Final Project Reflection Paper

The final project for this course was a comprehensive and challenging experience that allowed our team to explore the process of uncovering societal trends through data analysis. By working with datasets focused on global university rankings and country development metrics, we attempted to investigate the relationships between higher education institutions and societal well-being. This reflection provides an overview of the challenges faced, lessons learned, and skills gained during the project, spanning from data selection to the final analysis and presentation.

## I.    Challenges

The first stage—data selection and exploration—was pivotal in setting the foundation for the project. We chose the 2023 Global University Rankings Dataset and the 2023 Global Country Development & Prosperity Index because of their complementary nature. While the university dataset provided insights into globally ranked universities and the countries they are located in and international student statistics, the development dataset offered a broader view of societal metrics such as economic quality, education, and living conditions. However, aligning these datasets presented significant challenges, particularly regarding compatibility and variable selection. Differences within the structure between the datasets required thoughtful adjustments to ensure meaningful comparisons. Despite these obstacles, this phase highlighted the importance of thoroughly understanding datasets and establishing clear research objectives which ultimately strengthened the rationale for our analysis.

The ETL process was conceptually straightforward but introduced several technical difficulties. Extracting data in CSV format was relatively simple, but cleaning and transforming the data presented significant challenges. Addressing missing values and merging datasets based on a common country column required precise handling to prevent data loss or inaccuracies. Implementing the ETL pipeline in Python further added complexity, particularly ensuring reproducibility and scalability. Designing the pipeline to handle updates to the data sources required additional testing and validation.

Setting up cloud storage on Google Cloud for the transformed data posed technical and logistical challenges. Ensuring secure access controls while maintaining ease of use for all team

members required balancing competing priorities. Additionally, integrating cloud storage with the ETL pipeline necessitated thorough testing to ensure seamless data transfer and accessibility during the analysis phase.

## II.     Lessons

Throughout the project, our team learned valuable lessons about both technical and collaborative aspects of data science. First, the importance of careful dataset selection became evident. Choosing datasets that are not only relevant but also structurally compatible significantly reduces the complexity of subsequent steps. Second, documentation and visualization tools, such as flowcharts, proved invaluable during the ETL design phase. They provided clarity, enhanced communication among team members, and helped identify potential restriction early on. Third, the importance of flexibility in data analysis was reinforced as we worked to interpret outliers and unexpected results. This required us to think critically and consider external factors that may influence the data.

On a collaborative level, regular team check-ins and task delegation were crucial to managing the project's timeline and complexity. However, we recognized room for improvement in coordinating tasks more effectively, particularly during the ETL implementation phase, where overlapping efforts occasionally led to redundancies.

## III.    Skills

The project significantly enhanced our technical and analytical skill sets. Designing and implementing an ETL pipeline in Python deepened our understanding of data preprocessing, from handling missing values to merging datasets. Using MongoDB for data storage improved our proficiency in managing structured and semi-structured data, while setting up Google Cloud storage introduced us to secure data sharing and credential management. Additionally, conducting EDA and creating insightful visualizations honed our ability to extract meaningful trends and communicate them effectively.

Beyond technical skills, the project strengthened our collaborative abilities. Navigating the complexities of team coordination, particularly during phases like data transformation and analysis, taught us the importance of clear communication and task ownership. These skills are invaluable for future projects, where interdisciplinary collaboration is often a key component.

**IV. Areas for Improvement**

While the project was largely successful, there are several areas for improvement. Automating more of the data cleaning process could save time and reduce manual errors in future projects. Additionally, conducting more extensive exploratory data analysis at the outset would allow for a deeper understanding of potential limitations and relationships within the data. Improving team coordination through clearer role assignments and more frequent progress updates could also enhance efficiency and reduce overlapping efforts.

Another area for development is enhancing our ability to address anomalies in the data. While we identified and interpreted several outliers, incorporating external datasets or qualitative research could provide additional context to strengthen our analysis. Lastly, further training in cloud computing and database management could improve our ability to integrate these tools seamlessly into future projects.

The Data Science Final Project was a challenging yet rewarding experience that provided valuable insights into the technical and collaborative aspects of data science. From struggling with dataset selection to implementing a robust ETL pipeline and analyzing societal trends, the project underscored the complexity and potential of working with real-world data. The lessons learned and skills gained will undoubtedly prove beneficial in future endeavors, assisting us with the tools and knowledge to tackle even more ambitious projects. While there is room for growth, this project has laid a strong foundation for our continued development.