

STATS 3DA3

Homework Assignment 6

Pratheepa Jeganathan

04/04/2024

Submission Deadline

- All submissions must be made before 10:00 PM on Thursday, April 18, 2024.

Submission Guidelines

- Format: Submissions are to be made in PDF format via Avenue to Learn, either individually or as a group of up to three members.
 - GitHub Repository: Your submission must include a link to a public GitHub repository containing the assignment.
 - Team Submissions: For group submissions, Question 15 must detail each member's contributions. Note that while there are no points allocated to Question 15, failure to provide this information will result in the assignment not being graded.

Late Submissions

- 15% will be deducted from assignments each day after the due date (rounding up).
- Assignments won't be accepted after 48 hours after the due date.

Assignment Standards

Please ensure your assignment adheres to the following standards for submission:

- **Title Page Requirements:** Each submission must include a title page featuring your group members' names and student numbers. Assignments lacking a title page will not be considered for grading.
- **Individual Work:** While discussing homework problems with peers and group is permitted, the final written submission must be your group work.
- **Formatting Preferences:** The use of LaTeX for document preparation is highly recommended.
- **Font and Spacing:** Submissions must utilize an eleven-point font (Times New Roman or a similar font) with 1.5 line spacing. Ensure margins of at least 1 inch on all sides.
- **Submission Content:** Do not include the assignment questions within your PDF. Instead, clearly mark each response with the corresponding question number. Screenshots are not an acceptable form of submission under any circumstances.
- **Academic Writing:** Ensure that your writing and any references used are appropriate for an undergraduate level of study.
- **Originality Checks:** Be aware that the instructor may use various tools, including those available on the internet, to verify the originality of submitted assignments.
- Assignment policy on the use of generative AI:
 - Students are not permitted to use generative AI in this assignment. In alignment with [McMaster academic integrity policy](#), it “shall be an offence knowingly to ... submit academic work for assessment that was purchased or acquired from another source”. This includes work created by generative AI tools. Also state in the policy is the following, “Contract Cheating is the act of”outsourcing of student work to third parties” (Lancaster & Clarke, 2016, p. 639) with or without payment.” Using Generative AI tools is a form of contract cheating. Charges of academic dishonesty will be brought forward to the Office of Academic Integrity.

Chronic Kidney Disease Classification Challenge

Overview

Engage with the dataset from the [Early Stage of Indians Chronic Kidney Disease \(CKD\)](#) project, which comprises data on 250 early-stage CKD patients and 150 healthy controls.

For foundational knowledge on the subject, refer to “Predict, diagnose, and treat chronic kidney disease with machine learning: a systematic literature review” by [Sanmarchi et al., \(2023\)](#).

Objectives

Analyze the dataset using two classification algorithms, focusing on exploratory data analysis, feature selection, engineering, and especially on handling missing values and outliers. Summarize your findings with insightful conclusions.

Classifier Requirement: Ensure at least one of the classifiers is interpretable, to facilitate in-depth analysis and inference.

Guidelines

- **Teamwork:** Group submissions should compile the workflow (Python codes and interpretations) into a single PDF, including a GitHub repository link. The contributions listed should reflect the GitHub activity.
- **Content:** Address the following questions in your submission, offering detailed insights and conclusions from your analysis.

Assignment Questions

1. **Classification Problem Identification:** Define and describe a classification problem based on the dataset.
2. **Variable Transformation:** Implement any transformations chosen or justify the absence of such modifications.
3. **Dataset Overview:** Provide a detailed description of the dataset, covering variables, summaries, observation counts, data types, and distributions (at least three statements).

4. **Association Between Variables:** Analyze variable relationships and their implications for feature selection or extraction (at least three statements).
5. **Missing Value Analysis and Handling:** Implement your strategy for identifying and addressing missing values in the dataset, or provide reasons for not addressing them.
6. **Outlier Analysis:** Implement your approach for identifying and managing outliers, or provide reasons for not addressing them.
7. **Sub-group Analysis:** Explore potential sub-groups within the data, employing appropriate data science methods to find the sub-groups of patients and visualize the sub-groups. The sub-group analysis must not include the labels (for CKD patients and healthy controls).
8. **Data Splitting:** Segregate 30% of the data for testing, using a random seed of 1. Use the remaining 70% for training and model selection.
9. **Classifier Choices:** Identify the two classifiers you have chosen and justify your selections.
10. **Performance Metrics:** Outline the two metrics for comparing the performance of the classifiers.
11. **Feature Selection/Extraction:** Implement methods to enhance the performance of at least one classifier in (9). The answer for this question can be included in (12).
12. **Classifier Comparison:** Utilize the selected metrics to compare the classifiers based on the test set. Discuss your findings (at least two statements).
13. **Interpretable Classifier Insight:** After re-training the interpretable classifier with all available data, analyze and interpret the significance of predictor variables in the context of the data and the challenge (at least two statements).
14. **[Bonus] Sub-group Improvement Strategy:** If sub-groups were identified, propose and implement a method to improve one classifier performance further. Compare the performance of the new classifier with the results in (12).
15. **Team Contributions:** Document each team member's specific contributions related to the questions above.
16. **Link** to the public GitHub repository.

Notes

- This assignment encourages you to apply sophisticated machine learning methods to a vital healthcare challenge, promoting the development of critical analytical skills, teamwork, and

practical problem-solving abilities in the context of chronic kidney disease diagnosis and treatment.

- Students can choose one classifier not covered in the lectures.

Grading scheme

1. Answer [1]
2. Codes [2]
OR answer [2]
3. Codes [3] and answer [3]
4. Codes [2] and answer [3]
5. Codes [2]
OR answer [2]
6. Codes [2]
OR answer [2]
7. Codes [3] and Plot [1]
8. Codes [1]
9. Answers [2]
10. Describe the two metrics [2]
11. Codes [2]
these codes can be included in (12)
12. Codes (two classifiers training,
model selection for each classifier,
classifiers comparisons) [5] and answer [2]
13. Codes [1] and answers [2]
14. Codes and comparison will
give **bonus 2 points for the final grade.**

The maximum point for this assignment is 39. We will convert this to 100%.

All group members will receive the same grade if they contribute to the same.