

## Project Description

Prof. Dr. Isabel Valera, Jonas Klesen, Ayan Majumdar

Version 1.0

### 1 Introduction

In the following, we describe the practical project that we have prepared for the students of the Machine Learning (ML) core lecture of Summer 2024. The idea of this project is that the students bring into practice the theoretical concepts introduced during the lectures.

This year the project consists on predicting the healthcare usage of individuals, and will involve both regression and classification tasks. The evaluation of the healthcare project is divided into two main parts that contribute differently to your final grade for the course: the report and the challenge. Details on grading are provided in section 5.

#### 1.1 Participation in the project

The participation to the project is **fully voluntary but highly recommended**, as it provides an opportunity for students to gain hands-on experience in ML. Students are expected to work in groups of up to 3 members and must register their group via CMS (already possible). **Participating in the project means submitting both the report and the predictions for the challenge.** It is not possible to participate in the challenge but not submitting the report, or the other way around. Refer to Section 6 for the expected timeline.

### 2 Problem Description

In healthcare, effective allocation of limited resources is crucial to adequately address the needs of all patients optimally. To better inform healthcare allocation, predictive analytics have become crucial in recent times to predict the healthcare usage of individuals.

Your objective for this project is to develop machine learning models capable of forecasting patients' healthcare utilization and the total cost incurred by each individual. To enable learning such models, you are provided with a dataset compiled from a survey of patients in the United States. This dataset encompasses various features related to health data, demographics, and specific services utilized by individuals.

The specific target features in this dataset and the respective tasks pertaining to these target features are as follows.

1. **Total Medical Expenditure:** This is a continuous feature measuring the total healthcare expenditure in US dollars. Your goal is to build a regression model that, when deployed, can take the features of previously unseen patients as input and predict the total medical expenditure.
2. **Healthcare Utilization:** This is a binary categorical feature that can either indicate LOW usage or HIGH usage. Your goal is to build a binary classification model that, when deployed, can take the features of previously unseen patients as input and predict the utilization category of the healthcare system.

#### 2.1 Datasets

We share three datasets with you, one training set and two test sets. For each set, the data will contain 108 features and, in the case of the training set, two target features. The targets are TOT\_MED\_EXP indicating the continuous healthcare expenditure in US dollars and UTILIZATION indicating the binary utilization class (LOW or HIGH). The data features include:

- **Demographics:** RACE, SEX, AGE, PANEL (survey panel number), REGION (Census region), MARITAL\_STAT (marital status), POVRTY\_CAT (poverty category), POVRTY\_LEV (poverty level), EDU\_YRS (education years), EDU\_DEG (highest education degree), SPOUSE\_PRST (marital status with spouse present), STUDENT\_STAT (student status), UNION\_STAT (union status), NUM\_DEP\_OUT\_REP\_UNT (number of dependents outside survey)

reporting unit), EMPLOYMT (employment status), OCCUP (occupation), NON\_ENG\_LANG (non-English language spoken).

- **Personal:** PUB\_ASST (public assistance money), TAX\_FORM\_TYP (tax form type submitted), FOOD\_STMP\_MNTHS (number of months food stamps purchased), FOOD\_STMP\_VAL (value of food stamps), MIL\_ACTIV\_DUTY (military active duty), HON\_DISCHARGE (honorary discharge from army), INSUR\_COV (insurance coverage), TOT\_INCOME (total income), EMPLOYR\_INS (employer offers insurance), CHILD\_SUPP (child support), PROB\_WKIDS (problem with kids), FAM\_INCOME (family income), PROB\_BILL\_PAY (problem with bill payments), DELAY\_PRESC\_MED (delay getting prescription medication), DAYS\_CAREOTHR\_NOWORK (days not working due to care for others), PENSN\_PLAN (pension plan), NO\_WORK\_WHY (reason for not working).
- **Health-related:** WEIGHT, HEALTH\_STAT (perceived health status), MENTAL\_HLTH (perceived mental health), CHRON\_BRONCH (chronic bronchitis), JNT\_PAIN (joint pain), PREGNT (pregnant), WALK\_LIM (walking limitation), ACTIV\_LIM (activity limitation), SOCIAL\_LIM (social limitation), COGNTV\_LIM (cognitive limitation), BM\_IDX (BMI), MULT\_HIGHBP (multiple high blood pressure readings), HOUSEWRK\_LIM (housework limitation), SCHOOL\_LIM (school limitation), ADV\_NO\_FAT\_FOOD (advised to restrict high-fat food), ADV\_EXERCISE\_MORE (advised to exercise), ADV\_DNTL\_CKP (advised dental checkup), FREQ\_DNTL\_CKP (frequency of dental checkup), RSN\_NO\_DNTL\_CKP (reason for no dental checkup), RSN\_NO\_MED\_CKP (reason for no medical checkup), DOC\_CHK\_BP (doctor checked blood pressure), TAKE\_RISK (prone to taking risks), ADV\_BOOST\_SEAT (advised booster seat), WHEN\_ADV\_BOOST\_SEAT (when advised booster seat), FEEL\_DEPRS (feels depressed), ADV\_NO\_SMKG (advised no smoking), AGE\_DIAG\_ADHD (age when diagnosed ADHD), PROB\_WBHV (problem with home behavior), WEAR\_SEATBLT (wear seat belt), WHEN\_ADV\_LAP\_BLT (when advised lap belt), WHEN\_LST\_ASTHMA (when last asthma episode), ADV\_LAP\_BLT (advised lap belt), ADV\_EAT\_HLTHY (advised to eat healthy), DOC\_TIM\_ALN (doctor spent time alone), APPT\_REG\_MEDCARE (made routine appointment medical care), LOST\_ALL\_TEETH, ASPRN\_REG (regular aspirin usage), DIFF\_ERRND\_ALN (difficulty doing errands alone), DIAB\_KIDNY (diabetes-related kidney issue), DIAB\_INSLN (diabetes insulin use), DIAB\_MED (diabetes medicine use), DISPSN\_STAT (patient disposition status), TIME\_LAST\_PSA (last PSA test), WHEN\_ADV\_EXERCISE (when advised to exercise more), DEAF, BLIND, LAST\_FLU\_VAC (last flu vaccine), UNABL\_PRESC\_MED (unable to get proper medicine), HEAR\_AID (need hearing aid), LAST\_REG\_CKP (last regular checkup), DAYS\_ILL\_NOWORK (days miss work for illness), DAYS\_ILL\_NOSCHL (days miss school for illness), HIGH\_BP\_DIAG, COR\_HRT\_DIAG (coronary heart disease diagnosis), ANGINA\_DIAG, HRT\_ATT\_DIAG (heart attack diagnosis), OTH\_HRT\_DIAG (other heart-related diagnosis), STROKE\_DIAG (stroke diagnosis), EMPHYM\_DIAG (emphysema diagnosis), HIGHCHOL\_DIAG (high cholesterol diagnosis), CANCER\_DIAG, DIAB\_DIAG (diabetes), ARTHR\_DIAG (arthritis), ARTHR\_TYPE (arthritis type), ASTHM\_DIAG (asthma), ADHD\_DIAG, NUM\_PRESCR\_MEDS (number of prescription medicines), DIFFIC\_HEAR (difficulty hearing), DIFFIC\_SEE (difficulty seeing), SMOK (smoking), OVR\_FEEL\_14 (overall feeling rating 14 days), MENTAL\_HLTH\_SCR (mental health score), PHY\_HLTH\_SCR (physical health score), OVR\_FEEL\_30 (overall feeling rating 30 days).

The training data has information and labels for 15000 patients. The first test dataset, named **test\_public.csv**, has features of 4791 patients without labels. You can access it right now on CMS. The second test dataset, named **test\_private.csv** has features of 5000 patients without labels. This dataset will be accessible closer to the final submission deadline and will be used to compute the final leaderboard of the challenge.

Note that the datasets we provide you have been only partially cleaned. Thus, it is important to try different pre-processing approaches for feature selection, feature encoding, outliers, etc. You are also encouraged to explore different trustworthy aspects for your ML modeling and deployment, e.g., exploring fairness and explainability aspects (see point 5 in Section 3 below).

## 2.2 Jupyter Notebook Example

We have uploaded a Jupyter Notebook in CMS that contains a description of the data and an example on how to load and preprocess the data, visualize it, and use it to train a linear regression and classification model. The notebook also includes instructions on how to format and save your predictions for submission to CMS (both for the leaderboard and the final ranking). We highly recommend that you use this notebook as a starting point for your solution.

## 3 Report Instructions

Every student group/team participating in the project should submit a report (.pdf file) using the LaTeX template provided in CMS. The report must be **at most 6 pages long (references excluded)** and contain detailed information on the methodology applied to select the final model and make the necessary predictions to participate in the challenge. More specifically, the report should contain information on the following aspects:

1. **Data analysis & preprocessing:** The report should describe any considered approach used for data analysis and preprocessing to prepare the input data (features) to the ML model.
2. **ML modeling:** The report should include a short description of the different models (i.e. the classifiers and regression models) applied to the data, specifying the used python libraries (if any).
3. **Model selection:** The report should detail the methodology followed to compare the different ML models (and, if applicable, data preprocessing approaches), as well as to select the final model used to make the predictions for the challenge.
4. **Empirical results:** The report should provide a summary and description of the empirical results that have led the students to select the final classification and regression models for the challenge.
5. **Others:** The report may contain any additional analysis performed by the students that may be interesting from a practitioner point of view. Examples of such analysis may i) provide a thorough data analysis (for example, data visualization using unsupervised learning techniques); or account for the robustness, explainability or fairness considerations of the different models explored by the students.

### 3.1 Report grading

The report will be graded. There are four possible grades for the project report:

- **[0 (out of 10) points]** If a major methodological mistake (e.g., selecting the ML model on the data used to train it) is detected.
- **[5 points]** If a subset of the models introduced in the lectures and tutorials are correctly applied, evaluated and reported.
- **[7.5 points]** If a comprehensive application of the techniques covered in the lectures are correctly applied, evaluated and reported.
- **[10 points]** If the students go one step beyond the course material. They may, e.g., provide additional content in the report covering data analysis or robustness/fairness aspects (see point 5. above) and/or apply methodology that goes beyond what has been introduced in the lectures to train excellent regression and classification models.

## 4 Challenge

The ML project will be maintained in a challenge format similar to a Kaggle-like competition. This means you will see your performance and ranking in a leaderboard which will be updated at multiple time stamps during the semester. To begin with, we plan to update the leaderboard about once every week, and towards the end of the semester (i.e. after the main exam until the project deadline) we will update the leaderboard about once every 2 days. After each update of the leaderboard you will be able to see your model's performance on the public test set and your ranking in the whole competition. In the following part, we will deliver the details about the challenge.

Throughout the semester, you are encouraged to work with your training data, **train.csv**, to update your models, try different approaches, or perform better model selection and hyper-parameter tuning. Your participation in the challenge will be evaluated in two folds.

1. **Leaderboard:** For the first part of the challenge, we have a **leaderboard** in CMS which will be updated multiple times during the semester. Throughout the semester, you can submit your predictions on the test set, **test\_public.csv**, in CMS. In this leaderboard, you will be able to see the performance evaluation of your team's model and its ranking among the models of other teams (for both prediction tasks). The idea of the leaderboard is twofold: i) give you a realistic estimate of the team ranking for each of the tasks to incentivize healthy competition; and ii) get you familiar with the challenge evaluation process.
2. **Final Assessment:** For the second and final evaluation of the challenge, you are supposed to submit your predictions on the private test set, **test\_private.csv**, which will be shared with you towards the end of the semester. You need to submit your predictions again through CMS. Your participation in the challenge will be assessed by your model's performance in this second test set alone, independent of the performance on the data from **test\_public.csv**.

## 4.1 Performance evaluation

The teams will be ranked for the the classification and regression tasks based on the following metrics:

- For the regression task

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i^{\text{true}} - y_i^{\text{pred}})^2}. \quad (1)$$

- For the classification task

$$(\text{macro}) \text{ Average f1-score} = \frac{1}{2}(F_1(Y = 0) + F_1(Y = 1)). \quad (2)$$

Above,  $F_1(Y = k)$  denotes the f1-score computed for class  $k$ .

### Bonus points:

**The top 5% of the teams, ranked according to the above metrics evaluated on the test\_private.csv dataset for both tasks, will get a bonus (one extra point in the German grading system) on their final grade.**

## 4.2 Submission

To submit your predictions for both the leaderboard and final submission, put your .npy files containing the predictions for both regression and classification tasks in the **same zip file**. Please name files inside the .zip as stated at the end of the provided Jupyter notebook.

## 5 Effect on overall course grade

The project is **fully voluntary** but it can help your grade, and did so for many students of previous iterations. Thus, it is highly recommended to do the project. To pass the course, you will have to pass either the main or re-exam. If you do not submit a project report, then your exam grade (i.e. the better grade of main exam and re-exam) will be your final grade. If your team submits the report, your overall course grade will be the better of the following two grades: your exam grade alone, or a 75%exam+25%report mixture grade. To put in in math terms:

$$\text{Course grade} = \max(\text{Exam grade}; 0.75 * \text{Exam grade} + 0.25 * \text{project report grade})$$

All members of a team get an identical report grade. Note again that to participate in the project (and get a potential grade boost), your team must submit **both** the report and the predictions for the challenge.

### Bonus points:

As noted at the end of section 4.1, the best-performing teams in the challenge will get a bonus, which increases your grade by 1 point in the German grade system. An example, taken from the previous iteration: if you have a grade of 2.0 in the exam, and a good report (10 out of 10 points), and you also do well enough in the challenge to get the bonus, your final course grade would be a 1.3.

## 6 Timeline

In the following, we detail the key dates that should not be missed if interested in joining the ML project:

- **Team registration:** via CMS due by 20.06.24. Teams of up to 3 students should be registered in CMS by then. Later changes in a team will only be possible under request.
- **Final project submission** via CMS due by 19.08.24. the students will need to submit both their report and final predictions (on TWEETS\_TEST\_2) by this date. We will not accept any late submissions.