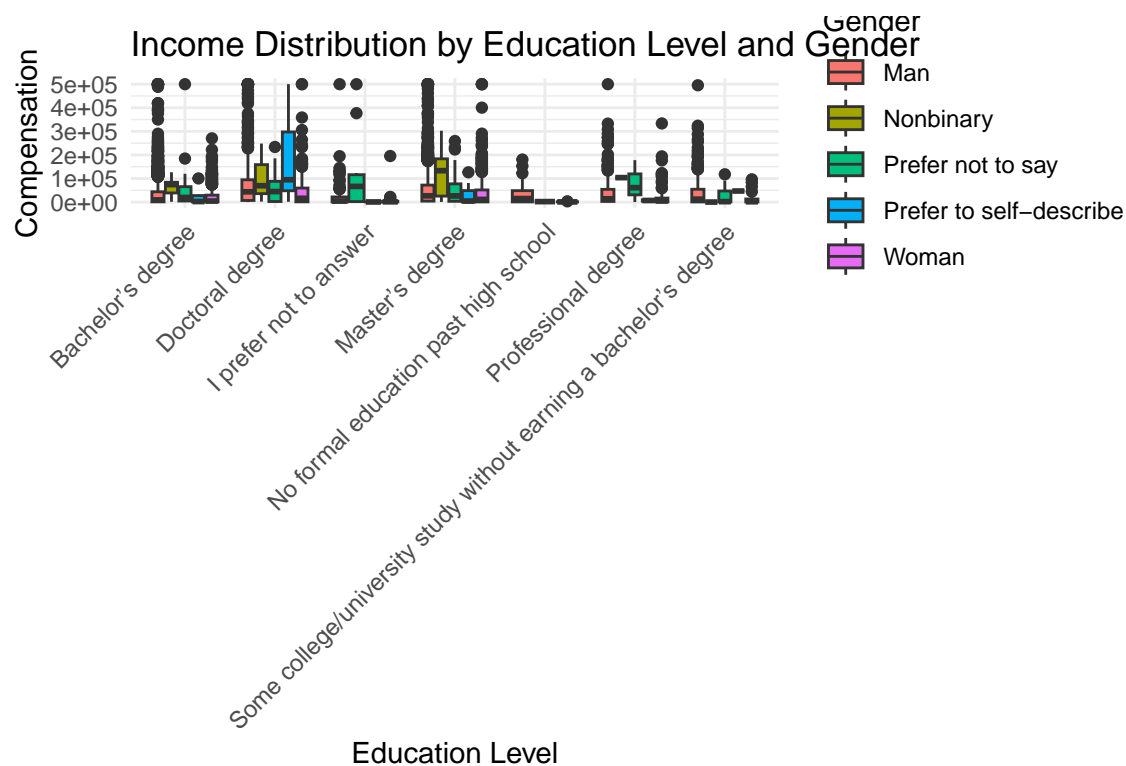
## Question 1

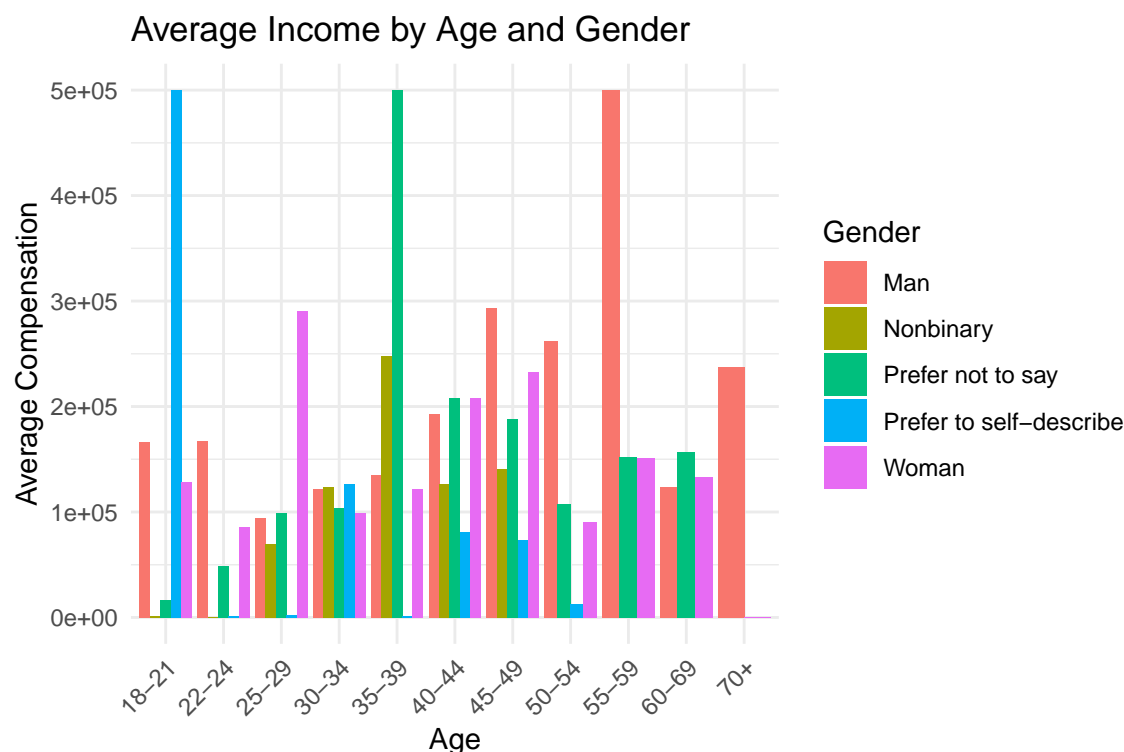
How do demographic factors (e.g., age, gender, country) correlate with income levels among data science professionals, and what role does the highest level of education play in these trends?

## Methods

The dataset provides information on respondents' age, gender, country, and highest level of education, alongside their reported income levels. Compensation data was cleaned and transformed by converting ranges into random values within specified bounds to enhance variability. Rows with missing or null compensation values were excluded to maintain data integrity. To explore trends, the data was aggregated and visualized using box plots and bar charts to examine correlations between income and demographic categories such as age, gender, country, and education.

## Analysis





The findings reveal distinct income patterns across demographic groups. Age strongly correlates with income, as older age groups (e.g., 55-59, 70+) generally report higher average compensation, likely reflecting accumulated experience. Gender disparities in income are evident, with men earning higher compensation on average across most age and education levels. The role of education is significant, as advanced degrees like doctoral and professional degrees correlate with higher median incomes. However, the variability within each educational category suggests other influencing factors. Country-level differences show significant disparities, with respondents from developed regions (e.g., Switzerland, the United States) reporting higher incomes compared to those from developing regions.

## Conclusion

The analysis demonstrates the substantial impact of age, gender, country, and education on income levels among data science professionals. Advanced education and experience play a critical role in achieving higher compensation, though gender disparities persist. Regional economic conditions further amplify these differences, emphasizing the importance of context in interpreting income trends. These insights are valuable for understanding global economic inequities in the data science field and can inform policies aimed at reducing disparities and promoting equitable professional development.

## Question 2

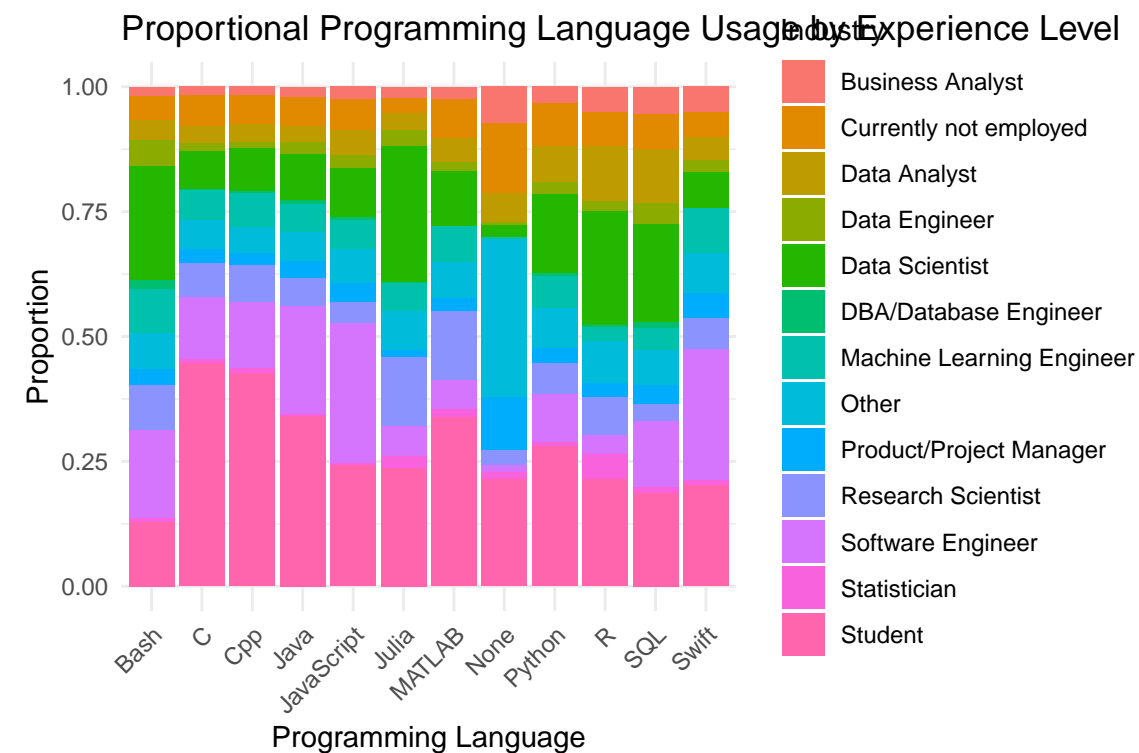
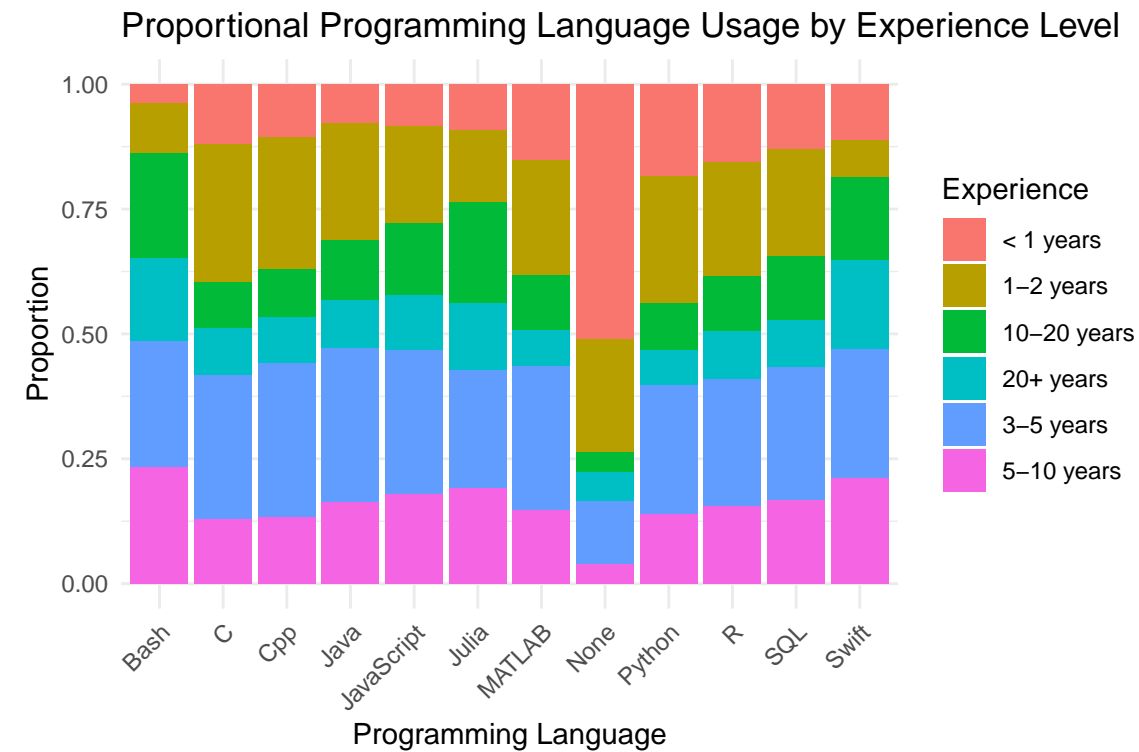
Are there distinct patterns in programming language usage between beginner, intermediate, and expert practitioners? How does this adoption vary across industries?

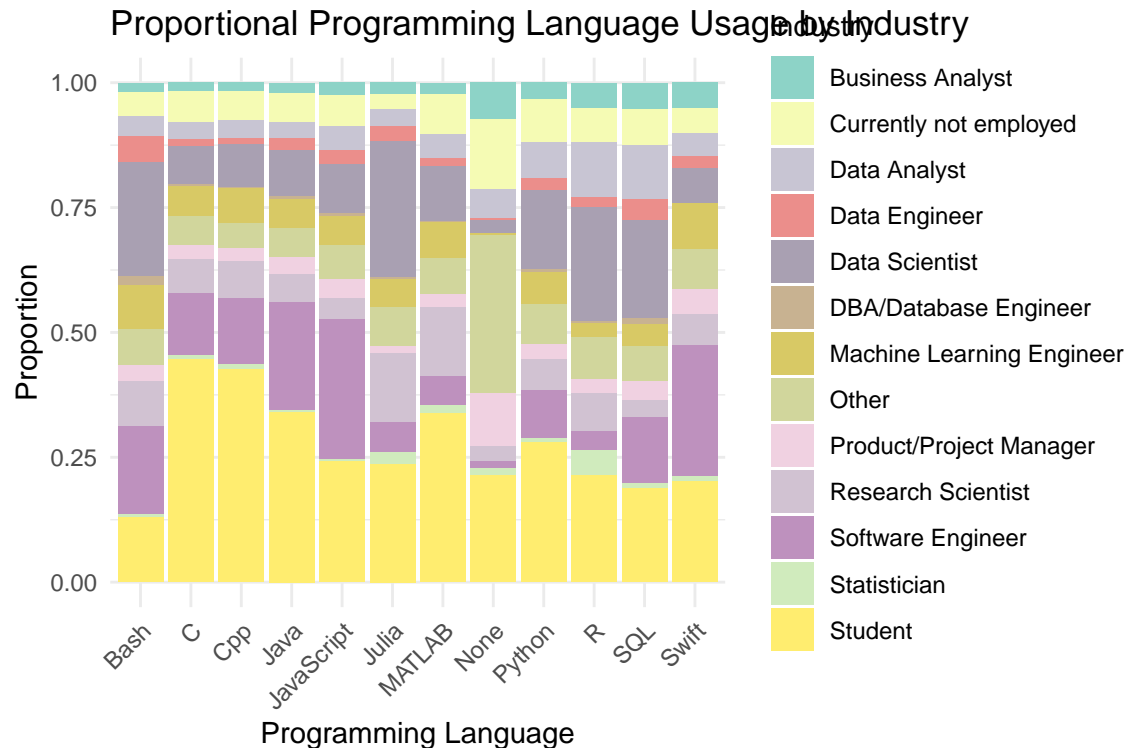
## Methods

Respondents' programming language usage and industry roles were derived from the survey's structured questions. Experience levels were classified into beginner, intermediate, and expert categories using the years of coding experience provided by the respondents. Proportional usage of programming languages across experience levels and industries was computed to identify trends. Visualizations were generated using stacked

bar charts to compare the proportional adoption of programming languages across both experience levels and industry categories.

## Analysis





The analysis revealed distinct trends in programming language preferences among practitioners of different experience levels. Beginners showed a strong reliance on foundational languages like Python and SQL, while experts exhibited broader adoption of languages like R, Java, and Bash for advanced tasks. Across industries, Python emerged as a dominant language for data-related fields like data science and machine learning, whereas SQL had widespread adoption in roles such as data engineering and database management. More niche languages, such as Julia and MATLAB, saw usage concentrated in specialized roles like research scientists.

## Conclusion

The study revealed distinct programming language preferences that evolve with practitioners' experience levels and vary by industry roles. Python's widespread adoption across all levels underscores its utility as a foundational language in the data science ecosystem. Meanwhile, experts and industry specialists integrate niche languages to address specific challenges. These findings provide actionable insights for curriculum design in education, workforce training, and recruitment strategies, emphasizing the importance of Python for beginners and advanced language fluency for experienced professionals in specialized roles.

## Question 3

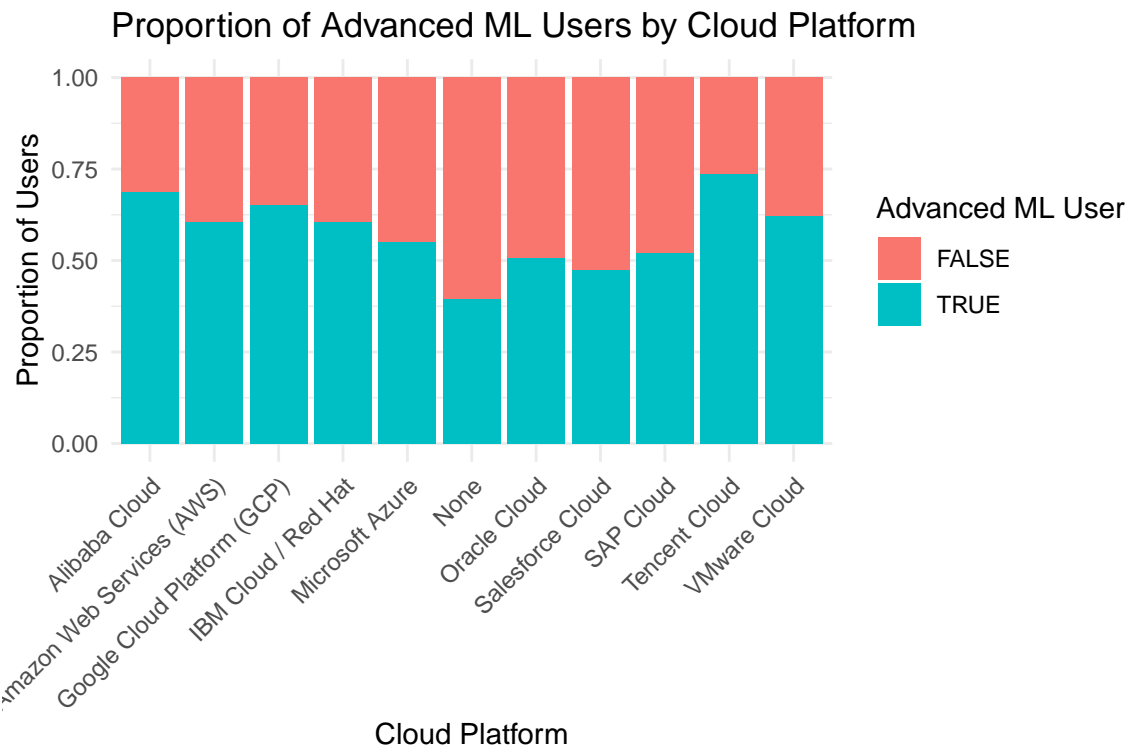
How does the adoption of advanced machine learning tools vary across different cloud platforms (e.g., AWS, GCP, Azure)? Which platforms show the highest proportion of advanced ML users, and what insights can this provide about user preferences or specialization trends in cloud services?

## Methods

We analyzed the relationship between cloud platform usage and advanced machine learning tool adoption using the Kaggle survey dataset. Advanced ML users were identified based on their reported use of tools such as Dense Neural Networks and Transformer Networks. Cloud platform codes were mapped to human-readable names, and the data was reshaped to analyze individual platform usage. We calculated the total users, advanced ML users, and their proportions for each platform, visualizing the results in a proportional bar

chart to highlight the platforms most associated with advanced ML adoption. This approach provided clear insights into user preferences and platform trends.

Analysis



The analysis reveals significant variation in the proportion of advanced machine learning (ML) users across different cloud platforms. The following insights can be drawn:

- **Tencent Cloud** has the highest proportion of advanced ML users, with approximately 73.7% of its users classified as advanced. This suggests that Tencent Cloud might be particularly well-suited for advanced ML tasks or attract users with more expertise.
- **Alibaba Cloud** follows closely, with 68.7% advanced ML users. Both Chinese platforms demonstrate a strong association with advanced ML activities, possibly due to their tailored offerings in AI and big data analytics.
- **Google Cloud Platform (GCP)** and **Amazon Web Services (AWS)** maintain strong proportions of advanced ML users at 65.2% and 60.4%, respectively. These platforms are renowned for their comprehensive ML services, such as TensorFlow on GCP and SageMaker on AWS, which cater to users with diverse levels of expertise.
- **Microsoft Azure** has a slightly lower proportion at 55.1%, reflecting its broad usage across both enterprise and technical user bases.
- **VMware Cloud** (62.1%) and **IBM Cloud / Red Hat** (60.5%) show notable engagement among advanced ML users. VMware’s focus on virtualization and IBM’s AI-driven services like Watson likely appeal to experienced users.
- **SAP Cloud**, **Oracle Cloud**, and **Salesforce Cloud** have lower proportions of advanced ML users, with values ranging from 47.4% to 52.1%. These platforms may cater more to enterprise and business use cases rather than technical ML applications.
- Interestingly, **None** (users not utilizing any cloud platform) constitutes a large group, but only 39.4% of these are advanced ML users. This indicates that cloud adoption is strongly linked to advanced ML capabilities.

## Conclusion

The results suggest a correlation between platform-specific features and their appeal to advanced ML users. Cloud providers like Tencent, Alibaba, GCP, and AWS dominate in attracting advanced practitioners, likely due to their robust AI and ML ecosystems. In contrast, platforms focused more on enterprise-level solutions (e.g., Oracle and Salesforce Cloud) show a weaker association with advanced ML users. These insights can help guide organizations in selecting cloud platforms and inform providers on how to refine their services to better support advanced ML workflows.

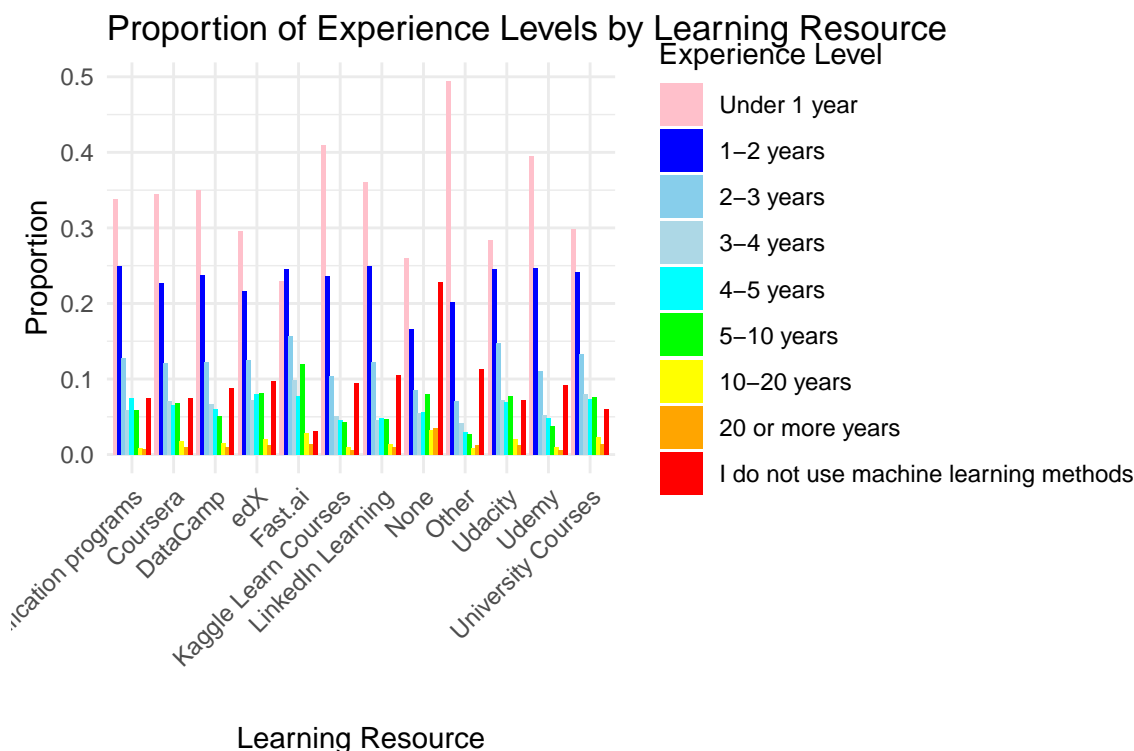
## Question 4

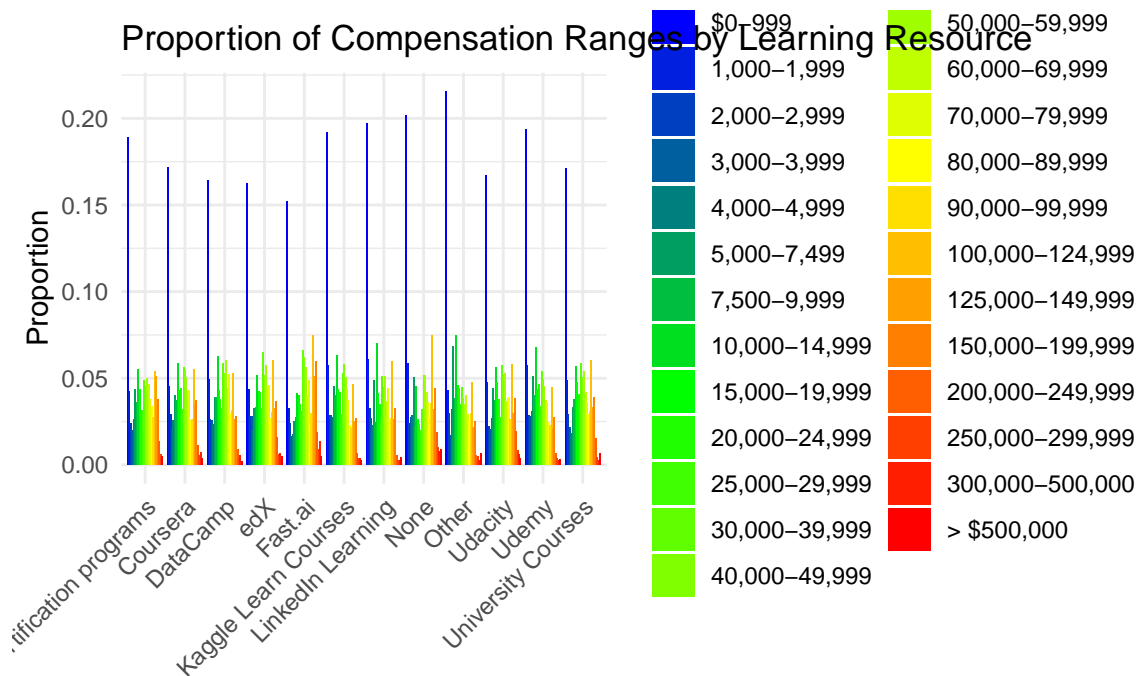
How do respondents' experience levels and compensation distributions vary across different educational resources (e.g., MOOCs, certifications, books)? Can we identify which resources are associated with higher compensation levels or a greater proportion of experienced professionals?

## Methods

Survey data was processed to analyze experience levels and compensation ranges across learning resources. Relevant columns were reshaped into a long format, with missing or empty values replaced appropriately. Experience levels (Q15) and compensation ranges (Q24) were standardized with ordered categories, while learning resources (Q37) were recoded into descriptive labels. Proportions of each category were calculated and visualized using bar plots to compare the distribution of experience and compensation across different learning resources.

## Analysis





## Learning Resource

The analysis highlights significant patterns in how learning resources are utilized across varying levels of experience and compensation. For experience levels, newer learners with less than a year of experience heavily rely on accessible platforms like Coursera and LinkedIn Learning, suggesting these platforms cater well to beginners. In contrast, advanced learners with over 10 years of experience are more inclined towards structured programs such as university courses or certification programs, reflecting their need for deeper, specialized knowledge.

Regarding compensation, resources like Cloud-certification programs and university courses showed a higher proportion of users in the upper compensation ranges, indicating their potential impact on career advancement. Conversely, platforms like Udemy and LinkedIn Learning had more users in lower compensation brackets, suggesting they are often leveraged for affordable, foundational learning. These insights demonstrate a clear connection between the type of learning resource, career stage, and compensation outcomes, emphasizing the diverse needs of learners across their career trajectories.

## Conclusion

In conclusion, the findings reveal distinct trends in how learning resources align with career stages and compensation levels in data science. Beginner-friendly platforms like Coursera and LinkedIn Learning are pivotal in attracting and supporting early-career professionals, while advanced learners gravitate towards more formalized education options such as certification programs and university courses, which are linked to higher compensation outcomes. These patterns underscore the importance of tailoring learning resources to meet the evolving needs of individuals at different experience levels. Furthermore, the observed correlation between resource utilization and compensation highlights the critical role of targeted skill development in career advancement, providing valuable insights for both learners and educators in shaping effective learning pathways.

## Advanced Analysis

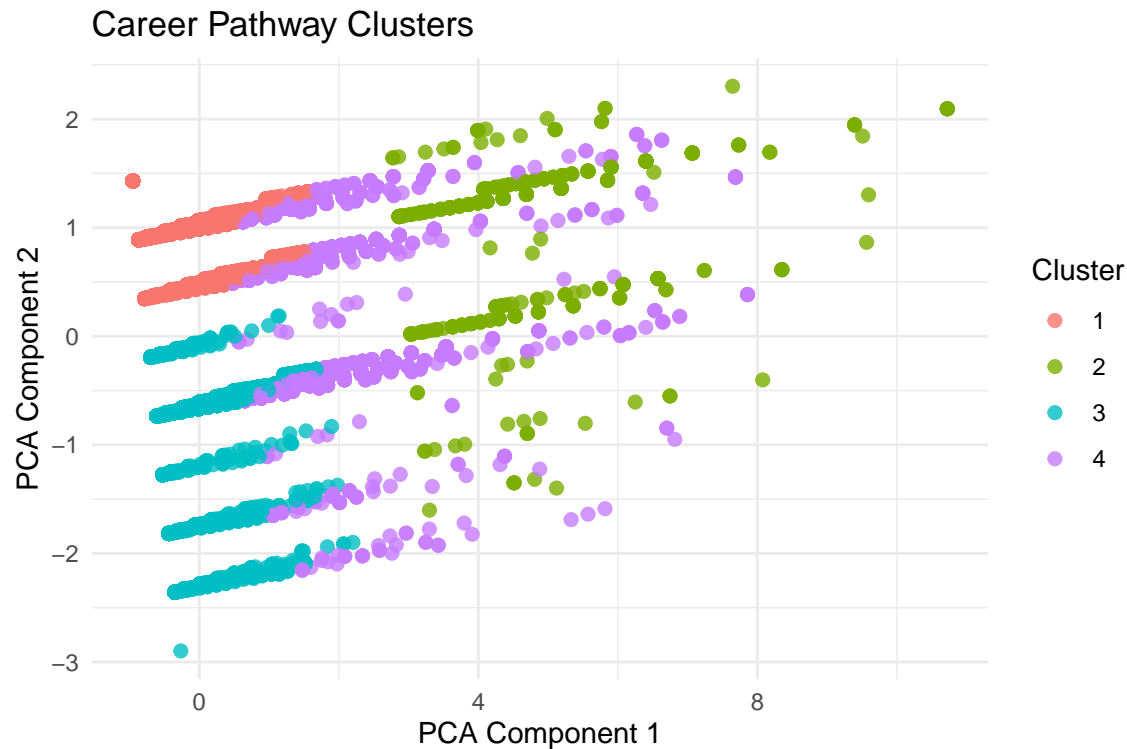
How do education levels, years of experience, and compensation interact to define distinct career stages in data science? Can we identify professional pathways, from entry-level roles to senior positions, and determine how education and experience contribute to progression between these stages?



## Methods

We analyzed career pathways in data science by clustering respondents based on experience (Q15), compensation (Q24), and education level (Q4). Experience levels and compensation ranges were mapped to numeric values, while education level was encoded as a numeric factor. Missing values were replaced with 0, and features were standardized to ensure comparability. Using K-Means clustering with four clusters, we identified distinct respondent groups. Principal Component Analysis (PCA) was used to visualize the clusters, and cluster characteristics—including average experience, compensation, and education level—were summarized to uncover patterns in career progression and demographics.

## Analysis



```
## # A tibble: 4 x 4
##   Cluster Avg_Experience Avg_Compensation Avg_Education
##   <fct>      <dbl>          <dbl>          <dbl>
## 1 1          0.962           7739.           2.13
## 2 2         17.1          104175.          3.89
## 3 3          1.07          10827.           5.47
## 4 4          3.67          138589.          4.13
```

The clustering analysis reveals distinct career pathways in data science based on experience, education, and compensation:

- **Cluster 1 (Early-Career Professionals):**

This cluster represents individuals with minimal experience (0.96 years on average) and low average compensation (\$7,739). Their average education level is 2.13, indicating most have some college experience without completing a bachelor's degree or are early in their educational journey. This group aligns with entry-level professionals or students who are just beginning their data science careers.

- **Cluster 2 (Seasoned Experts):**

Individuals in this cluster have significant experience, with an average of 17.1 years, and are compensated

accordingly (\$104,175 on average). The average education level is 3.89, suggesting that most respondents hold a master's degree or higher. This cluster represents established professionals who have progressed to senior or specialized roles in their careers.

- **Cluster 3 (Highly Educated Newcomers):**

Cluster 3 consists of individuals with low experience (1.07 years on average) and modest compensation (\$10,827). However, their standout feature is the high average education level of 5.47, indicating advanced degrees such as master's, doctoral, or professional degrees. These individuals are likely transitioning from academia to professional data science roles, leveraging their education to build careers.

- **Cluster 4 (Mid-Career Professionals):**

This cluster includes individuals with 3.67 years of experience and the highest average compensation (\$138,589). Their average education level is 4.13, which corresponds to respondents holding master's or doctoral degrees. These individuals are likely in mid-career positions, experiencing accelerated career growth and benefiting from the combination of advanced education and practical experience.

This analysis underscores the interplay between education, experience, and compensation in shaping career trajectories in data science. The insights highlight opportunities for professionals at various career stages to enhance their growth, whether through advanced education or gaining practical experience.

## Conclusion

This project provided an in-depth analysis of the Kaggle Survey 2020 dataset, exploring the dynamics of data science and machine learning professionals worldwide. By addressing questions on demographic and economic factors, programming language preferences, cloud platform adoption, educational resource effectiveness, and career pathways, we derived meaningful insights into the evolving landscape of the field.

1. **Demographic and Economic Correlation (Q1):**

The analysis revealed that factors such as age, gender, country, and education significantly influence income levels. Older respondents with advanced degrees earned higher compensation, reflecting the value of experience and education in career progression. However, notable gender disparities persisted across all levels, highlighting ongoing inequities. Regional differences further underscored the economic advantages of developed countries, where professionals reported substantially higher earnings.

2. **Programming Language Adoption (Q2):**

Trends in programming language usage varied by experience and industry. Beginners predominantly used Python and SQL, foundational tools for entering the data science domain. In contrast, experienced professionals adopted niche languages like R, Bash, and Julia for specialized tasks. Industry-specific preferences further emerged, with Python dominating data science roles and SQL widely used in engineering and database-focused positions.

3. **Impact of Cloud Computing on Advanced ML Usage (Q3):**

The analysis of cloud platform adoption revealed that platforms such as Tencent Cloud, Alibaba Cloud, and Google Cloud Platform demonstrated strong associations with advanced ML users, reflecting their robust AI ecosystems. Conversely, platforms like Salesforce Cloud and Oracle Cloud catered more to enterprise needs, with lower adoption by advanced ML practitioners. The results highlighted the role of cloud computing in enabling advanced analytics and AI capabilities.

4. **Effectiveness of Educational Resources (Q4):**

Educational resources showed varying impacts on experience levels and compensation outcomes. Beginner-friendly platforms like Coursera and LinkedIn Learning were popular among early-career professionals, while university courses and cloud certification programs were associated with higher compensation, reflecting their appeal to advanced learners. The findings emphasized the importance of tailoring resources to meet the diverse needs of learners at different career stages.

5. **Advanced Analysis – Career Pathways:**

Clustering analysis identified four distinct career pathways based on experience, compensation, and education. Cluster 1 represented entry-level professionals with minimal experience and education, while Cluster 2 highlighted seasoned experts with high compensation and advanced degrees. Cluster 3 captured highly educated newcomers transitioning from academia to industry, and Cluster 4 showcased mid-career professionals with rapid compensation growth. These clusters underscored the interplay between education, experience, and career progression in shaping professional trajectories.

This study illuminated the diverse dynamics of the data science field, providing actionable insights into the factors driving professional growth and specialization. Demographic and economic disparities emphasized the importance of addressing gender inequities and regional differences. Programming language adoption highlighted the utility of Python for beginners and the strategic integration of niche languages by experienced professionals. Cloud platform trends showcased the critical role of advanced ML capabilities in shaping platform preferences. Educational resources emerged as key enablers of career growth, with specific platforms and programs linked to higher compensation and advanced expertise.

The clustering analysis reinforced the significance of education and experience in navigating career pathways, offering a roadmap for aspiring professionals and organizations to strategize skill development and role transitions. Together, these findings contribute to a nuanced understanding of the global data science ecosystem, paving the way for more informed decision-making in education, recruitment, and professional development.