# load-diabetes

June 20, 2024

```
[5]: # Import necessary libraries
     from sklearn.datasets import load_diabetes
     import pandas as pd
     import matplotlib.pyplot as plt
     import seaborn as sns
     from sklearn.model_selection import train_test_split
     from sklearn.linear_model import LinearRegression
     from sklearn.metrics import mean_squared_error

     # Load dataset
     diabetes = load_diabetes()
     X, y = diabetes.data, diabetes.target

     # Convert data to DataFrame for easier analysis
     df = pd.DataFrame(data=X, columns=diabetes.feature_names)
     df['target'] = y

     # Display basic statistics and information
     print(f"Dataset shape: {df.shape}")
     print(f"Columns: {df.columns}")
     print(f"Target variable summary:\n{df['target'].describe()}")

     # Display correlation heatmap
     plt.figure(figsize=(10, 8))
     sns.heatmap(df.corr(), annot=True, cmap='coolwarm', center=0)
     plt.title('Correlation Heatmap')
     plt.show()

     # Pairplot for visualizing relationships and distributions
     sns.pairplot(df, diag_kind='hist')
     plt.suptitle('Pairplot of Diabetes Dataset Features', y=1.02)
     plt.tight_layout()
     plt.show()

     # Train-test split
     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,␣
       ↪random_state=42)
```

```python
# Train Linear Regression model
lr = LinearRegression()
lr.fit(X_train, y_train)

# Predict and evaluate
y_pred = lr.predict(X_test)
print(f"\nMean Squared Error: {mean_squared_error(y_test, y_pred)}")
```
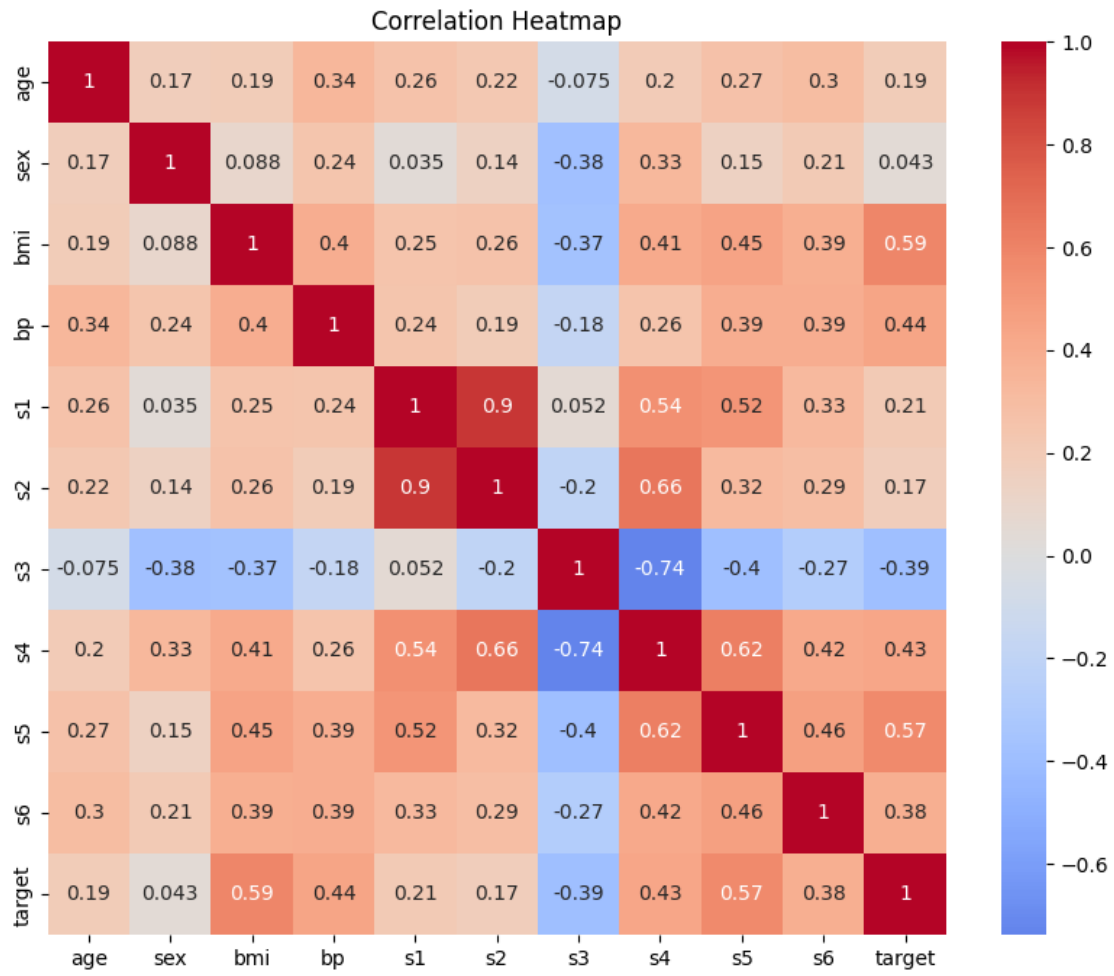
```
Dataset shape: (442, 11)
Columns: Index(['age', 'sex', 'bmi', 'bp', 's1', 's2', 's3', 's4', 's5', 's6',
       'target'],
      dtype='object')
Target variable summary:
count    442.000000
mean     152.133484
std       77.093005
min       25.000000
25%       87.000000
50%      140.500000
75%      211.500000
max      346.000000
Name: target, dtype: float64
```
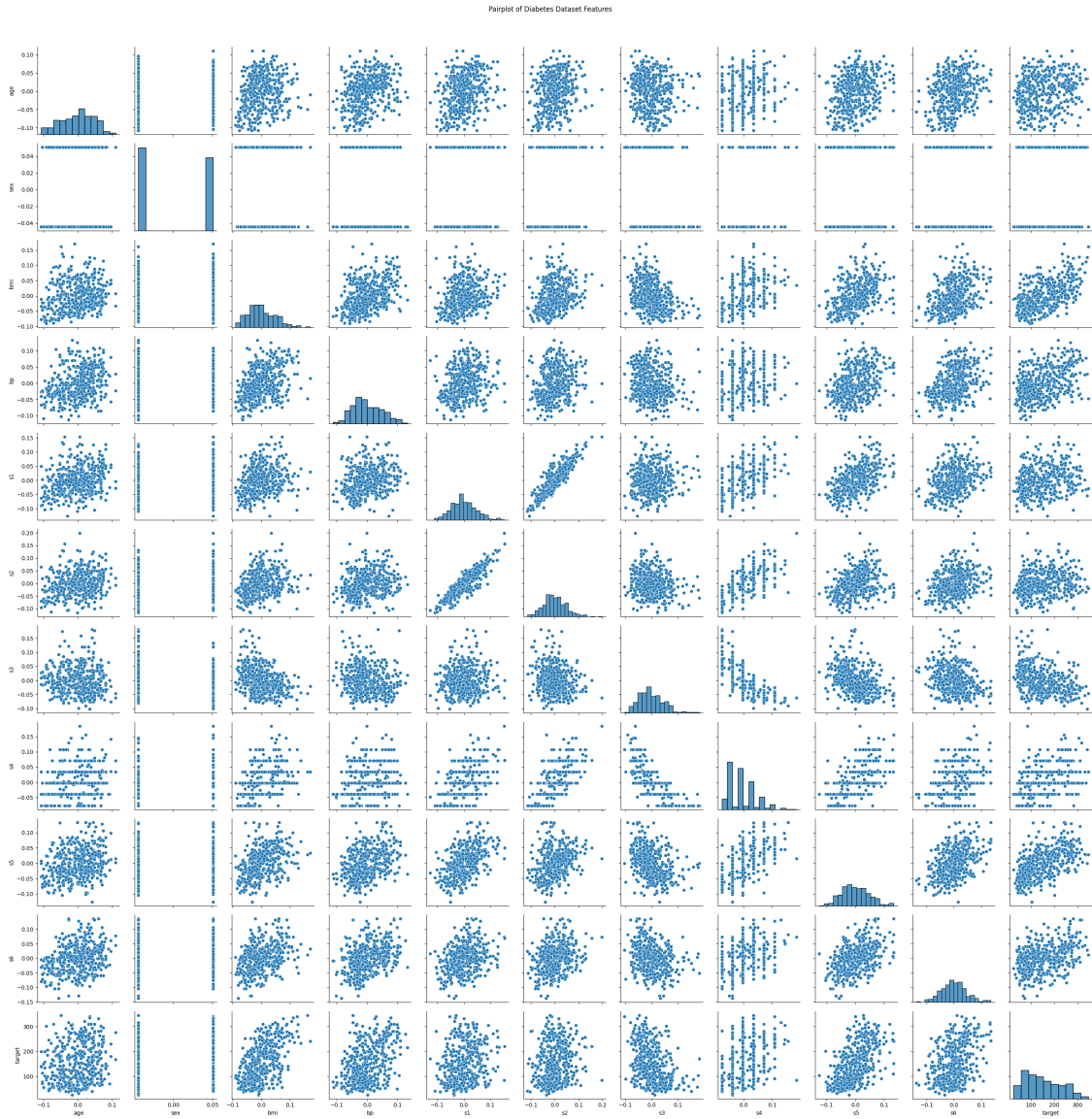
Correlation Heatmap

Pairplot of Diabetes Dataset Features

Mean Squared Error: 2900.193628493482