

Problem Set 01: Data Exploration

Johannes C. B. Dietschreit, Sascha Mausenberger

Due: 14.10.2024

To solve the questions in this PSet, download the file `solubility_dataset.csv`. The problems are meant to be solved using Python functions and libraries. You can do so by either writing short scripts or using Jupyter notebooks. **Report your results in the Quiz on Moodle** to obtain a grade.

1 Introduction

In this PSet you will have the chance to familiarize yourself with libraries used when dealing with data sets. We will present you with some introductory code for some important packages, however, we expect you to look for functions you might need yourself (the documentations of all mayor Python packages we use in this course are very good, further Google and ChatGPT are of help).

1.1 Pandas

Pandas is a powerful and flexible open-source data analysis and manipulation library for Python, widely used for working with structured data. It provides data structures like DataFrames and Series, which allow for efficient handling, analysis, and manipulation of large datasets. It has extensive functionality for data cleaning, merging, and reshaping. Pandas is essential for data scientists and analysts.

```
1 import pandas as pd
2
3 # Load a CSV file into a DataFrame
4 df = pd.read_csv('data.csv')
5
6 # Display the first few rows of the DataFrame
7 print(df.head())
8 # Get basic statistical summary of the data
9 print(df.describe())
10
11 # Calculate mean, median, and standard deviation for a specific column
12 column_name = 'column_of_interest'
13 mean_value = df[column_name].mean()
14 median_value = df[column_name].median()
15 std_dev_value = df[column_name].std()
16
17 # Identify correlations between columns
18 print(df.corr())
19
20 # Count the occurrences of unique values in a specific column
21 print(df[column_name].value_counts())
```

1.2 Matplotlib and Seaborn

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python, offering a variety of plotting functions and customization options.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Load a CSV file into a DataFrame
5 df = pd.read_csv('data.csv')
6
7 column_name = 'column_of_interest'
8 # Histogram
9 plt.figure(figsize=(10, 6))
10 plt.hist(df[column_name], bins=30, edgecolor='k', alpha=0.7)
11 plt.title(f'Histogram of {column_name}')
12 plt.xlabel(column_name)
13 plt.ylabel('Frequency')
14 plt.show()
15 # Scatter plot
16 plt.scatter(df['x_column'], df['y_column'], alpha=0.7)
17 plt.title('Scatter Plot of x_column vs y_column')
18 plt.xlabel('x_column') # Replace with your x column name
19 plt.ylabel('y_column') # Replace with your y column name
20 plt.show()
21 # Line plot
22 plt.plot(df['x_column'], df['y_column'], marker='o')
23 ...
```

Seaborn is a statistical data visualization library based on Matplotlib, designed to create informative and attractive visualizations with minimal code.

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 df = pd.read_csv('data.csv')
6 # Plotting with Seaborn
7 column_name = 'column_of_interest'
8
9 # Histogram with Seaborn
10 sns.histplot(df[column_name], bins=30, kde=True)
11 plt.title(f'Histogram of {column_name}')
12 # labels are often auto-generated by Seaborn
13 plt.show()
14
15 # Scatter plot with regression line
16 sns.scatterplot(x='x_column', y='y_column', data=df)
17 sns.regplot(x='x_column', y='y_column', data=df, scatter=False, color='
    red')
18 plt.title('Scatter Plot with Regression Line')
19 plt.xlabel('x_column')
20 plt.ylabel('y_column')
21 plt.show()
22
23 # Pairplot for visualizing pairwise relationships in a dataset
24 sns.pairplot(df)
25 plt.show()
```

1.3 RDKit

RDKit is an open-source toolkit for cheminformatics that supports various functionalities for chemical informatics, such as reading and writing molecular data, calculating molecular properties, performing substructure searches, and generating 2D and 3D molecular representations. It is widely used in computational chemistry and bioinformatics for drug discovery and molecular modeling tasks.

```
1 # Import RDKit modules
2 from rdkit import Chem
3 from rdkit.Chem import Descriptors
4 from rdkit.Chem import Draw
5
6 # Create a molecule from a SMILES string
7 smiles = 'CCO' # Ethanol (hydrogens are omitted)
8 molecule = Chem.MolFromSmiles(smiles)
9
10 # Print formular, calculate weight
11 print("Molecular Formula:", Chem.rdMolDescriptors.CalcMolFormula(molecule
12 ))
13 mol_weight = Descriptors.MolWt(molecule)
14 print("Molecular Weight:", mol_weight)
15
16 # Draw the molecule
17 Draw.MolToImage(molecule)
18 # In a notebook one can simply do
19 display(molecule)
20
21 # Convert a molecule to a SMILES string
22 smiles_back = Chem.MolToSmiles(molecule)
23
24 # Calculate other properties, based on heuristics
25 num_h_donors = Descriptors.NumHDonors(molecule)
26 num_h_acceptors = Descriptors.NumHAcceptors(molecule)
27 logp = Descriptors.MolLogP(molecule)
28
29 # Perform a substructure search
30 substructure = Chem.MolFromSmarts('CO') # Look for an alcohol group
31 matches = molecule.GetSubstructMatches(substructure)
```

2 The Problems

2.1 Inspect the Column Names

Which is not a column in `solubility_dataset.csv`, multiple answers allowed? (1 P)

- (a) NumBonds
- (b) Solubility
- (c) InChIKey
- (d) SELFIES
- (e) SMILES

2.2 Size of a DataFrame

How many rows does the table (the CSV-file) have? (1 P)

2.3 Basic Stats

1. Compute the mean and standard deviation of the column `Solubility`. Round your answer to 2 digits after the decimal point. (each 0.5 P)
2. Compute the covariance of `Solubility` and `MolWt`. Round your answer to 2 digits after the decimal point. (1 P)
The general formula is

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

the normalized one

$$\text{cov}(X, Y) = \frac{\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]}{\sqrt{\mathbb{E}[(X - \mathbb{E}[X])^2]\mathbb{E}[(Y - \mathbb{E}[Y])^2]}}$$

3. Compute the covariance of the number of heavy atoms and the number of all atoms (including H) for the data set (hint: you will need to use RDKit for this). Round your result up to two digits after the decimal point. (1 P)

2.4 Getting the Index of a Value

1. The column `MolWt` gives the Molar weight of each compound and `RingCount` the total number of rings (aliphatic and aromatic). How many rings are in the haviest compound? (0.5 P)
2. The column `NumValenceElectrons` contains, as the name suggests, the number of valence electrons of each compound. Find the InChIKey of compound with the fewest valence electrons. (0.5 P)

2.5 Display Lewis Structure

Display the Lewis structure of the 56th molecule in the table (index 55, remember Python starts counting at 0). (1 P)

2.6 Plotting

1. Make a histogram of the data in the column MolLogP, you can use `matplotlib` or `seaborn`. (1 P)
2. Make a scatter plot of Solubility (x-axis) and MolLogP (y-axis). (1 P)

2.7 Searching Strings

How many compounds contain at least one zink ion/atom (hint: use the SMILES)? (1 P)