

```

1 #Mount Google drive
2 from google.colab import drive
3 drive.mount('/content/drive')

```

↗ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True)

```

1 import pandas as pd
2 import numpy as np
3
4 # Load the dataset from the datafile provided
5 df = pd.read_csv('/content/drive/MyDrive/jyoti/PupilBioTest_PMP_revA.csv')
6 print(f"DataFrame size: {df.size}")
7
8 # Shuffling Dataset to diversify the Tissue part
9 df_shuffled = df.sample(frac=1)
10
11 #Reducing the size of the Dataframe to 1 million consume less memory.
12 df = df_shuffled.head(1000000)
13
14 #clearing the useless dataframe
15 df_shuffled = None
16 print(f"New Data Size: {df.size}")
17

```

↗ DataFrame size: 200098379  
New Data Size: 13000000

```

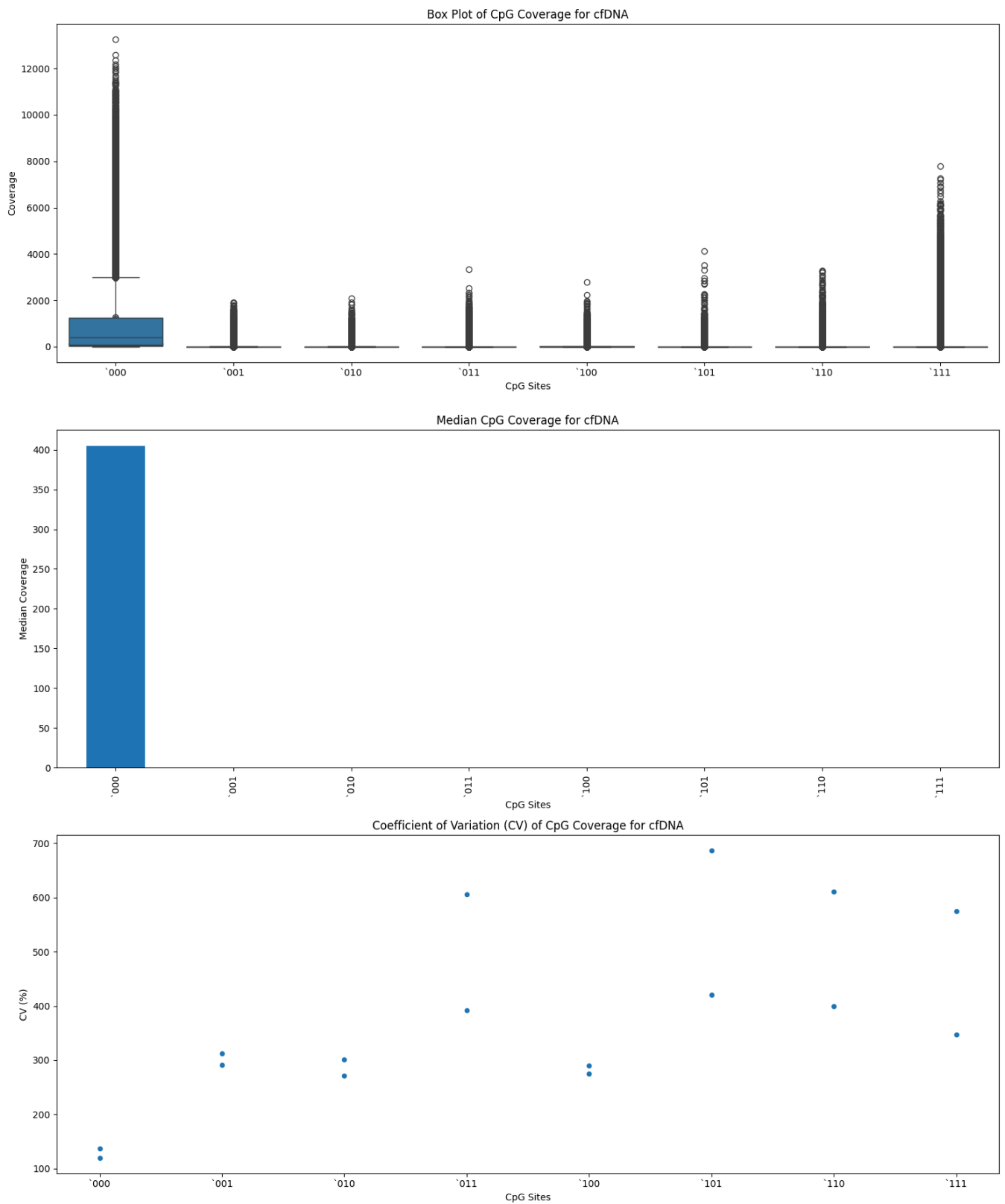
1 #Question 1.1.1
2 # Ensure data is clean and organized by removing any missing values
3 df.dropna(inplace=True)
4
5 # Group the data by Tissue to calculate statistics for each tissue type separately
6 grouped = df.groupby('Tissue')
7
8 # Function to calculate the Coefficient of Variation (CV)
9 def calculate_cv(series):
10     mean = series.mean()
11     std_dev = series.std()
12     cv = (std_dev / mean) * 100 if mean != 0 else 0
13     return cv
14
15 # Dictionary to store results for each tissue type
16 results = {}
17
18 # Loop through each tissue type and calculate median and CV for each CpG site
19 for tissue, group in grouped:
20     stats = {}
21     for col in group.columns[2:10]: # Coverage columns `000` to `111`
22         median = group[col].median() # Calculate the median
23         cv = calculate_cv(group[col]) # Calculate the CV
24         stats[col] = {'median': median, 'cv': cv}
25     results[tissue] = stats
26
27 # Display the results
28 for tissue, stats in results.items():
29     print(f"Tissue: {tissue}")
30     for cp_g_site, metrics in stats.items():
31         print(f"  CpG Site: {cp_g_site}, Median: {metrics['median']}, CV: {metrics['cv']:.2f} %")
32

```

↗ Tissue: Islet  
 CpG Site: `000`, Median: 62.0, CV: 119.72%  
 CpG Site: `001`, Median: 0.0, CV: 290.77%  
 CpG Site: `010`, Median: 0.0, CV: 270.79%  
 CpG Site: `011`, Median: 0.0, CV: 391.62%  
 CpG Site: `100`, Median: 0.0, CV: 274.56%  
 CpG Site: `101`, Median: 0.0, CV: 421.08%  
 CpG Site: `110`, Median: 0.0, CV: 400.01%  
 CpG Site: `111`, Median: 0.0, CV: 346.63%  
 Tissue: cfDNA  
 CpG Site: `000`, Median: 405.0, CV: 137.03%  
 CpG Site: `001`, Median: 0.0, CV: 312.65%  
 CpG Site: `010`, Median: 0.0, CV: 301.40%  
 CpG Site: `011`, Median: 0.0, CV: 605.23%  
 CpG Site: `100`, Median: 0.0, CV: 290.19%  
 CpG Site: `101`, Median: 0.0, CV: 686.40%  
 CpG Site: `110`, Median: 0.0, CV: 611.41%  
 CpG Site: `111`, Median: 0.0, CV: 575.27%

```
1 %matplotlib inline
```

```
1 #Question 1.1.2
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 grouped = df.groupby('Tissue')
6 # Function to calculate CV (reused from above)
7 def calculate_cv(series):
8     mean = series.mean()
9     std_dev = series.std()
10    cv = (std_dev / mean) * 100 if mean != 0 else 0
11    return cv
12
13 # Initialize a figure for plots
14 fig, axes = plt.subplots(3, 1, figsize=(15, 18))
15
16 # Plot for each tissue type
17 for i, (tissue, group) in enumerate(grouped):
18     # Box plot for distribution of coverage
19     sns.boxplot(data=group.iloc[:, 2:10], ax=axes[0])
20     axes[0].set_title(f'Box Plot of CpG Coverage for {tissue}')
21     axes[0].set_xlabel('CpG Sites')
22     axes[0].set_ylabel('Coverage')
23
24     # Calculate median and CV for each CpG site
25     medians = group.iloc[:, 2:10].median()
26     cvs = group.iloc[:, 2:10].apply(calculate_cv)
27
28     # Bar plot for median coverage
29     medians.plot(kind='bar', ax=axes[1])
30     axes[1].set_title(f'Median CpG Coverage for {tissue}')
31     axes[1].set_xlabel('CpG Sites')
32     axes[1].set_ylabel('Median Coverage')
33
34     # Convert CVs Series to DataFrame for scatter plot
35     cvs_df = cvs.reset_index()
36     cvs_df.columns = ['CpG Site', 'CV']
37
38     # Scatter plot for CV
39     cvs_df.plot(kind='scatter', x='CpG Site', y='CV', ax=axes[2])
40     axes[2].set_title(f'Coefficient of Variation (CV) of CpG Coverage for {tissue}')
41     axes[2].set_xlabel('CpG Sites')
42     axes[2].set_ylabel('CV (%)')
43
44 # Adjust the layout
45 plt.tight_layout()
46 plt.show()
47
```



```

1 # Question 1.2.1
2 from scipy.stats import ttest_ind
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import classification_report, confusion_matrix
6
7 # Function to calculate p-values for each PMP using t-test
8 def calculate_p_values(df, tissue1_label):
9     p_values = {}
10    tissue1 = df[df['Tissue'] == tissue1_label]
11    other_tissues = df[df['Tissue'] != tissue1_label]
12    for col in df.columns[2:10]: # PMP coverage columns
13        _, p_value = ttest_ind(tissue1[col], other_tissues[col], equal_var=False)
14        p_values[col] = p_value
15    return p_values
16
17 # Calculate p-values for Tissue #1 (assuming Tissue #1 is labeled as 0)
18 p_values = calculate_p_values(df, tissue1_label=0)
19
20 # Filter PMPs with significant p-values (e.g., p < 0.05)
21 significant_pmps = {pmp: p for pmp, p in p_values.items() if p < 0.05}
22 print("Significant PMPs with p-values:", significant_pmps)
23
24 # Prepare the data for machine learning
25 X = df.iloc[:, 2:10] # PMP coverage columns
26 y = df['Tissue']
27 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
28
29 # Train a RandomForestClassifier
30 model = RandomForestClassifier(random_state=42)
31 model.fit(X_train, y_train)
32
33 # Predict on test data
34 y_pred = model.predict(X_test)
35
36 # Evaluate model performance
37 print("Classification Report:\n", classification_report(y_test, y_pred))
38 print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
39
40 # Identify important features (PMPs) based on feature importances
41 fi = model.feature_importances_
42 important_pmps = [X.columns[i] for i in range(len(fi)) if fi[i] > 0.05]
43 print("Important PMPs:", important_pmps)
44

```

```

↗ Significant PMPs with p-values: {}
Classification Report:
              precision    recall  f1-score   support

      Islet         0.69       0.40       0.51       74219
      cfDNA         0.83       0.94       0.88      225781

 accuracy         0.76       0.67       0.81      300000
 macro avg         0.76       0.67       0.69      300000
weighted avg         0.79       0.81       0.79      300000

Confusion Matrix:
[[ 29802  44417]
 [ 13448  212333]]
Important PMPs: ['`000', '`001', '`010', '`100', '`111']

```

```
1 print(df.head())
```

```

↗ <bound method NDFrame.head of
2300065    r    8736:8740:8807  2774    1    2    0    0    0    0    0    0
12296074    f  13360:13368:13409  2924    1    0    0    15    0    0
6021080    r    7235:7301:7364    47    0   12    0    0    0    0    0
8746229    f    8565:8573:8682   826    7    1    0    0    0    0    0
12144718    f    8500:8553:8598   608    0   18    0    0    0    0    0
...
15135198    f    8464:8479:8517    42    0    0    0    0    0    0    0
6738561    r  13415:13498:13505    72    4    0    1    1    8    0
14322977    f    8333:8388:8420    23    2    0    0    0    0    0    0
9347203    f    8591:8632:8692   442    0    0    0    0    0    0    0
8408148    f    9873:9877:9957  1027    0    0    0    0    0    0    0

      `111  Sample_ID  Replicate  Tissue
2300065    0          37          Rep2  cfDNA

```

```

12296074      0      38      Rep1      cfDNA
6021080       0      50      Rep1      Islet
8746229       0      15      Rep2      cfDNA
12144718      0      47      Rep1      cfDNA
...          ...      ...      ...      ...
15135198      0      76      Rep1      Islet
6738561       20     62      Rep1      Islet
14322977      0      62      Rep1      Islet
9347203       0      25      Rep2      cfDNA
8408148       0       9      Rep2      cfDNA

```

```
[1000000 rows x 13 columns]>
```

```

1 #Question 1.2.2
2
3 # Specify the range of columns that represent the PMPs
4 pmp_columns = df.columns[2:10]
5
6 # Calculate variant read fractions (VRF) for each PMP
7 for pmp in pmp_columns:
8     df[f'vrf_{pmp}'] = df[pmp] / df[pmp_columns].sum(axis=1)
9
10 # Calculate mean VRF for each PMP in both tissues
11 mean_vrf_per_pmp_tissue = df.groupby('Tissue')[[f'vrf_{pmp}' for pmp in pmp_columns]].mean().reset_index()
12
13 print(mean_vrf_per_pmp_tissue)
14
15
16

```

```

→ Tissue  vrf_`000  vrf_`001  vrf_`010  vrf_`011  vrf_`100  vrf_`101  \
0  Islet   0.844559  0.024278  0.020356  0.014883  0.021653  0.012468
1  cfDNA   0.907256  0.015745  0.013760  0.007599  0.016339  0.004126

      vrf_`110  vrf_`111
0   0.014122  0.047681
1   0.007447  0.027728

```

### Question 1.3.

a. How does sequencing depth affect specificity confidence?

Sequencing depth affects the confidence in detecting true methylation patterns by reducing noise and increasing the accuracy of measurements. Higher sequencing depth allows for better discrimination between true signals and sequencing errors, thus improving specificity.

In practical terms, higher sequencing depth ensures that even low-abundance methylation patterns are detected with greater reliability, minimizing false positives and increasing the specificity of biomarkers.

b. For the top 10 PMPs, estimate the threshold of reads required to confidently call Tissue #2 at a sequencing depth of 1 million reads

We can estimate the threshold of reads required to confidently call Tissue #2 by analyzing the distribution of read counts and calculating the minimum number of reads needed for reliable detection. This involves looking at the read coverage for the top 10 PMPs and determining the cutoff where the detection is robust.

```

1 #Question 1.3.3
2 # Function to evaluate specificity
3 def evaluate_specificity(model, X_test, y_test):
4     y_pred = model.predict(X_test)
5     return classification_report(y_test, y_pred, output_dict=True)
6
7 # Evaluate specificity for the top 10 PMPs
8 X_top_pmps = X[top_10_pmps]
9 X_train_top_pmps, X_test_top_pmps, y_train_top_pmps, y_test_top_pmps = train_test_split(X_top_pmps, y, test_size=0.3, rand
10 model_top_pmps = RandomForestClassifier(random_state=42)
11 model_top_pmps.fit(X_train_top_pmps, y_train_top_pmps)
12 specificity_top_pmps = evaluate_specificity(model_top_pmps, X_test_top_pmps, y_test_top_pmps)
13
14 # Evaluate specificity for individual CpG sites
15 specificity_individual = {}
16 for col in X.columns:
17     X_individual = X[[col]]
18     X_train_individual, X_test_individual, y_train_individual, y_test_individual = train_test_split(X_individual, y, test_
19     model_individual = RandomForestClassifier(random_state=42)
20     model_individual.fit(X_train_individual, y_train_individual)
21     specificity_individual[col] = evaluate_specificity(model_individual, X_test_individual, y_test_individual)
22
23 # Print the results
24 print("Specificity for top 10 PMPs:\n", specificity_top_pmps)
25 print("Specificity for individual CpG sites:\n", specificity_individual)

```