# Capstone Project Assignment: ML Model Development and Optimization

## Overview

The capstone project is a comprehensive assessment designed to integrate and apply the knowledge and skills you have acquired throughout the course. This project involves selecting a real-world dataset, conducting thorough analyses, building machine learning models, and addressing ethical considerations related to data usage. By completing this project, you will develop practical skills in data analysis, machine learning, and critical thinking, preparing you for professional challenges in the field of data science.

## Project Objectives

- **Application of Skills**: Utilize machine learning techniques learned in class to analyze a dataset of your choice.
- **Problem Solving**: Develop a structured approach to identify, analyze, and solve a real-world problem using data-driven methods.
- **Critical Thinking**: Evaluate various machine learning models and determine the most appropriate for your dataset.
- **Ethical Considerations**: Understand and analyze the ethical implications of your work, including issues of bias, fairness, and data privacy.

## Dataset Selection

For your capstone project, selecting the right dataset is crucial for conducting a meaningful analysis. The dataset you choose should meet the following criteria:

1. **Complexity**:
   o The dataset should contain **at least 500 rows** (observations) to ensure that there is enough data for meaningful analysis and model training.
   o It should have **at least 10 features** (variables), which can be a combination of numerical (e.g., age, income, temperature) and categorical variables (e.g., gender, country, product type). This diversity will allow you to explore various machine learning techniques and feature engineering methods.
   o **Example**: A dataset containing customer information for an e-commerce site with features such as customer_id, age, gender, purchase_amount, product_category, and purchase_date meets these requirements.

2. **Relevance**:
   o Your dataset should address a **real-world problem** that interests you or is significant in your field of study. This relevance will enhance your engagement with the project and make your findings more applicable and impactful.
   o **Example**: Possible topics include:
      o **Healthcare**: Patient records, disease outcomes, or healthcare access data.
      o **Finance**: Stock market data, loan approval datasets, or credit card transactions.
      o **E-commerce**: Customer reviews, sales data, or user interaction data.
      o **Environmental Studies**: Climate data, pollution levels, or wildlife tracking data.

3. **Public Availability**:
   o Ensure that your dataset is publicly available and can be accessed without any restrictions. You should source your data from reputable repositories, which often provide datasets that are cleaned and preprocessed to some extent.
   o **Reputable Sources Include**:
      o **Kaggle**: A platform with a vast collection of datasets across various domains. Check for datasets that are actively used and have accompanying kernels (notebooks) for reference.
      o **UCI Machine Learning Repository**: A well-known repository that contains a wide range of datasets for machine learning research.
      o **Government Open Data Portals**: Many governments provide open access to datasets for transparency and public interest. Examples include data.gov (U.S.) or the European Union Open Data Portal.

4. **Potential for Exploration**:
   o Choose a dataset that offers **opportunities for exploration and analysis**. Datasets that contain missing values, interesting relationships between variables, or are suitable for feature engineering will provide you with more avenues to demonstrate your analytical skills.
   o Look for datasets that allow for the application of various machine learning techniques (e.g., classification, regression, clustering).
   o **Example Considerations**:
      o Does the dataset contain **missing values** that you can address through imputation or other techniques?
      o Are there **interesting correlations** or trends that can be visualized?
      o Can you create new features from existing ones, such as extracting the month from a date field or combining multiple categorical features into a single variable?

**Additional Tips**
- **Preview the Dataset**: Before finalizing your selection, load a sample of the data to ensure it has the expected structure and quality.
- **Check Documentation**: Read any accompanying documentation or data dictionaries to understand the context and meaning of the features.
- **Consider Your Interests**: Choose a dataset related to a topic that you are passionate about or curious about; this will make the project more enjoyable and meaningful.

# Personal Reflection and Metacognition

As you progress through the capstone project, it is essential to engage in personal reflection and document your metacognitive processes. Metacognition refers to "thinking about thinking" and involves being aware of your cognitive processes during problem-solving and learning tasks. Regularly reflecting on your choices, challenges, and strategies will enhance your learning experience, help you identify areas for improvement, and develop a deeper understanding of your work.

**Reflective Questions to Guide Your Metacognitive Process**:
- What challenges did I face during this milestone?
- Which strategies were most effective in overcoming these challenges?
- How did my understanding of the project topic evolve?
- What have I learned about my own learning process, and how can I apply this in the future?

# Milestones and Detailed Steps

## Milestone 1: Data Preparation Report

**Due Date**: 11.10 (Sun)

**Purpose**: To prepare your dataset for analysis by performing exploratory data analysis (EDA), data cleaning, and feature engineering.

**Steps**:
1. **Exploratory Data Analysis (EDA)**
   - Load the dataset and display the first few rows to understand its structure.
   - Analyze data distributions for numerical features (e.g., using histograms, box plots).
   - Visualize relationships between variables (e.g., scatter plots, correlation matrices).
   - Identify any outliers or anomalies in the data.
2. **Data Cleaning**
   - Identify and handle missing values (e.g., imputation or removal).

- o Standardize data formats (e.g., date formats, categorical encoding).
- o Remove or address outliers identified in the EDA.
- o Document all steps taken during the cleaning process, explaining your reasoning.
3. **Feature Engineering**
   - o Create new features that could enhance the predictive power of your models (e.g., combining features, extracting date components).
   - o Transform categorical variables into numerical ones using techniques like one-hot encoding.
   - o Normalize or standardize numerical features if necessary.
   - o Provide a clear explanation of the rationale behind each feature created or modified.
4. **Report Preparation**
   - o Compile all findings, methods, and visualizations into a comprehensive Jupyter Notebook.
   - o Ensure the notebook is well-structured, with clear markdown sections explaining each part of your analysis.

**Deliverables**:
- **Jupyter Notebook**:
  - o A Jupyter Notebook containing all code for data preparation, including EDA, cleaning steps, and feature engineering. The notebook should be well-documented with comments explaining each step, along with markdown cells summarizing key insights and findings.

---

## Milestone 2: Model Development Report

**Due Date**: 11.24 (Sun)

**Purpose**: To apply machine learning algorithms to your dataset, evaluate model performance, and document the modeling process.

**Steps**:
1. **Model Selection**
   - o Research and select appropriate machine learning models based on your data and problem type (e.g., regression, classification).
   - o Justify your model choices by referencing their strengths and weaknesses in relation to your dataset.
2. **Model Training**
   - o Split the dataset into training and testing sets (e.g., 80/20 split).
   - o Train your selected models on the training set, documenting the processes used (e.g., hyperparameters, training algorithms).
   - o Use techniques such as cross-validation to optimize model parameters.

3. **Model Evaluation**
   - Evaluate model performance using suitable metrics (e.g., accuracy, precision, recall, F1-score for classification; RMSE, $R^2$ for regression).
   - Compare the performance of different models and justify your findings.
   - Visualize results using appropriate plots (e.g., confusion matrix, ROC curve).
4. **Code Quality Review**
   - Ensure your code is clean, well-organized, and well-commented for clarity.
   - Follow best practices in coding standards.
5. **Report Preparation**
   - Document your model selection process, training methods, evaluation metrics, and results in a well-structured Jupyter Notebook.
   - Include visualizations and interpretations of your model's performance in markdown cells.

**Deliverables**:
- **Jupyter Notebook**:
  - A Jupyter Notebook containing all code for model selection, training, evaluation, and visualizations. The notebook should include comments and markdown cells explaining your thought process and results.

---

## Milestone 3: Final Project Report

**Due Date**: **12.9 (Mon)**

**Purpose**: To analyze the performance of your final model, outline a comprehensive deployment plan, and discuss expanded ethical considerations.

**Steps**:
1. **Final Model Performance Analysis**
   - Analyze the performance of the final model chosen from Milestone 2.
   - Include a comparative analysis with alternative models, discussing performance metrics and insights from model interpretation techniques (e.g., SHAP values, LIME).
   - Discuss any limitations of your model and potential areas for improvement.
2. **Deployment Plan**
   - Outline a detailed deployment architecture that includes the model, data pipelines, monitoring systems, and user interfaces.
   - Simulate a mock deployment and provide a video demonstration or detailed description of the process.
   - Discuss considerations such as scalability, maintenance, and compliance with legal regulations.
3. **Ethical Considerations**
   - Conduct an in-depth analysis of the ethical implications of your project, focusing on a broader range of issues such as societal impacts and stakeholder effects.

o   Provide recommendations for ensuring fairness, accountability, and compliance with regulations.
4. **Report Preparation**
   o   Compile all analyses, plans, and discussions into a comprehensive Jupyter Notebook.
   o   Ensure the notebook is well-structured, clear, and professionally presented.

**Deliverables**:
- **Jupyter Notebook**:
   o   A Jupyter Notebook containing all code for the final model performance analysis, deployment steps, and any relevant visualizations. Include thorough documentation and explanations in markdown cells.

---

## Grading Summary

- **Milestone 1**: 2.5% of total course grade
- **Milestone 2**: 5% of total course grade
- **Milestone 3**: 7.5% of total course grade

## Conclusion

This capstone project will provide you with hands-on experience in machine learning while fostering critical thinking, problem-solving, and ethical considerations in your work. The grading rubric and timeline are designed to support you in successfully completing your projects.

# Milestone 1: Data Preparation Report (20 points)

| Criteria | Weight (%) | Excellent (5) | Good (4) | Satisfactory (3) | Needs Improvement (2) | Unsatisfactory (1) |
|---|---|---|---|---|---|---|
| **Data Exploration** | 40% | Comprehensive EDA with insightful visualizations. | Good EDA with appropriate visualizations. | Basic EDA with minimal insights. | Incomplete EDA with unclear visualizations. | No EDA conducted. |
| **Data Cleaning** | 30% | Thoroughly addresses all missing values; clear justification. | Addresses most missing values; some justification. | Addresses some missing values; minimal justification. | Limited addressing of missing values; unclear methods. | No data cleaning performed. |
| **Feature Engineering** | 20% | Creative and effective feature engineering; clear rationale. | Some useful features engineered; mostly clear explanations. | Limited feature engineering; lacks clarity. | Minimal or ineffective feature engineering; unclear rationale. | No feature engineering attempted. |
| **Reflection and Metacognition** | 10% | Comprehensive reflection on learning processes; insightful documentation of strategies and challenges. | Good reflection with clear documentation; minor gaps in detail. | Basic reflection present; lacks depth or clarity in documentation. | Limited reflection; minimal documentation of strategies. | No reflection or documentation of metacognitive processes. |

# Milestone 2: Model Development Report (20 points)

| Criteria | Weight (%) | Excellent (5) | Good (4) | Satisfactory (3) | Needs Improvement (2) | Unsatisfactory (1) |
|---|---|---|---|---|---|---|
| **Model Selection** | 30% | Clearly justified selection of multiple models. | Justifies selection of models; some explanation. | Basic model selection; limited justification. | Limited model selection with unclear justification. | No models selected. |
| **Model Training and Evaluation** | 40% | Comprehensive training process; thorough evaluation with appropriate metrics. | Good training process; evaluation mostly appropriate. | Basic training and evaluation; lacks depth. | Limited training and evaluation; missing metrics. | No model training or evaluation performed. |
| **Code Quality** | 10% | Clean, well-commented code; easy to follow. | Good code quality with some comments. | Basic code quality; limited comments. | Poor code quality; difficult to follow. | No code submitted. |
| **Reflection and Metacognition** | 20% | Comprehensive reflection on model selection and training; insightful documentation of strategies and challenges. | Good reflection with clear documentation; minor gaps in detail. | Basic reflection present; lacks depth or clarity in documentation. | Limited reflection; minimal documentation of strategies. | No reflection or documentation of metacognitive processes. |

# Milestone 3: Final Project Report (20 points)

| Criteria | Weight (%) | Excellent (5) | Good (4) | Satisfactory (3) | Needs Improvement (2) | Unsatisfactory (1) |
|---|---|---|---|---|---|---|
| **Final Model Performance** | 30% | Detailed analysis of final model performance, including comprehensive metrics and comparative analysis with alternative models. | Good analysis; includes relevant metrics and some comparison to alternative models. | Basic performance analysis; some metrics included; limited comparisons. | Limited performance analysis; few metrics provided; no comparative analysis. | No performance analysis included. |
| **Deployment Plan** | 30% | Comprehensive deployment plan; addresses all necessary steps, including architecture, scalability, and compliance with regulations. | Good deployment plan; most steps addressed; some considerations for scalability. | Basic deployment plan; missing significant details regarding architecture or compliance. | Limited deployment plan; lacks clarity; minimal discussion of real-world constraints. | No deployment plan provided. |
| **Ethical Considerations** | 20% | Thorough analysis of ethical implications, including societal impacts and stakeholder analysis, with actionable recommendations for | Good discussion; mostly relevant points covered, including some ethical considerations. | Basic mention of ethics; lacks depth; minimal discussion of societal impacts or stakeholder effects. | Limited discussion of ethics; minimal relevant points; no recommendations provided. | No ethical considerations discussed. |

| | | fairness and accountability. | | | | |
|---|---|---|---|---|---|---|
| **Reflection and Metacognition** | 20% | Comprehensive reflection on learning processes; insightful documentation of strategies and challenges throughout the project. | Good reflection with clear documentation; minor gaps in detail. | Basic reflection present; lacks depth or clarity in documentation. | Limited reflection; minimal documentation of strategies. | No reflection or documentation of metacognitive processes. |