

The "CancerDrugAnalysis(1).ipynb" notebook provides a comprehensive exploratory data analysis (EDA) of a cancer drug sensitivity dataset. The analysis focuses on understanding the distributions of key variables, their correlations, and the relationships between drug responses and various cancer and tissue descriptors. Here's a detailed breakdown of the analysis:

1. Data Loading and Preprocessing

- The notebook begins by loading the dataset, "GDSC_DATASET.csv", into a pandas DataFrame and importing the necessary libraries: pandas, numpy, seaborn, and matplotlib.
- It then checks for missing values in the dataset.
- The notebook preprocesses the dataset by removing rows with any missing values using `df.dropna()`. This step ensures that subsequent analyses are conducted on complete data.

2. Distribution Analyses

- **LN_IC50 Distribution:**
 - The distribution of LN_IC50 values is analyzed using box plots, grouped by the Cancer Type (matching TCGA label). This allows for visualization of the spread and central tendency of drug sensitivity, as measured by LN_IC50, across different cancer types.
 - Additionally, the notebook examines the distribution of LN_IC50 using box plots, this time grouped by GDSC Tissue descriptor 1. This step reveals how drug sensitivity varies across different tissue types.
 - Histograms are also used to show the distribution of LN_IC50 values.
- **Z-score and AUC Distribution:**
 - Box plots are generated to show the distribution of Z_SCORE by Cancer Type (matching TCGA label). This is done to understand how the normalized drug response varies across cancer types.
 - Box plots are also used to show the distribution of AUC (Area Under the Curve) by Cancer Type (matching TCGA label).
 - Histograms are used to show the distributions of AUC and Z_score.

3. Correlation Analysis

- **AUC vs. Z-score:** The notebook uses a scatter plot to visualize the relationship between AUC and Z_SCORE. This step helps in understanding if there is any correlation between these two drug response metrics.
- **Correlation Heatmap:** A heatmap is generated to display the correlation matrix between LN_IC50, AUC, and Z_SCORE. This visualization provides a clear picture of how these variables are related to each other.

4. Tissue Descriptor Analysis

- The notebook explores the GDSC Tissue descriptor 1 and GDSC Tissue descriptor 2 columns by printing the rows where GDSC Tissue descriptor 1 has null values.
- The distributions of GDSC Tissue descriptor 1 and GDSC Tissue descriptor 2 are analyzed using count plots. These plots show the frequency of different tissue descriptors in the dataset.
- A scatter plot is generated to visualize the relationship between GDSC Tissue descriptor 1 and GDSC Tissue descriptor 2.

5. Drug Analysis

- **Drug Frequency:**
 - The notebook analyzes the frequency of drug usage in the dataset, plotting the top 10 most frequent drug names.
 - It also generates a list of the unique drugs present in the dataset.
- **Drug Usage by Tissue Descriptor:**
 - The notebook creates bar plots showing the top 10 most frequent drug names for each GDSC Tissue descriptor 1 and GDSC Tissue descriptor 2. This analysis reveals which drugs are most commonly associated with each tissue type.
 - The code generates a visualization of the occurrences of the tissue descriptors by drug name, but it uses a sample dataset.

6. Additional Visualizations

- **Data Subsetting:** The notebook shows how to subset the data using .loc to display a specific range of rows.
- It also prints the head and tail of the dataset.
- **Drug Counts:** The notebook plots the count of each drug used, and it also prints the summary statistics of the counts of the drugs. It also identifies and prints the least 30 drugs used.
- **Drug and Tissue Relationships:** The notebook visualizes the relationship between top drugs and tissue descriptors using bar plots. It shows the top 10 most used drugs and also the 8 least used drugs, all grouped by tissue descriptor.
- **Pie Chart:** The notebook generates a pie chart of the 8 least used drugs.
- **Histograms of Tissue Descriptors:** The notebook generates a histogram and a bar chart of the tissue descriptors.

Summary of Analysis

The "CancerDrugAnalysis(1).ipynb" notebook conducts a thorough exploratory analysis of a cancer drug sensitivity dataset. It investigates:

- **Distributions of key variables:** Including LN_IC50, AUC, and Z_SCORE across different cancer types and tissue descriptors.

- **Correlations:** Between AUC and Z_SCORE.
- **Relationships:** Between drug responses and various cancer and tissue descriptors.
- **Drug Usage:** It shows the frequency of drugs, and how they relate to tissue types.

This detailed analysis provides a solid foundation for further statistical testing and modeling. The visualizations and summary statistics help to reveal patterns and relationships within the data, which are essential for understanding drug sensitivity in cancer research.