

The "Cancerdrug2.ipynb" notebook performs an analysis of a drug sensitivity dataset, focusing on data cleaning, exploratory data analysis, and feature engineering. Here's a breakdown of the key steps and findings:

1. Data Loading and Initial Inspection

- The analysis begins by importing necessary libraries such as pandas, numpy, matplotlib, and seaborn.
- The dataset is loaded from a CSV file named "GDSC2-dataset.csv" into a pandas DataFrame.
- The code then checks for missing values using `df.isnull().sum()`, revealing that the TCGA_DESC column and PUTATIVE_TARGET have a large number of null values.

2. Data Cleaning

- Rows with any missing values are dropped using `df.dropna()` to create a clean DataFrame called `df_clean`. This ensures that further analysis is performed on complete records.
- The data types of the columns in `df_clean` are examined using `df_clean.dtypes`. This step confirms the types of data in each column (e.g. integers, floats, or objects).

3. Descriptive Statistics

- The `df_clean.describe()` method provides descriptive statistics for the numerical columns of the DataFrame, such as count, mean, standard deviation, minimum, and maximum values.
- These statistics offer insights into the distribution and range of values for columns such as LN_IC50, AUC, RMSE, and Z_SCORE.

4. Exploratory Data Analysis (EDA)

- Histograms are generated for the numerical columns (LN_IC50, AUC, RMSE, Z_SCORE) to visualize the distribution of the data.
- These histograms are enhanced with Kernel Density Estimation (KDE) plots to show the probability density of each variable.
- The notebook performs additional EDA, calculating and displaying mean values for the LN_IC50 column grouped by the DRUG_NAME.
- The average RMSE and Z_SCORE are also calculated.

5. Feature Engineering

- **Categorical Variable Encoding:** The notebook uses two techniques to transform the categorical columns into a numerical format:
- **One-Hot Encoding:** The columns DATASET, CELL_LINE_NAME, SANGER_MODEL_ID, TCGA_DESC, DRUG_NAME, PUTATIVE_TARGET, PATHWAY_NAME, and COMPANY_ID are one-hot encoded using "OneHotEncoder". This creates a new binary column for each unique value in the categorical column.

- **Label Encoding:** Simultaneously, the original categorical columns in *df_clean* are 'label encoded' using LabelEncoder. This maps each unique value to an integer.
- The original dataframe *df_clean* is concatenated with the one hot encoded dataframe *one_hot_encoded_df* to create a new dataframe *df_clean_one_hot*.
- **Data Splitting:** The dataset is split into training and testing sets (80% training, 20% testing).
- The shape of the training and test datasets are printed to confirm the size of the split datasets.

6. Drug Effectiveness Analysis

- The mean LN_IC50 values are calculated for each drug and sorted in descending order to determine the most effective drugs based on their mean LN_IC50 values.
- Bar charts are generated to visualize the top 10 and top 20 most effective drugs, as ranked by mean LN_IC50.

Summary of Analysis The notebook provides a comprehensive overview of the dataset. It highlights the importance of preprocessing steps such as data cleaning and feature engineering. It uses visualization and summary statistics to offer insights into the data and begins to examine drug effectiveness based on the mean of LN_IC50 values. It also highlights common methods for preparing a data set for downstream model training tasks.