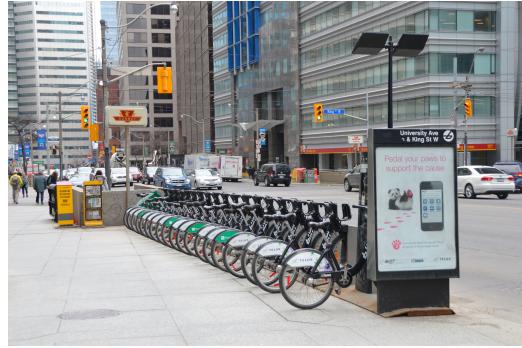


IML'22 Project: Bike-Sharing Demand Dataset

Background. Predicting the number of bikes available for sharing in urban cities is important for many reasons. Bike-sharing apps have become increasingly popular in recent years as a way to provide convenient and environmentally friendly transportation options in urban areas. However, managing a bike-sharing app can be challenging, as demand for bikes often varies throughout the day and week. Predicting the number of bikes available for sharing allows cities and bike-sharing companies to better plan for and manage their bike fleet. Cities, for example, can deploy additional bikes to a particular location if they predict high demand at a particular time. As a result, users will be able to find a bike more easily when they need one.



Dataset. The [Seoul Bike Sharing Demand Dataset](#) is a publicly available UCI Machine Learning Repository dataset. The dataset, which can be downloaded from [here](#), contains 8760 records of 14 attributes of weather information (Temperature, Humidity, Windspeed, Visibility, Dewpoint, Solar radiation, Snowfall, and Rainfall), relevant information about the season, holidays, and functional hours, along with the number of bikes rented per hour (target). Your task is to predict the number of bikes rented per hour given the other attributes. For the sake of simplicity, you can ignore the dates.

Requirements.

1. Formulate the machine learning problem by identifying the task (classification or regression), data, and challenges.
2. Perform Exploratory Data Analysis. This includes:
 - o Plot relationship between variables
 - o Identify correlated features or the ones that highly correlated with the label/outcome, if any
 - o Perform data imputation for missing variables
 - o Encode your outcome to a one-hot vector, if needed
 - o Remove redundant variables by dimensionality reduction techniques
3. Consider splitting the data into Training, Validation, and Testing sets. The suggested split is 70%, 10%, and 20% held-out testing set, respectively. (same split for all tasks)
4. Build and develop the following models (tasks):
 - o **Task#1:** Clustering in an unsupervised fashion
 - o **Task#2:** Logistic/Linear Regression according to your formulation in step 1
 - o **Task#3:** Neural Network for Classification/Regression according to your formulation in step 1

5. For each task, run the model *with* and *without* processing the data, e.g., without normalization or dimensionality reduction, and compare the model's performance after you normalize and/or reduce the dimensionality of the data.
6. For each task, show how you performed the *model selection*. For example, demonstrate the performance of variants of your model with different hyper-parameters, e.g., number of clusters and initialization when it comes to clustering methods.
7. For each task, perform a 5-fold Cross Validation.
8. For each task, run the corresponding evaluation metrics on each fold to demonstrate the performance of your model on the held-out testing set

Submission.

- **Prepare** a well-documented Jupyter notebook demonstrating the following sections:
 - Your name along with your Uni-ID and task assignment within the group
 - The background of the given task
 - Sample of your dataset
 - The rationale behind your ML formulation
 - The exploratory data analysis and data split
 - The developed models (Task #1-#3) including the model selection (hyper-parameters) and their results
 - A brief comparison between different models + concluding remarks
 - Lessons learned from this project
- **Consider** the following naming convention when you name your Jupyter notebook, BikeSharingDemandDataset_G##, and consider replacing ## with your group number, e.g, if your group number is 5, then the name of your file should be BikeSharingDemandDataset_G05.ipynb
- **Submit** the Jupyter notebook to the following folder
<https://uni-bonn.sciebo.de/s/t00R2IXI1RamUWB> which allows you to only upload files. You can upload the same file multiple times in case updates/modifications were required.
- **The submission deadline** is 18th Jan. 2023 at 11:59 PM

Grading.

- **Individual-level grading (65%) + Group-level grading (35%)**
 - The problem addressed (clearly stated and understood) – *Group-level*
 - ML problem formulation (rationale) – *Group-level*
 - Exploratory Data Analysis – *Group-level*
 - Task Completion and Findings – *Individual-level*
 - Interpretation and conclusion – *Individual-level*
 - Limitations – *Group-level*
 - Documentation, Submission, and Presentation – *Group-level*