

Introduction

Statistical analysis is a fundamental aspect of data science and machine learning. By understanding and applying statistical methods, you can derive meaningful insights from datasets, make predictions, and inform decision-making. In this session, we will explore the basics of statistical analysis using Python, leveraging libraries like pandas, numpy, and scipy. We will use various datasets, including the built-in Diabetes dataset from the sklearn library, and explore how to apply statistical methods to datasets from sources like Kaggle and the UCI Machine Learning Repository.

Basic Concepts in Statistics

Population and Sample

- Population: The complete set of all possible observations or measurements.
- Sample: A subset of the population selected for the actual study.

Need for Sampling in Statistics

- Cost Efficiency
- Time Efficiency
- Feasibility
- Manageability

Benefits of Sampling

- Accuracy
- Less Data Overhead
- Focused Research

Descriptive Statistics

Measures of Central Tendency

- Mean: The average value.
- Median: The middle value.
- Mode: The most frequent value.

Measures of Variability

- Standard Deviation: Measures the amount of variation.
- Variance: The average of the squared differences from the mean.
- Range: The difference between the maximum and minimum values.

Additional Descriptive Measures

- Percentiles and Quartiles
- Interquartile Range (IQR)
- Skewness
- Kurtosis

Inferential Statistics

- Hypothesis Testing
 - Confidence Intervals
 - Regression Analysis
 - ANOVA (Analysis of Variance)
 - Chi-Square Test

Getting Started with Python for Statistical Analysis

Before we dive into coding, make sure you have Python installed along with the necessary libraries. You can install the required libraries using the following commands:

```
In [1]: pip install scikit-learn
Requirement already satisfied: scikit-learn in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (1.5.1)
Requirement already satisfied: numpy>=1.19.5, <=1.8 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from scikit-learn) (2.1.0)
Requirement already satisfied: scipy>=1.6.0 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from scikit-learn) (1.14.1)
Requirement already satisfied: joblib>=1.2.0 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from scikit-learn) (3.5.0)
[notice] A new release of pip is available: 24.0 -> 24.2
[notice] To update, run: python.exe -m pip install --upgrade pip

In [46]: pip install statsmodels
Requirement already satisfied: statsmodels in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (0.14.2)
Requirement already satisfied: numpy>=1.22.3 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from statsmodels) (2.1.0)
Requirement already satisfied: matplotlib>=3.6.1, <=3.4 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from statsmodels) (3.9.2)
Requirement already satisfied: pandas>=2.1.0, <=1.4 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from statsmodels) (2.2.2)
Requirement already satisfied: patsy>=0.5.4 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from statsmodels) (0.5.2.0)
Requirement already satisfied: packaging>=21.3 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from statsmodels) (24.1)
Requirement already satisfied: python-dateutil>=2.8.0, <=2.9 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from pandas=>2.1.0, <=1.4->statsmodels) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from pandas=>2.1.0, <=1.4->statsmodels) (2024.1)
Requirement already satisfied: six in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from patsy=>0.5.4->statsmodels) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
[notice] A new release of pip is available: 24.0 -> 24.2
[notice] To update, run: python.exe -m pip install --upgrade pip

In [8]: pip install seaborn
Collecting seaborn
  Downloading seaborn-0.12.2-py3-none-any.whl.metadata (5.4 KB)
Requirement already satisfied: numpy>=1.24.0, <=1.20 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from seaborn) (2.1.0)
Requirement already satisfied: pandas>=1.2 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from seaborn) (2.2.2)
Requirement already satisfied: matplotlib>=3.6.1, <=3.4 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from seaborn) (3.9.2)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (1.3.0)
Requirement already satisfied: cycler>=0.10 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (4.53.1)
Requirement already satisfied: kiwisolver>=3.1 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (1.4.7)
Requirement already satisfied: packaging>=20.0 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (24.1)
Requirement already satisfied: pyparsing>=2.3.1 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (3.1.4)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from matplotlib=>3.6.1, <=3.4->seaborn) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from pandas=>1.2, <=1.4->seaborn) (2024.1)
Requirement already satisfied: tzdata>=2022.1 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from pandas=>1.2, <=1.4->seaborn) (2024.1)
Requirement already satisfied: six>=1.5 in c:\users\amrut\appdata\local\programs\python\python32\lib\site-packages (from python-dateutil=>2.7, <=2.9->matplotlib=>3.6.1, <=3.4->seaborn) (1.16.0)
Downloading seaborn-0.13.2-py3-none-any.whl (1294 KB)
-----0.0/294.9 KB ? eta -:-:--
-----71.7/294.9 KB 603.6 KB/s eta 0:00:01
-----143.4/294.9 KB 1.1 MB/s eta 0:00:01
-----194.6/294.9 KB 1.1 MB/s eta 0:00:01
-----286.7/294.9 KB 1.3 MB/s eta 0:00:01
-----294.9/294.9 KB 1.2 MB/s eta 0:00:00
Installing collected packages: seaborn
Successfully installed seaborn-0.13.2
[notice] A new release of pip is available: 24.0 -> 24.2
[notice] To update, run: python.exe -m pip install --upgrade pip
```

Applying Statistical Methods to Datasets

Descriptive Statistics in Python

Descriptive statistics help us summarize and describe the main features of a dataset.

Inferential Statistics in Python

Inferential statistics involve drawing conclusions about a population based on a sample.

Applying Statistical Methods

- Load the dataset into Python using pandas.

```
In [11]: # Sample Code for Statistical Analysis in Python
# Loading the Diabetes Dataset
import pandas as pd
from sklearn.datasets import load_diabetes

# Load the dataset
diabetes = load_diabetes()
df = pd.DataFrame(diabetes.data, columns=diabetes.feature_names)
df['target'] = diabetes.target

# Display the first few rows
print(df.head())

age    sex    bmi    bp    a1    a2    a3    \
0    0.03876  0.050480  0.061496  0.021872 -0.044223 -0.034821 -0.043401
1    -0.030182  0.044441 -0.031474 -0.025228 -0.004449 -0.033063 -0.074412
2    0.085299  0.050680  0.044451 -0.005670 -0.048599 -0.034194 -0.032356
3    -0.089043 -0.044441 -0.011395 -0.033656  0.011191  0.024991 -0.034038
4    0.005383 -0.044441 -0.030385  0.021872  0.003935  0.050596  0.008142

age    sex    bmi    bp    a1    a2    a3    \
0    -0.002592  0.013907 -0.017646   151.0
1    -0.030493 -0.048332 -0.092054    75.0
2    -0.002592  0.002861 -0.002930   141.0
3    0.034359  0.021688 -0.009382   205.0
4    -0.002592 -0.011881 -0.046561   135.0

# Additional descriptive statistics
print("Mean:", df.mean())
print("Median:", df.median())
print("Standard Deviation:", df.std())
print("Variance:", df.var())
print("Range:", df.max() - df.min())
print("Skewness:", df.skew())
print("Kurtosis:", df.kurt())

Mean:
age    -1.444235e-18
sex     2.143215e-18
bmi     2.155925e-16
bp      -4.854084e-17
a1      -1.428584e-17
a2      3.198818e-17
a3      -6.028330e-18
a4      -1.788100e-17
a5      9.243484e-17
a6      1.103715e-17
target  1.502315e+02
dtype: float64

Median:
age    0.005383
sex    -0.044442
bmi    -0.007248
bp     -0.005870
a1     -0.004323
a2     -0.003819
a3     -0.005854
a4     -0.002592
a5     -0.001847
a6     -0.001078
target  140.500000
dtype: float64

Variance:
age    0.016281
sex    0.046462
bmi    0.000786
bp     0.000099
a1     0.000344
a2     0.000101
a3     0.000345
a4     0.000493
a5     0.001114
a6     0.000064
target  72.000000
dtype: float64

Standard Deviation:
age    0.047619
sex    0.047619
bmi    0.047619
bp     0.047619
a1     0.047619
a2     0.047619
a3     0.047619
a4     0.047619
a5     0.047619
a6     0.047619
target  77.000000
dtype: float64

Range:
age    0.217952
sex    0.093322
bmi    0.260831
bp     0.264462
a1     0.280494
a2     0.314402
a3     0.281486
a4     0.261629
a5     0.259494
a6     0.273719
target  321.000000
dtype: float64

Skewness:
age    -0.231382
sex     0.127385
bmi     0.580145
bp      0.290658
a1      0.770108
a2      0.436592
a3      0.195220
a4      0.735714
a5      0.281794
a6      0.070717
target  0.440563
dtype: float64

Kurtosis:
age    -0.671224
sex    -1.992811
bmi     0.950594
bp     -0.532797
a1     -0.322648
a2     0.491381
a3     0.902507
a4     0.444002
a5     -0.124607
a6     0.230571
target -0.883057
dtype: float64
```

calculate T-statistic and p_value

```
In [13]: # Performing Inferential Statistics
from scipy import stats

# Example data: BMI values
bmi_values = df['bmi']

# Hypothesized population mean for BMI
population_mean = 0.05

# Perform one-sample t-test
t_stat, p_value = stats.ttest_1samp(bmi_values, population_mean)

print(f"T-statistic: {t_stat}")
print(f"P-value: {p_value}")

T-statistic: -22.0740809487010174
P-value: 2.7634312235044838e-73
```

Calculate a 95% confidence interval for the mean of a selected feature.

```
In [14]: # 95% Confidence Interval
from scipy import stats

# Sample mean and standard error for BMI
sample_mean = np.mean(bmi_values)
standard_error = stats.sem(bmi_values)

# Compute 95% confidence interval for BMI
confidence_interval = stats.norm.interval(0.95, loc=sample_mean, scale=standard_error)

print(f"95% Confidence Interval for BMI: {confidence_interval}")

95% Confidence Interval for BMI: (np.float64(0.004439332370169161), np.float64(0.0044393323701686915))
```

calculate regression

```
In [15]: # OLS Regression Analysis
import statsmodels.api as sm

# Define independent variable (add constant for intercept)
X = sm.add_constant(df['age'])

# Define dependent variable
y = df['target']

# Fit linear regression model
model = sm.OLS(y, X).fit()

# Print model summary
print(model.summary())

OLS Regression Results
=====
Dep. Variable:    target    R-squared:    0.348
Model:            OLS     Adj. R-squared:  0.342
Date:            Thu, 05 Sep 2024    Prob (F-statistic):    3.47e-42
Time:            19:47:33    Log-likelihood:    -2481.0
No. Observations:    482    AIC:            4921.
DF Residuals:        480    BIC:            4920.
DF Model:            1
Covariance Type:    nonconstant

=====
coef    std err    t    P>|t|    [0.025    0.975]
-----
const    152.1335    2.974    51.162    0.000    146.289    157.978
age      849.4359    62.515    13.597    0.000    826.570    1072.301
=====
Omnibus:    11.674    Durbin-Watson:    1.848
Prob(Omnibus):    0.003    Jarque-Bera (JB):    7.105
Skew:       0.156    Prob(JB):    0.0259
Kurtosis:    3.060    Cond. No.      11.19
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
```

Create Visualizations to illustrate the relationships between variables

```
In [19]: # 4.6 Data Visualization
import seaborn as sns
import matplotlib.pyplot as plt

# Pairplot with regression lines
sns.pairplot(df, kind='reg')

plt.figure(figsize=(10, 8))
sns.heatmap(df.corr(), annot=True, cmap='magma', cbar=True)

plt.figure(figsize=(8, 6))
sns.histplot(df['age'], kde=True)
plt.title('Distribution of Age')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()

plt.figure(figsize=(8, 6))
sns.regplot(df['age'], df['target'], dataargs=(df, 'age', 'target'), scatter_kws={'s': 100, 'line_kws': {'color': 'red'}})
plt.xlabel('Age')
plt.ylabel('Target')
plt.show()
```



