

C-HACK 2022 – PROBLEM 1

General instructions

Note: Required submissions and documents are highlighted in red

- Take time to meet with your team and decide how you plan to solve the problem. Define your subtasks and decide who will take care of which subtask and when you will meet to update. Note that **teamwork will be evaluated** at the end of the hackathon.
- Work with your team to **generate a single Jupyter notebook** containing your solution to the given problem. You are free to include additional files needed to run your solution in the same folder, but we should be able to repeat your evaluation simply by running your notebook. If your notebook requires files/modules that are not turned in, we won't be able to run it!
- Your solution should be written in Python 3.
- Make sure that your notebook runs on Google Colab in your Google Drive team folder and is well-documented (i.e. provide explanatory and clarifying comments throughout your code).
- **(REQUIRED)** Write a **1 page summary** which describes your work and your results. You will be evaluated on your summary. The summary should include a title, an introduction, a paragraph which describes your results and a conclusion which summarizes your findings. See the template in the shared google drive folder **C-HACK 2022 EVENT/Templates/**
- Send us pictures of your team and/or your work as you hack your problem! We will be happy to post these on your team webpage and the C-HACK website news feed.
- **(REQUIRED)** Upload your **notebook, all files it depends on and a one page summary of your work** to the google drive folder we have created and shared with your team (this folder will not be seen by any other participants) **by 11:59 PM PDT on Wednesday January 12th 2021**
- **(REQUIRED)** Prepare a powerpoint or google drive presentation with your team (see the template in **C-HACK 2022 EVENT/Templates/**). **Submit the final version by uploading your presentation** to the shared google drive folder **by 12pm (noon) on Friday January 14th**. Recall that we will be posting your final presentations on the C-HACK website.
- Participate in the dry run on Thursday January 13th to practice presenting with your team.
- **(REQUIRED)** Present your final work with your team (**max 10min** to present) at the Final Presentations and Award ceremony on January 14th.
- Review and follow the [code of conduct](https://www.c-hack.org/), which is also posted on the c-hack webpage - <https://www.c-hack.org/>

Problem 1: Description

As announced this morning, DOW has shared *.csv format files containing data that was gathered from a rectangular waste treatment basin using a drone. In the data files, the positions at which the drone took samples within the basin (x, y position) at three depths (z position) are described by integers. For each sample there are a number of measurements taken: abundance data for different species of bacteria (labeled with “C”), concentrations of different metals (labeled with “FM”), and basin output variables (labeled with “O”). A more detailed description of the data is given below in the provided data section. You can find the files in the shared Google Drive folder: **C-HACK 2022 EVENT/Projects/DOW_data**. Note that you will only be able to view the data. To work with it you must **copy the files to your team’s Google Drive directory**.

IMPORTANT: **Do not share the DOW data with anyone outside the hackathon**, this data is proprietary to DOW and has kindly been provided **for the sole purpose of the hackathon**.

In this problem you will work on using **visualization** to analyze and represent the dataset.

Tasks for problem 1

- Analyze the data you are given and identify outliers and missing data.
- Produce a visualization of the basin’s performance
 - a. Think of a way to quantify **the basin’s performance**. Hint: there are 12 performance variables (“O”) which you could consider. We know that dissolved oxygen and MLSS ratio are correlated to the bacteria performing well, but there may be correlations to other variables.
 - b. Use your performance metric to **visualize** how well the basin is operating over space.

As you may have noted, the tasks are described in a broad context. We leave you the flexibility to decide...

- Which performance variables are important?
- How will you combine or incorporate different variables to determine performance?
- What is the clearest way to spatially visualize performance in the basin?

Provided data and code

Dow Dataset

Copy the contents of the “DOW_data” folder to your team’s folder!

The dataset comes in two csv files:

- CHACK2022_abundanceData.csv
- CHACK2022_OutputsMetals.csv

In both files, the position of each sample is given as columns ‘location x’, ‘location y’, and ‘location z’. For example $(x,y,z) = 1,1,1$ means basin section 1,1 at depth 1. Please see Figure 1 for a visualization of basin positions. In the first file, species abundance data is provided as columns, with labels starting with “C”. In the second file, additional measurements are provided as columns. The concentrations of various metals appear as labels starting with “FM” and “DM” (the metal identities have been hidden for proprietary reasons). Note that “FM” are the fixed metal concentrations, which represent the amount of metal that has been consumed by the bacteria, and “DM” are dissolved metals, so the concentration of metal dissolved in the solution. Other measurements starting with label “O” are described in the Table 1 below.

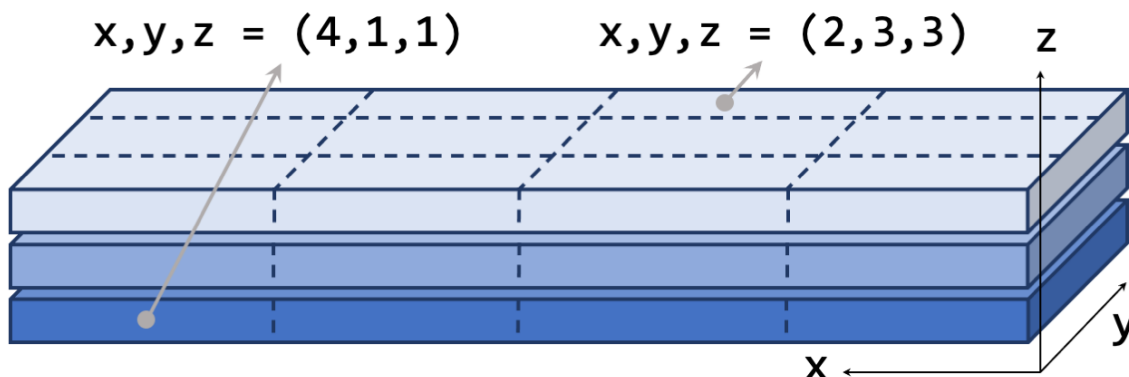


Figure 1: Illustration of how x , y , z positions are defined in the basin. A set of x , y , z defines one sector of the basin at a certain depth. Two examples are shown. Position 1,1,1 corresponds to the bottom right corner (the origin).

Table 1: Description of the output labels in the DOW dataset

Output	Full name	Units	Description
O_pH	pH	NA	Measure of how acidic or basic the basin is pH in range (7, 14] means basic pH in range [0, 7) means acidic pH = 7 means neutral

O_dO2	Dissolved Oxygen	mg/L	Amount of dissolved oxygen present in the water
O_conductivity	Conductivity	mS/cm	How well the water conducts an electric current
O_TDS	Total Dissolved Solids	mg/L	Combined amount of organic and inorganic substances dissolved in the water (Note: if the solid is present but not dissolved it is not accounted for in this quantity)
O_salinity	Salinity	Unknown	Concentration of salts in the water
O_T	Temperature	Celsius	Temperature in degrees Celsius
O_proteins	Protein content	mg/L	Amount of proteins present in the water
O_MLSS	Mixed Liquor Suspended Solids	Unknown	Total amount of solids suspended in the liquid – i.e. what's left if you filter the water
O_MLVSS	Mixed Liquor Volatile Suspended Solids	Unknown	Total amount of organic solids suspended in the water – i.e. what gets burned off if you filter the water then burn the solids in a furnace
O_MLSSratio	MLVSS/MLSS	Unitless	Ratio of MLVSS to MLSS

Code

To get you started on the right foot, we have provided you with a starter notebook, located in **C-HACK 2022 EVENT/Projects/Provided_Code**. This notebook will walk you through the process of mounting your google drive in the notebook and loading the data. We have also provided you with a python module called `CHACKutils.py`, located in the same directory as the starter notebook. Copy these files to your team's drive folder. In your notebook, you can import this

module by following the directions in the notebook to set up your google drive and executing `import CHACKutils`. A general guide on how to do these things can be found in the following videos ([Video 1](#), [Video 2](#)).

After you have imported the module with the above command, the function you will use to load the data is:

CHACKutils.dataloader()

Purpose: Load the DOW data.

- Inputs: You do not need to pass any arguments to execute this function. However, there is an optional argument `multiIndex`, which specifies if the abundance dataframe should be returned as a multi index dataframe. You will not need to use this option for this problem however feel free to look at it if it is of interest to you.
- Outputs: Two dataframes. The first dataframe is the DOW bacterial abundance data, and the second is the metals and output measurements data (Table 1).

Note: Make sure you have copied the 2 csv files from the **DOW_data** folder to your team's folder. The data files must be in the same folder as your notebook and the `CHACKutils.py` file for the data loading to work correctly.

Example use of data loading function, execute the following command:

```
species_data, outputs_and_metals = CHACKutils.dataloader()
```

`species_data` is the dataframe of abundance data, and `outputs_and_metals` is the extra measurements dataframe.

Terminology

- *Taxonomy* - Classification of organisms. All organisms fall within a set of taxonomic classes in a hierarchy of levels. For example, one taxonomic level is “kingdom” and all animals from us humans down to fleas are part of the “animalia” kingdom. Other kingdoms include “plantae” (plants). Other taxonomic levels include genus and species. For example, brown bears (species, latin name “arctos”) are part of the “Ursus” genus, still part of the animalia kingdom.
- *Performance* - This term describes how well the basin is doing what it is supposed to do. We know that it is meant to use bacterial biomass to convert toxic biproducts such as formaldehydes into something that can be discharged safely. A metric of performance is a number that describes on a scale from “bad” to “good” how well that task is being accomplished. Note that different metrics may be used to define performance.

- *Bacterial abundance* - This is a measure of portion (think of a pie chart) of the amount of bacteria in the sample, i.e. how much of it is a specific bacteria out of 1.0. For example, if one quarter of the bacterial mass in the sample is a specific species, that species has an abundance of 0.25. Note that this does provide a measure of how much total bacteria is in the sample.

Suggestions / References for problem 1

- Gather with your team and develop a plan! Not just a technical plan, but a working plan: how will you work together, who will work on what, how will you communicate? (text, zoom, slack ?). You must submit one working final notebook, however you may find it easier to work on separate notebooks for different portions of the project, keep each other updated, and create the final notebook as you go. Up to you!
- What does it mean for a column of data to be important? Are the values well distributed - do you have a homogeneous set of data? Are there relationships between variables?
- When is a plot representative? Would someone understand it without having any explanation? Does it show the behavior which you believe is there? Are the things which are significant labeled? Does the color scheme or visualization method emphasize the fact that, e.g. you have a subset of values which are larger or different from the rest?
- A 2D plot might be useful to show a trend, a 3D plot could show how points are distributed in space, a histogram could show how values are grouped (how frequently they occur)
 - <https://www.machinelearningplus.com/plots/top-50-matplotlib-visualizations-the-master-plots-python/>
 - <https://medium.com/swlh/python-data-visualization-with-matplotlib-for-absolute-beginner-part-iii-three-dimensional-8284df93dfab>
 - <https://towardsdatascience.com/data-visualization-how-to-choose-the-right-chart-part-1-d4c550085ea7>
- When looking at single columns of data you might want to see whether there are outliers - see e.g. <https://towardsdatascience.com/ways-to-detect-and-remove-the-outliers-404d16608dba>

GOOD LUCK!