

# Reflective Report

## Process of Solving Problems and Learning to Use Notebooks

### Process of Solving Problems

At the start of this unit, my approach to solving problems was relatively basic, focusing on straightforward techniques and limited exploration. As I progressed, I learned to adopt a more systematic and comprehensive methodology. This involved:

1. **Understanding the Problem:** Initially, I spent more time grasping the core of the problem. This meant reading the problem statement multiple times, understanding the data, and defining the objectives clearly.
2. **Data Exploration and Preprocessing:** I realized the importance of thorough data exploration and cleaning. Techniques such as handling missing values, outlier detection, and feature engineering became integral parts of my workflow. I used visualization tools like Matplotlib and Seaborn extensively to understand data distributions and relationships.
3. **Model Selection and Evaluation:** I moved from using basic models to experimenting with a variety of algorithms, understanding their strengths and weaknesses. I also learned the significance of model evaluation metrics beyond accuracy, such as precision, recall, F1-score, and ROC-AUC, to ensure a more robust assessment of model performance.

### Learning to Use Notebooks

Using Jupyter Notebooks as a primary tool for this unit was both a learning curve and an enlightening experience:

1. **Structured Workflow:** Notebooks allowed me to structure my code in a more readable and organized manner. The ability to mix code with markdown explanations helped in documenting the process and making the analysis more understandable.
2. **Interactive Development:** The interactive nature of notebooks facilitated quick experimentation. I could run individual cells, inspect outputs immediately, and make incremental changes without having to re-run the entire codebase. This was particularly useful for debugging and refining my approach.
3. **Visualization and Reporting:** Integrating visualizations directly into the notebook made it easier to interpret results and adjust my approach accordingly. This also enhanced the presentation aspect, as I could include plots and tables alongside explanations, making my work more communicative.

# Progress from the Start of the Unit

At the beginning of the unit, my understanding of data science and machine learning was fairly rudimentary. I had basic knowledge of Python and some experience with data analysis, but my approach was not systematic or thorough.

## Key Areas of Progress

1. **Technical Proficiency:** My coding skills in Python, especially with libraries like Pandas, NumPy, Scikit-learn, and visualization tools, have improved significantly. I now write more efficient, clean, and modular code.
2. **Analytical Thinking:** My ability to think analytically about data and problems has sharpened. I have learned to ask the right questions, design experiments, and interpret results critically.
3. **Machine Learning Expertise:** I have gained a deeper understanding of various machine learning algorithms, their applications, and their limitations. I am more confident in selecting and tuning models based on the problem at hand.
4. **Project Management:** Working on portfolios has improved my project management skills. I am better at planning, executing, and documenting my projects, ensuring a logical flow from problem statement to solution.

## Future Interests

Moving forward, I am particularly interested in:

1. **Advanced Machine Learning:** Delving deeper into advanced machine learning techniques, including ensemble methods, neural networks, and deep learning.
2. **Specialized Domains:** Applying my skills to specialized domains such as healthcare, finance, or environmental science, where data-driven solutions can have significant real-world impact.
3. **Big Data and Cloud Computing:** Exploring big data technologies and cloud computing platforms to handle and analyze large-scale datasets more efficiently.

## Discussion Points Based on Portfolio 4

### Why I Chose the Dataset for Portfolio 4

The dataset for Portfolio 4 was selected because of its critical importance and the richness of its features. Specifically, I chose a dataset on stroke prediction, which included a variety of features such as age, gender, hypertension, heart disease, and lifestyle factors. This dataset was appealing for several reasons:

1. **Relevance:** Stroke prediction is a significant and practical problem in the healthcare industry. Accurate predictions can lead to early interventions and potentially save

lives.

2. **Complexity and Depth:** The dataset had multiple features that required comprehensive preprocessing and feature engineering. This complexity provided an excellent opportunity to practice and refine these crucial skills.
3. **Potential for Insights:** Analyzing stroke data can yield significant insights into the factors that contribute to stroke risk, making the project both technically challenging and highly impactful.

## Identifying the Problem

The primary problem targeted in Portfolio 4 was predicting the likelihood of a stroke based on patient data. Identifying this problem involved several steps:

1. **Understanding the Dataset:** Initially, I performed exploratory data analysis (EDA) to understand the dataset's structure, the types of features it contained, and the relationships between these features and the target variable (stroke occurrence).
2. **Defining Objectives:** The objective was clearly defined as predicting stroke occurrences accurately. This involved understanding the medical context—how accurate predictions could benefit healthcare providers and patients.
3. **Challenges and Considerations:** Identifying potential challenges such as class imbalance, missing values, and feature importance was crucial. These considerations guided the choice of preprocessing techniques and models.

## Reason for Choosing Specific Machine Learning Models

- **Linear Regression:** Chosen for its simplicity in modeling linear relationships between health factors and stroke risk, offering interpretable insights into the impact of each predictor variable.
- **Polynomial Regression:** Extends linear regression to capture nonlinear relationships, potentially improving prediction accuracy by accommodating complex interactions among health factors.
- **Logistic Regression:** Well-suited for binary classification tasks like stroke prediction, providing probabilistic outputs for risk assessment and aiding clinical decision-making.
- **KNN Classifier:** Offers flexibility in classifying individuals based on similarity to neighbors, making it useful for personalized risk assessment without assuming underlying data distributions.

## Insights and Conclusions

- **Age and Stroke Risk:** Clear association found between age and stroke risk, consistent with expectations of higher risk in older individuals due to age-related health issues.
- **Biological Factors:** Hypertension, heart disease, and average glucose level showed notable correlations with stroke occurrence, supporting existing medical literature linking these factors to cardiovascular health and stroke risk.

