

# MA124 Maths by Computer: Assignment 4

## Machine Learning Applied to Bike Sharing Demand Data (20 Marks)

---

In a recent research article published in the journal Computer Communications, authors Sathishkumar V E, Jangwoo Park, and Yongyun Cho sought to predict the "bike count required at each hour for the stable supply of rental bikes"[1]. They employed a number of regression models, including linear regression. The dataset used in the original study is available [here](#).

**Assignment:** Apply machine learning to a modified version of the original dataset and report the results.

[1] Sathishkumar V E, Jangwoo Park, and Yongyun Cho. 'Using data mining techniques for bike sharing demand prediction in metropolitan city.' Computer Communications, Vol.153, pp.353-366, March, 2020. [web link](#).

---

The original research article and a modified dataset are posted on the Moodle page. You will need to refer to the article for some of the tasks below. You will need to download SeoulBikeData\_mod.csv and put it into the folder with your assignment notebook. You do not need to submit SeoulBikeData\_mod.csv with your assignment (see below).

SeoulBikeData\_mod.csv has been modified from the original dataset to remove the categorical variables, and to convert dates to months. Months have been coded by number, e.g. 1 = January, etc. Only the first six months are included in the modified dataset.

---

While the number of tasks is large, this is in part because the instructions are rather specific. Many of them follow directly from the notebooks this week and last week. (Do not make the mistake of starting this assignment before doing all the notebooks from Week 7.)

Computational tasks:

1. Import needed libraries. (You will need pandas, seaborn, as well as things from sklearn, and of course numpy and matplotlib.)
2. Using pandas, read SeoulBikeData\_mod.csv into a Dataframe.
3. describe the Dataframe.
4. Plot a histogram of Rented Bike Count . Do not plot this as a density, but as a count. See Fig. 3 of Ref. [1].

(Optional, produce a box plot similar to that in Fig. 3 of Ref. [1]. If you produce the box plot, you may want both the histogram and box plot to be in approximately the aspect ratio of Ref.[1])

5. Produce two violin plots: one showing `Rented Bike Count` for different values of the `Month` and the other showing `Rented Bike Count` for different values of the `Hour`.
  6. From the full Dataframe, create a new Dataframe `X` containing all the columns except `Rented Bike Count` and a Series `y` containing only the `Rented Bike Count` column. These are your design matrix and target respectively.
  7. Perform a test-train split to create `X_train`, `X_test`, `y_train` and `y_test`. You **must** use the same percentage of data for testing and training as was used in Ref. [1] and you **must** state what they are. You can find these in the article.
  8. Create and train a linear regression model.
  9. Use the trained model to obtain `y_pred`, the prediction on the test data `X_test`. Form the residual `resid = y_test - y_pred`.
  10. Compute and report the: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and Rsquared (R2). These results should be similar to those on the top, right of Table 4 of Ref. [1]. (Note, the modified dataset we are studying is different from that used in the article. Hence the results will not be identical. However, the procedure is very close to that used in the article.)
  11. Produce the following plots.
    - A scatter plot of `resid` as a function of `y_test` corresponding to Fig. 9 of Ref. [1]. (Recall what `y_test` represents and label the plot appropriately.) Unlike Fig. 9 of the paper, you should use a colormap to plot the different months in different colours.
    - Histograms of `y_test` and of `y_pred` (on the same plot).
    - A scatter plot of `resid` as a function of `X_test['Month']`. Use a colormap to indicate the absolute value of `resid`.
    - A scatter plot of `resid` as a function of `X_test['Hour']`. Use a colormap to indicate the absolute value of `resid`.

(For all of the scatter plots, feel free to also vary the point size to make attractive and informative plots. Choose a colormap that looks good to you.)
  1. (Optional) You will see in Ref. [1] that most of the results involve "Trees". There are a number of types of trees used in machine learning. Sklearn provides a [DecisionTreeRegressor](#).
    - Create and train a `DecisionTreeRegressor` with `max_depth=6`. [Here](#) is an example that you will want to use.
    - The train model to obtain `y_pred` in this case.
    - Compute and report the: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R squared (R2).
    - Plot histograms of `y_test` and `y_pred` (on the same plot).
-

## Further 5 marks

A further 5 marks will be awarded for each assignment based on overall quality and clarity of the submission; the level of understanding demonstrated; originality, creativity and engagement.

---

## Submission

**You should not submit the SeoulBikeData\_mod.csv file.** You will submit **one Jupyter notebook**. **This must be a .ipynb file, not a pdf file or any other file type.** You should make sure that your notebook **executes without error before submitting it**. From the Kernel menu, you should: Restart Kernel and Run All Cells as the last step before saving your notebook and submitting it.

Clearly this assignment lends itself to producing a nice document. Such a document might be useful to you in the future, for example in applying for summer placements.

As usual, if the notebook is run and all code cells are collapsed, the notebook should be readable as a well-formatted report, primarily consisting of

- A short introduction making reference to the original research article. (Approximately 100 words might be appropriate for this assignment. Restate from the Abstract or Introduction what the article is about. Obviously you won't understand all the details, but in a few words you should be able to summarise the motivation and goals.)
- Computational tasks. Descriptions of these can be brief (from one to a few sentences, enough for the reader to follow without looking at the code.)
- Properly labelled figures and descriptions of each (this is very important).
- Ending summary of the results making connection back to the original research article, for example by comparing your results to those in the article.
- Somewhere in the report a full citation to the research article should appear. Use this assignment sheet as a model.

Use the example notebooks as a guide for Python style. One assumes the reader understands Python. Add comments to set off blocks of code or to note anything tricky. In most cases Python code explains itself.

---