

## Week 4 — Data Acquisition & Description Importing all necessary and Collection of Data

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.preprocessing import OneHotEncoder, MinMaxScaler, StandardScaler, PowerTransformer
6 from sklearn.pipeline import Pipeline
7 from sklearn.impute import SimpleImputer
8 from scipy.sparse import hstack
9 from sklearn.compose import ColumnTransformer
10 import seaborn as sns
11 from sklearn.model_selection import GridSearchCV #for hypertuning
12 from sklearn.linear_model import LinearRegression, LogisticRegression, Lasso, Ridge
13 from lightgbm import LGBMRegressor
14

```

### Importing csv File

```

1 # Uploading A csv file
2 import os, pandas as pd
3
4 # CANDIDATES = ["exams.csv", "./exams.csv", "/content/exams.csv", "/mnt/data/exams.csv"]
5 # df = None
6 # for p in CANDIDATES:
7 #     if os.path.exists(p):
8 #         df = pd.read_csv(p); print(f"Loaded: {p}"); break
9
10 # if df is None:
11 try:
12     from google.colab import files
13     up = files.upload()
14     fname = next(iter(up.keys()))
15     # df = pd.read_csv(fname) # This was trying to read an Excel file as CSV
16     df = pd.read_excel(fname) # Use pd.read_excel for Excel files
17     print(f"Uploaded and loaded: {fname}")
18 except Exception as e:
19     # raise FileNotFoundError("Upload exams.csv or put it beside the notebook.")
20     print(f"An error occurred: {e}")
21
22
23 df.head()

```

Choose Files cvs-student...0\_2025.xlsx

**cvs-student\_performance\_10\_30\_2025.xlsx**(application/vnd.openxmlformats-officedocument.spreadsheetml.sheet) - 64561 bytes, last modified: 11/20/2025 - 100% done

Saving cvs-student\_performance\_10\_30\_2025.xlsx to cvs-student\_performance\_10\_30\_2025.xlsx

Uploaded and loaded: cvs-student\_performance\_10\_30\_2025.xlsx

|   | math score | reading score | writing score | pass_flag | gender_encoded | lunch_standard | test preparation course_none | race/ethnicity_group | race/eB |
|---|------------|---------------|---------------|-----------|----------------|----------------|------------------------------|----------------------|---------|
| 0 | 72         | 72            | 74            | 1         | 1              | 1              | 1                            | 1                    | 1       |
| 1 | 69         | 90            | 88            | 1         | 1              | 1              | 0                            | 0                    | 0       |
| 2 | 90         | 95            | 93            | 1         | 1              | 1              | 1                            | 1                    | 1       |
| 3 | 47         | 57            | 44            | 0         | 0              | 0              | 1                            | 0                    | 0       |
| 4 | 76         | 78            | 75            | 1         | 0              | 1              | 1                            | 0                    | 0       |

Next steps: [Generate code with df](#) [New interactive sheet](#)

Double-click (or enter) to edit

```
1 from google.colab import sheets
2 sheet = sheets.InteractiveSheet(df=df)
```

[https://docs.google.com/spreadsheets/d/1Geg2P4IGDUBHL\\_5TpZufGcfkAmZq4AUGoJnhbVoCgkY/edit#gid=0](https://docs.google.com/spreadsheets/d/1Geg2P4IGDUBHL_5TpZufGcfkAmZq4AUGoJnhbVoCgkY/edit#gid=0)

|    | A  | B  | C  | D | E | F | G | H |
|----|----|----|----|---|---|---|---|---|
| 35 | 74 | 81 | 83 | 1 | 1 | 1 | 1 | 1 |
| 39 | 50 | 64 | 59 | 1 | 1 | 0 | 1 |   |
| 40 | 75 | 90 | 88 | 1 | 1 | 0 | 0 | 0 |
| 41 | 57 | 56 | 57 | 1 | 0 | 0 | 1 |   |
| 42 | 55 | 61 | 54 | 1 | 0 | 0 | 1 |   |
| 43 | 58 | 73 | 68 | 1 | 1 | 1 | 1 |   |
| 44 | 53 | 58 | 65 | 1 | 1 | 1 | 1 |   |
| 45 | 59 | 65 | 66 | 1 | 0 | 0 | 0 |   |
| 46 | 50 | 56 | 54 | 1 | 1 | 0 | 1 |   |
| 47 | 65 | 54 | 57 | 1 | 0 | 1 | 1 |   |
| 48 | 55 | 65 | 62 | 1 | 1 | 1 | 0 |   |
| 49 | 66 | 71 | 76 | 1 | 1 | 1 | 1 |   |
| 50 | 57 | 74 | 76 | 1 | 1 | 0 | 0 |   |
| 51 | 82 | 84 | 82 | 1 | 0 | 1 | 0 |   |
| 52 | 53 | 55 | 48 | 0 | 0 | 1 | 1 |   |
| 53 | 77 | 69 | 68 | 1 | 0 | 0 | 0 |   |
| 54 | 53 | 44 | 42 | 0 | 0 | 1 | 1 |   |

+    Sheet1 ▾

First let's analyze the dataset And Upload The csv File

```
1 from google.colab import sheets
2 sheet = sheets.InteractiveSheet(df=df)
```

[https://docs.google.com/spreadsheets/d/1NKB-D02xAm97xF0mHxAfiYYtAReASMs\\_yisZWjv8tE/edit#gid=0](https://docs.google.com/spreadsheets/d/1NKB-D02xAm97xF0mHxAfiYYtAReASMs_yisZWjv8tE/edit#gid=0)

G Get professional email like "@your-company.com"

Plus 2 TB of storage per user, longer video calls, and more with Google Workspace.

Try Workspace



File Edit View Insert Format Data Tools Extensions Help

Q Menus ⌘ ⌘ ⌘ 100% \$ % 0 .00 123 Default... - 10 + B I S A

Here the dataset have categorical as well as numerical data. Categorical data includes race/ethnicity, parental level of education, lunch, test preparation course while numerical data includes math score, reading score, writing score.

```
1 df.shape, df.dtypes, df.isna().sum(), df.describe(include="all").T
2
```

|                              | 10 | 22 | 20    | 0  | 1 | 0 | 1 | 1 |
|------------------------------|----|----|-------|----|---|---|---|---|
| ((1000, 11),                 | 46 | 42 | 46    | 0  | 0 | 0 | 0 | 0 |
| math score                   | 54 | 58 | int64 | 61 | 1 | 1 | 0 | 1 |
| reading score                | 66 | 69 | int64 | 63 | 1 | 0 | 1 | 1 |
| writing score                | 65 | 75 | int64 | 70 | 1 | 1 | 0 | 0 |
| pass_flag                    | 44 | 54 | int64 | 53 | 0 | 0 | 1 | 1 |
| gender_encoded               | 69 | 73 | int64 | 73 | 1 | 1 | 1 | 1 |
| lunch_standard               | 74 | 71 | int64 | 80 | 1 | 0 | 0 | 0 |
| test preparation course_none | 74 | 74 | int64 | 72 | 1 | 0 | 0 | 1 |
| race/ethnicity_group B       | 73 | 74 | int64 | 55 | 1 | 0 | 1 | 1 |
| race/ethnicity_group C       | 66 | 69 | int64 | 75 | 1 | 1 | 1 | 1 |
| race/ethnicity_group D       | 69 | 70 | int64 | 65 | 1 | 0 | 1 | 1 |
| race/ethnicity_group E       | 69 | 70 | int64 | 70 | 0 | 1 | 1 | 1 |
| dtype: object,               | 70 | 70 |       | 65 | 1 | 0 | 1 | 1 |
| math score                   | 62 | 70 |       | 75 | 1 | 1 | 1 | 1 |
| reading score                | 69 | 74 |       | 74 | 1 | 1 | 1 | 1 |
| writing score                | 74 | 0  |       | 0  | 0 | 1 | 1 | 1 |
| pass_flag                    | 0  |    |       |    |   |   |   |   |
| gender_encoded               | 0  |    |       |    |   |   |   |   |
| lunch_standard               | 0  |    |       |    |   |   |   |   |
| test preparation course_none | 0  |    |       |    |   |   |   |   |
| race/ethnicity_group B       | 0  |    |       |    |   |   |   |   |
| race/ethnicity_group C       | 0  |    |       |    |   |   |   |   |
| race/ethnicity_group D       | 0  |    |       |    |   |   |   |   |
| race/ethnicity_group E       | 0  |    |       |    |   |   |   |   |
| dtype: int64,                |    |    |       |    |   |   |   |   |

|                              | count  | mean   | std       | min  | 25%   | 50%  | \ |
|------------------------------|--------|--------|-----------|------|-------|------|---|
| math score                   | 1000.0 | 66.089 | 15.163080 | 0.0  | 57.00 | 66.0 |   |
| reading score                | 1000.0 | 69.169 | 14.600192 | 17.0 | 59.00 | 70.0 |   |
| writing score                | 1000.0 | 68.054 | 15.195657 | 10.0 | 57.75 | 69.0 |   |
| pass_flag                    | 1000.0 | 0.812  | 0.390908  | 0.0  | 1.00  | 1.0  |   |
| gender_encoded               | 1000.0 | 0.518  | 0.499926  | 0.0  | 0.00  | 1.0  |   |
| lunch_standard               | 1000.0 | 0.645  | 0.478753  | 0.0  | 0.00  | 1.0  |   |
| test preparation course_none | 1000.0 | 0.642  | 0.479652  | 0.0  | 0.00  | 1.0  |   |
| race/ethnicity_group B       | 1000.0 | 0.190  | 0.392497  | 0.0  | 0.00  | 0.0  |   |
| race/ethnicity_group C       | 1000.0 | 0.319  | 0.466322  | 0.0  | 0.00  | 0.0  |   |
| race/ethnicity_group D       | 1000.0 | 0.262  | 0.439943  | 0.0  | 0.00  | 0.0  |   |
| race/ethnicity_group E       | 1000.0 | 0.140  | 0.347161  | 0.0  | 0.00  | 0.0  |   |

|                              | 75%  | max   |
|------------------------------|------|-------|
| math score                   | 77.0 | 100.0 |
| reading score                | 79.0 | 100.0 |
| writing score                | 79.0 | 100.0 |
| pass_flag                    | 1.0  | 1.0   |
| gender_encoded               | 1.0  | 1.0   |
| lunch_standard               | 1.0  | 1.0   |
| test preparation course_none | 1.0  | 1.0   |
| race/ethnicity_group B       | 0.0  | 1.0   |
| race/ethnicity_group C       | 1.0  | 1.0   |
| race/ethnicity_group D       | 1.0  | 1.0   |
| race/ethnicity_group E       | 0.0  | 1.0   |

## Week 5 — Preprocessing & Feature Engineering

```
1 # Week 5 – Feature Engineering, Encoding, Split, Save
2 import pandas as pd
3 from sklearn.preprocessing import LabelEncoder
```

```

4 from sklearn.model_selection import train_test_split
5
6 # start from original data
7 work = df.copy()
8
9 # sanity check: required columns present
10 # The original categorical columns are already one-hot encoded, so this assertion is not needed.
11 # required = [
12 #     "gender", "race/ethnicity", "parental level of education",
13 #     "lunch", "test preparation course",
14 #     "math score", "reading score", "writing score"
15 # ]
16 # missing = [c for c in required if c not in work.columns]
17 # assert not missing, f"Missing columns: {missing}"
18
19 # safe engineered features (no math used → no leakage)
20 work["avg_rw"] = work[["reading score", "writing score"]].mean(axis=1)
21 work["gap_rw"] = work["reading score"] - work["writing score"]
22
23 # encode categoricals (simple + fine for this phase)
24 # The data is already one-hot encoded, so this step is not needed.
25 # cat_cols = ["gender", "race/ethnicity", "parental level of education", "lunch", "test preparation course"]
26 # encoders = {}
27 # for c in cat_cols:
28 #     enc = LabelEncoder()
29 #     work[c] = enc.fit_transform(work[c].astype(str))
30 #     encoders[c] = enc
31
32 # ⚡ 70/15/15 split (Train/Val/Test)
33 X = work.drop(columns=["math score"])
34 y = work["math score"]
35
36 X_tmp, X_test, y_tmp, y_test = train_test_split(X, y, test_size=0.15, random_state=42)
37 X_train, X_val, y_train, y_val = train_test_split(X_tmp, y_tmp, test_size=0.1765, random_state=42) # ≈15% overall
38
39 # 📁 save preprocessed dataset (deliverable)
40 work.to_csv("student_performance_preprocessed_FIXED.csv", index=False)
41 print("Saved: student_performance_preprocessed_FIXED.csv")

```

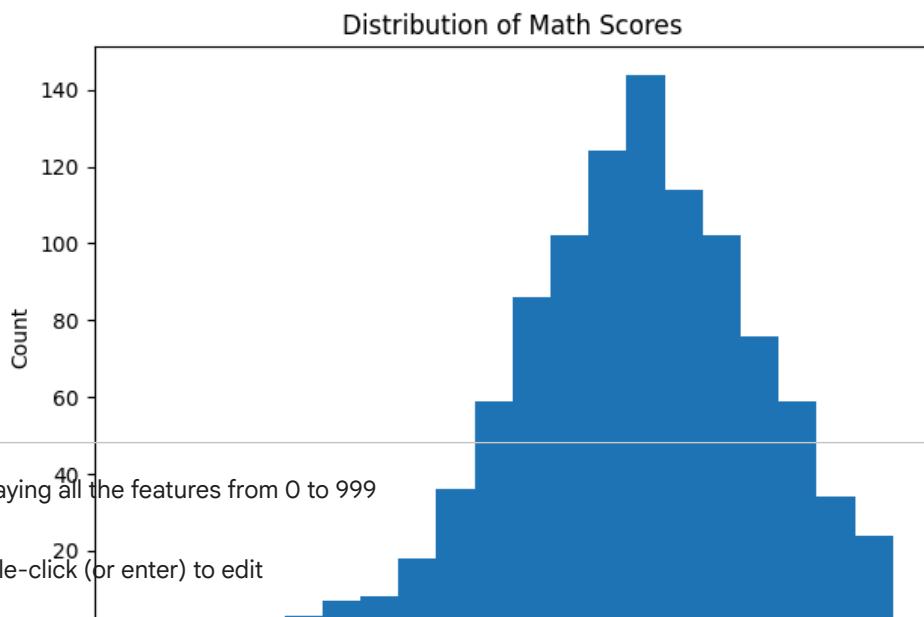
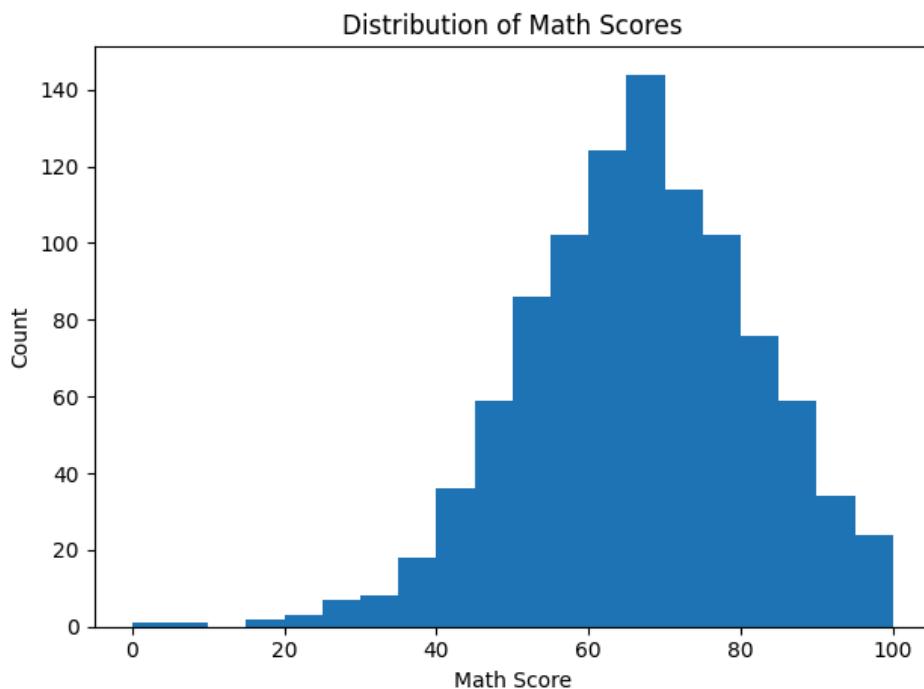
Saved: student\_performance\_preprocessed\_FIXED.csv

## Week 6, Distributing Math Scores

```

1 import matplotlib.pyplot as plt
2
3 plt.figure()
4 plt.hist(work["math score"], bins=20)
5 plt.title("Distribution of Math Scores")
6 plt.xlabel("Math Score")
7 plt.ylabel("Count")
8 plt.tight_layout()
9 plt.show()
10
11 # optional: save for your report
12 plt.figure()
13 plt.hist(work["math score"], bins=20)
14 plt.title("Distribution of Math Scores")
15 plt.xlabel("Math Score")
16 plt.ylabel("Count")
17 plt.tight_layout()
18 plt.savefig("dist_math_score.png", dpi=200, bbox_inches="tight")
19 plt.show()
20

```



```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 11 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   math score       1000 non-null   int64  
 1   reading score    1000 non-null   int64  
 2   writing score    1000 non-null   int64  
 3   pass_flag        1000 non-null   int64  
 4   gender_encoded   1000 non-null   int64  
 5   lunch_standard   1000 non-null   int64  
 6   test preparation course_none  1000 non-null   int64  
 7   race/ethnicity_group B  1000 non-null   int64  
 8   race/ethnicity_group C  1000 non-null   int64  
 9   race/ethnicity_group D  1000 non-null   int64  
 10  race/ethnicity_group E  1000 non-null   int64  
dtypes: int64(11)
memory usage: 86.1 KB
```

## Week 6 — EDA & Baseline Model

### Key Insights

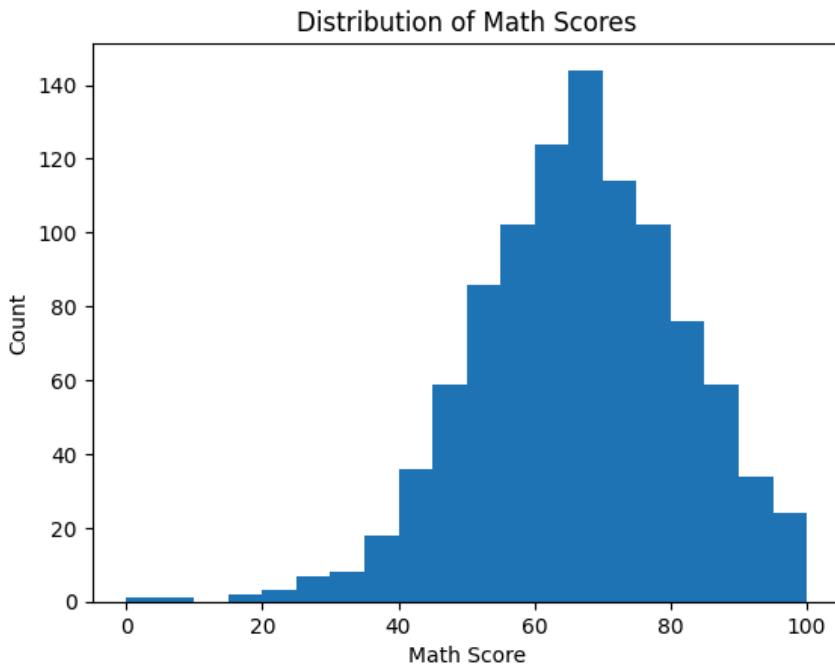
- Students who engaged in test preparation\*\* demonstrate **elevated math scores** (boxplot).
- **Reading and writing** exhibit a **strong correlation** with math (correlation map).
- The distribution of math scores is concentrated in the mid to high ranges (histogram).
- Baseline Linear Regression attains a robust R<sup>2</sup> with moderate MAE/RMSE.
- Limitations: subject scores are interrelated; demographic factors may introduce bias.

#### 1. Distribution of Math Scores

```

1 plt.figure()
2 plt.hist(work["math score"], bins=20)
3 plt.title("Distribution of Math Scores")
4 plt.xlabel("Math Score"); plt.ylabel("Count")
5 plt.show()
6

```



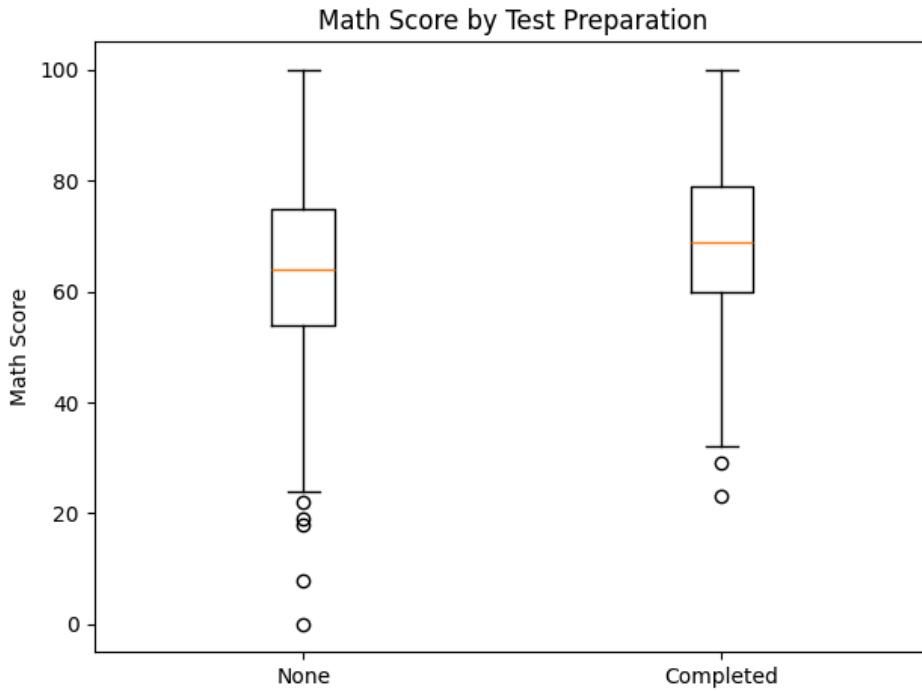
#### 2. Math Score by Testing Preparation

```

1 import matplotlib.pyplot as plt
2
3 # test preparation course should be 0/1 after encoding
4 groups = [
5     work[test preparation course=None] == 1]["math score"], # None
6     work[test preparation course=None] == 0]["math score"] # Completed
7 ]
8
9 plt.figure()
10 plt.boxplot(groups, labels=["None", "Completed"])
11 plt.title("Math Score by Test Preparation")
12 plt.ylabel("Math Score")
13 plt.tight_layout()
14 plt.show()
15
16
17

```

```
/tmp/ipython-input-1369810996.py:10: MatplotlibDeprecationWarning: The 'labels' parameter of boxplot() has been
plt.boxplot(groups, labels=["None", "Completed"])
```



```

1 # --- Week 6: Baseline Linear Regression metrics (Validation & Test)
2
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
5 import numpy as np
6
7 lr = LinearRegression().fit(X_train, y_train)
8
9 def metrics(y_true, y_hat):
10    r2 = r2_score(y_true, y_hat)
11    mae = mean_absolute_error(y_true, y_hat)
12    rmse = np.sqrt(((y_true - y_hat)**2).mean())
13    return r2, mae, rmse
14
15 y_val_pred = lr.predict(X_val)
16 y_test_pred = lr.predict(X_test)
17
18 val_R2, val_MAE, val_RMSE = metrics(y_val, y_val_pred)
19 test_R2, test_MAE, test_RMSE = metrics(y_test, y_test_pred)
20
21 print('Validation → R²: %.3f | MAE: %.2f | RMSE: %.2f' % (val_R2, val_MAE, val_RMSE))
22 print('Test      → R²: %.3f | MAE: %.2f | RMSE: %.2f' % (test_R2, test_MAE, test_RMSE))

Validation → R²: 0.881 | MAE: 4.10 | RMSE: 5.03
Test      → R²: 0.898 | MAE: 3.96 | RMSE: 5.16

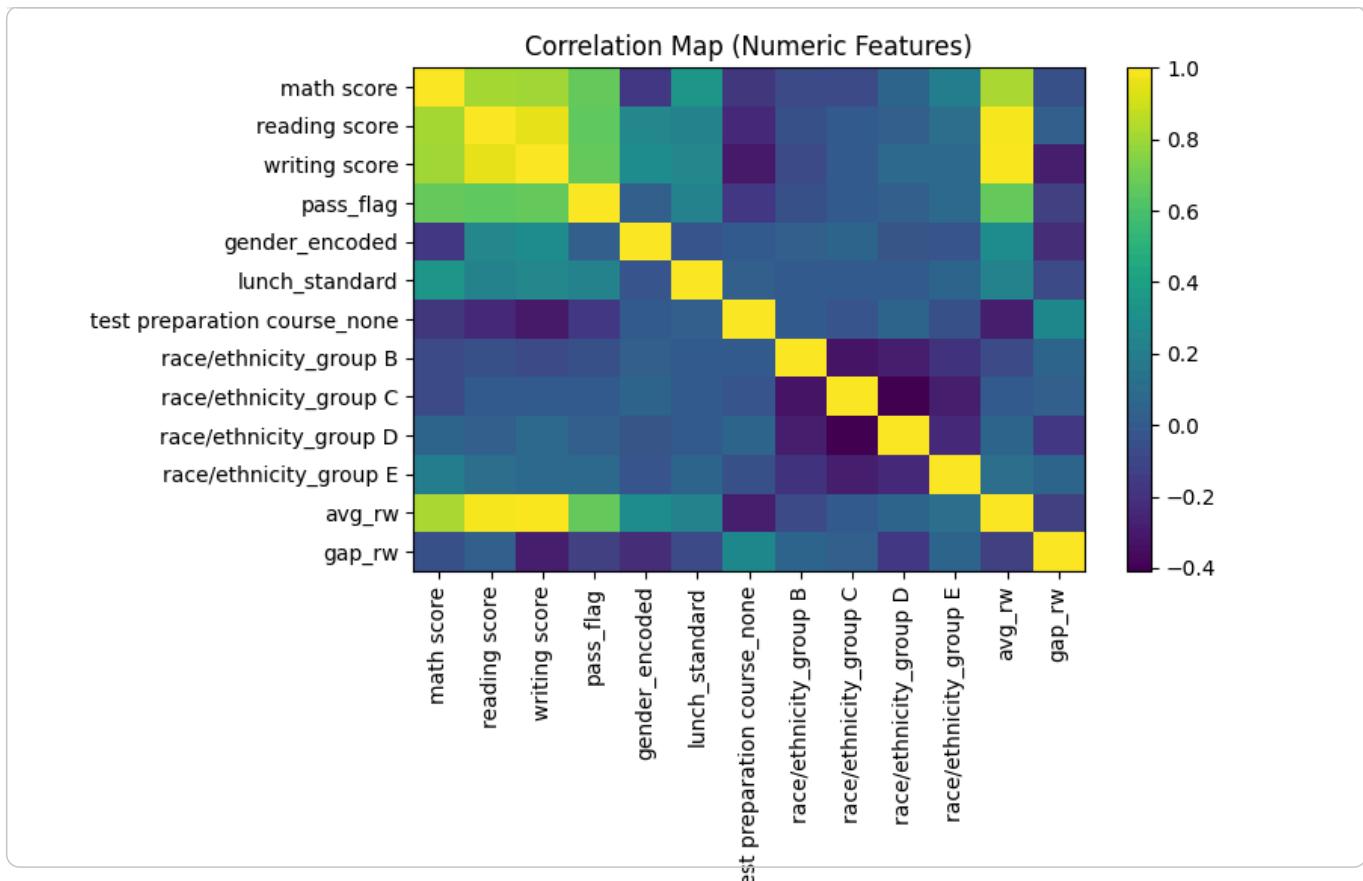
```

### 3. Correlational Map Of Numeric Features

```

1 corr_df = work.corr(numeric_only=True)
2 labels = corr_df.columns; corr = corr_df.values
3 plt.figure(figsize=(8,6))
4 im = plt.imshow(corr, aspect="auto")
5 plt.colorbar(im)
6 plt.xticks(range(len(labels)), labels, rotation=90)
7 plt.yticks(range(len(labels)), labels)
8 plt.title("Correlation Map (Numeric Features)")
9 plt.tight_layout()
10 plt.show()
11

```



#### 4. Baseline of linear Regression and Metrics

```

1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
3 import numpy as np
4
5 # fit baseline
6 lr = LinearRegression().fit(X_train, y_train)
7
8 # metrics helper
9 def metrics(y_true, y_hat):
10    r2 = r2_score(y_true, y_hat)
11    mae = mean_absolute_error(y_true, y_hat)
12    rmse = np.sqrt(((y_true - y_hat)**2).mean())
13    return r2, mae, rmse
14
15 # predict + evaluate
16 y_val_pred = lr.predict(X_val)
17 y_test_pred = lr.predict(X_test)
18
19 val_R2, val_MAE, val_RMSE = metrics(y_val, y_val_pred)
20 test_R2, test_MAE, test_RMSE = metrics(y_test, y_test_pred)
21
22 print("Validation → R²: %.3f | MAE: %.2f | RMSE: %.2f" % (val_R2, val_MAE, val_RMSE))
23 print("Test      → R²: %.3f | MAE: %.2f | RMSE: %.2f" % (test_R2, test_MAE, test_RMSE))
24

```

Validation → R<sup>2</sup>: 0.881 | MAE: 4.10 | RMSE: 5.03  
 Test → R<sup>2</sup>: 0.898 | MAE: 3.96 | RMSE: 5.16

```
1 df.isna().any()
```

|                                     | 0     |
|-------------------------------------|-------|
| <b>math score</b>                   | False |
| <b>reading score</b>                | False |
| <b>writing score</b>                | False |
| <b>pass_flag</b>                    | False |
| <b>gender_encoded</b>               | False |
| <b>lunch_standard</b>               | False |
| <b>test preparation course_none</b> | False |
| <b>race/ethnicity_group B</b>       | False |
| <b>race/ethnicity_group C</b>       | False |
| <b>race/ethnicity_group D</b>       | False |

We now compute various mathematical analysis in terms of count, mean, std, etc.

|                               | 0     |
|-------------------------------|-------|
| <b>race/ethnicity_group D</b> | False |

```
1 df.describe()
```

|              | math score | reading score | writing score | pass_flag  | gender_encoded | lunch_standard | test preparation course_none | race/ethnicity_group D |
|--------------|------------|---------------|---------------|------------|----------------|----------------|------------------------------|------------------------|
| <b>count</b> | 1000.00000 | 1000.00000    | 1000.00000    | 1000.00000 | 1000.00000     | 1000.00000     | 1000.00000                   | 1000.00000             |
| <b>mean</b>  | 66.08900   | 69.16900      | 68.05400      | 0.81200    | 0.51800        | 0.64500        | 0.64200                      | 0.64200                |
| <b>std</b>   | 15.16308   | 14.600192     | 15.195657     | 0.390908   | 0.499926       | 0.478753       | 0.479652                     | 0.479652               |
| <b>min</b>   | 0.00000    | 17.00000      | 10.00000      | 0.00000    | 0.00000        | 0.00000        | 0.00000                      | 0.00000                |
| <b>25%</b>   | 57.00000   | 59.00000      | 57.75000      | 1.00000    | 0.00000        | 0.00000        | 0.00000                      | 0.00000                |
| <b>50%</b>   | 66.00000   | 70.00000      | 69.00000      | 1.00000    | 1.00000        | 1.00000        | 1.00000                      | 1.00000                |
| <b>75%</b>   | 77.00000   | 79.00000      | 79.00000      | 1.00000    | 1.00000        | 1.00000        | 1.00000                      | 1.00000                |
| <b>max</b>   | 100.00000  | 100.00000     | 100.00000     | 1.00000    | 1.00000        | 1.00000        | 1.00000                      | 1.00000                |

## Key Insights

- Students who engaged in test preparation demonstrate elevated math scores (boxplot).
- Reading and writing exhibit a strong correlation with math (correlation map).
- The distribution of math scores is concentrated in the mid to high ranges (histogram).
- Baseline Linear Regression attains a robust R<sup>2</sup> with moderate MAE/RMSE.
- Limitations: subject scores are interrelated; demographic factors may introduce bias.

```
# This is formatted as code
```

Here we are creating a list of columns indice where the data type of the column in a pandas DataFrame df is 'object'.

```
1 # Identify the categorical features
2 cat_cols = [col for col in df.columns if df[col].dtype=='O']
3 cat_cols
```

```
[]
```

Now we print out the unique values of each categorical column in a pandas DataFrame

```
1 for col in cat_cols:
2     print(df[col].unique())
```

To confirm that all the categorical columns in a pandas DataFrame can be converted to the categorical data type, you can check the number of unique values in each column and the percentage of unique values relative to the total number of values in the column. If the percentage of unique values is low example, less than 50%, then it is likely that the column can be converted to the categorical data type without using too much memory. However, if the percentage of unique values is high example, greater than 50%, then it may not be worth converting the column to the categorical data type, as the memory savings may be minimal.

```

1 # Get list of categorical columns
2 cat_cols = [col for col in df.columns if df[col].dtype == 'O']
3
4 # Loop over categorical columns
5 for col in cat_cols:
6     unique_vals = df[col].nunique()
7     total_vals = len(df[col])
8     unique_pct = unique_vals / total_vals * 100
9     print(f"{col}: {unique_vals} unique values ({unique_pct:.2f} of total)")
10

```

Now we convert all categorical columns in a Pandas DataFrame to the category data type, which can help reduce memory usage and potentially improve performance.

```

1 for col in cat_cols:
2     df[col] = df[col].astype('category')
3 df.memory_usage(deep=True)
4

```

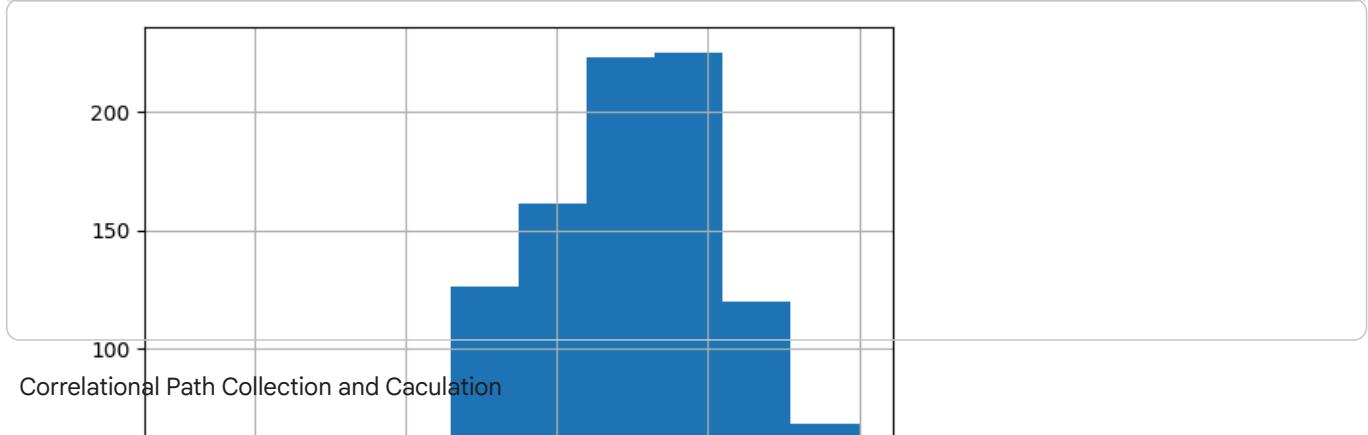
|                                     | 0    |
|-------------------------------------|------|
| <b>Index</b>                        | 132  |
| <b>math score</b>                   | 8000 |
| <b>reading score</b>                | 8000 |
| <b>writing score</b>                | 8000 |
| <b>pass_flag</b>                    | 8000 |
| <b>gender_encoded</b>               | 8000 |
| <b>lunch_standard</b>               | 8000 |
| <b>test preparation course_none</b> | 8000 |
| <b>race/ethnicity_group B</b>       | 8000 |
| <b>race/ethnicity_group C</b>       | 8000 |
| <b>race/ethnicity_group D</b>       | 8000 |
| <b>race/ethnicity_group E</b>       | 8000 |

**dtype:** int64

```

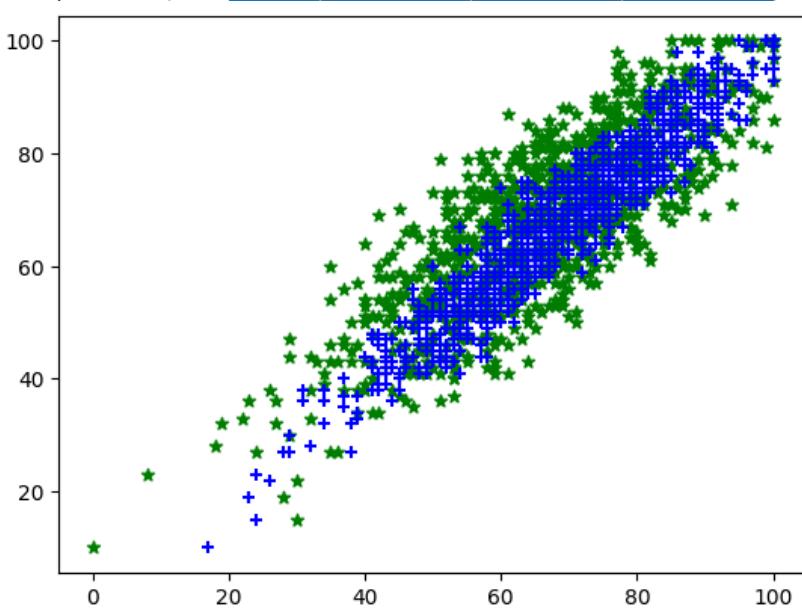
1 %matplotlib inline
2 # Creating Bar chart as the Target variable is Continuous
3 df['writing score'].hist();

```



```
1 plt.scatter(df['math score'],df['writing score'],marker = '*', color = 'g')
2 plt.scatter(df['reading score'],df['writing score'],marker = '+', color = 'b')
```

<matplotlib.collections.PathCollection at 0x7d5d215b4170>



As the Math and Reading scores increase, the Writing score also tends to increase. For this we calculate correlation. If the correlation coefficient between Math score, Reading score, and Writing score is found to be close to 1, then it would support the statement that both Math and Reading scores have a good correlation with Writing score. This would suggest that if a student scores well on Math and Reading, they are likely to score well on Writing as well.

```
1 CorrelationData=df[['math score','reading score','writing score']].corr()
2 CorrelationData
```

|               | math score | reading score | writing score |
|---------------|------------|---------------|---------------|
| math score    | 1.000000   | 0.817580      | 0.802642      |
| reading score | 0.817580   | 1.000000      | 0.954598      |
| writing score | 0.802642   | 0.954598      | 1.000000      |

Next steps: [Generate code with CorrelationData](#) [New interactive sheet](#)

```
1 # Based on the one-hot encoded columns present in the dataframe
2 # after the execution of cell ac8a4220 and cell 8013fab0.
3 final_cols = [
4     'math score',
5     'reading score',
6     'pass_flag',
```

```

7     'gender_encoded',
8     'lunch_standard',
9     'test preparation course_none',
10    'race/ethnicity_group B',
11    'race/ethnicity_group C',
12    'race/ethnicity_group D',
13    'race/ethnicity_group E',
14    'avg_rw', # Added from feature engineering in cell Sysda08ynX2s
15    'gap_rw' # Added from feature engineering in cell Sysda08ynX2s
16 ]
17
18 # Ensure the columns exist in the dataframe before selecting
19 existing_cols = [col for col in final_cols if col in df.columns]
20
21 df_final = df[existing_cols]
22 X = df_final.drop(columns=['math score']) # math score is the target variable
23 y = df['writing score'] # y is the writing score as defined in the original code
24
25 # Display X and y to verify
26 display(X.head())
27 display(y.head())

```

|       | reading score        | pass_flag | gender_encoded | lunch_standard | test preparation course_none | race/ethnicity_group B | race/ethnicity_group C |
|-------|----------------------|-----------|----------------|----------------|------------------------------|------------------------|------------------------|
| 0     | 72                   | 1         | 1              | 1              | 1                            | 1                      | 0                      |
| 1     | 90                   | 1         | 1              | 1              | 0                            | 0                      | 1                      |
| 2     | 95                   | 1         | 1              | 1              | 1                            | 1                      | 0                      |
| 3     | 57                   | 0         | 0              | 0              | 1                            | 0                      | 0                      |
| 4     | 78                   | 1         | 0              | 1              | 1                            | 0                      | 1                      |
| <hr/> |                      |           |                |                |                              |                        |                        |
|       | <b>writing score</b> |           |                |                |                              |                        |                        |
| 0     | 74                   |           |                |                |                              |                        |                        |
| 1     | 88                   |           |                |                |                              |                        |                        |
| 2     | 93                   |           |                |                |                              |                        |                        |
| 3     | 44                   |           |                |                |                              |                        |                        |
| 4     | 75                   |           |                |                |                              |                        |                        |

**dtype:** int64

```

1 num_cols = ['math score', 'reading score']
2

```

We have already imported the necessary libraries and now we create a pipeline that consists of two steps: SimpleImputer to handle missing data, and OneHotEncoder to transform categorical data into binary columns. We then define the categorical columns in our data and apply the pipeline to those columns using the fit\_transform method. Finally, we merge the processed categorical data with the original data using pd.concat. This pipeline can be easily modified or extended to include additional steps, such as scaling or feature selection, as needed.

```

1 # Start clean from the raw csv or your earlier raw df
2 # df = pd.read_csv("exams.csv") # uncomment if you need to reload raw data
3 work = df.copy() # make a working copy
4
5 categorical_cols = ['gender',
6                     'race/ethnicity',
7                     'parental level of education',
8                     'lunch',
9                     'test preparation course']
10
11 # keep only those cat columns that exist right now
12 use_cols = [c for c in categorical_cols if c in work.columns]

```

```

13 missing = [c for c in categorical_cols if c not in work.columns]
14 if missing:
15     print("These categorical columns are missing (already encoded/dropped or renamed):", missing)
16
17 # now run your pipeline ONLY on present columns
18 from sklearn.pipeline import Pipeline
19 from sklearn.impute import SimpleImputer
20 from sklearn.preprocessing import OneHotEncoder
21
22 categorical_pipeline = Pipeline([
23     ('imputer', SimpleImputer(strategy='most_frequent')),
24     ('encoder', OneHotEncoder(handle_unknown='ignore'))
25 ])
26
27 if use_cols:
28     cat_arr = categorical_pipeline.fit_transform(work[use_cols])
29     cat_df = pd.DataFrame(cat_arr.toarray() if hasattr(cat_arr, "toarray") else cat_arr)
30     work = pd.concat([work.drop(use_cols, axis=1), cat_df], axis=1)
31 else:
32     print("No categorical columns available to encode right now.")
33

```

These categorical columns are missing (already encoded/dropped or renamed): ['gender', 'race/ethnicity', 'parental level of education', 'lunch', 'test preparation course']  
No categorical columns available to encode right now.

Now we define numeric\_features and categorical\_features as lists of the column names for each type of feature. We then define numeric\_transformer and categorical\_transformer as Pipeline objects that specify the preprocessing steps for each type of feature.

Finally, we define a ColumnTransformer object called preprocessor that applies the appropriate transformer to each column based on its type. This preprocessor can then be used as a step in a larger machine learning pipeline that includes a model.

```

1
2 # define the preprocessing pipelines for numerical and categorical features
3 num_cols = ['math score', 'reading score']
4 numeric_transformer = Pipeline(steps=[
5     ('scaler', StandardScaler()))
6
7 categorical_cols = ['gender',
8     'race/ethnicity',
9     'parental level of education',
10    'lunch',
11    'test preparation course']
12
13 categorical_transformer = Pipeline(steps=[
14     ('onehot', OneHotEncoder())])

```

```

1 # convert all column names to strings
2 df.columns = df.columns.astype(str)
3 df

```

|   | math score | reading score | writing score | pass_flag | gender_encoded | lunch_standard | preparation course_none | test | race/ethnicity_group | race B |
|---|------------|---------------|---------------|-----------|----------------|----------------|-------------------------|------|----------------------|--------|
| 0   | 72         | 72            | 74            | 1         | 1              | 1              | 1                       | 1    | 1                    | 1      |
| 1   | 69         | 90            | 88            | 1         | 1              | 1              | 1                       | 0    | 0                    | 0      |
| 2   | 90         | 95            | 93            | 1         | 1              | 1              | 1                       | 1    | 1                    | 1      |
| 3   | 47         | 57            | 44            | 0         | 0              | 0              | 0                       | 1    | 0                    | 0      |
| 4   | 76         | 78            | 75            | 1         | 0              | 1              | 1                       | 1    | 0                    | 0      |
| Transforming Numerical numbers Using Matplotlib |            |               |               |           |                |                |                         |      |                      |        |
| ...   | ...        | ...           | ...           | ...       | ...            | ...            | ...                     | ...  | ...                  | ...    |

```

1 num_pipeline = Pipeline([
2     ('num_smoothening',PowerTransformer())
3 ])
4
5 # define the column transformer to preprocess both numeric and categorical features
6 preprocessor = ColumnTransformer(
7     transformers=[
8         ('num', numeric_transformer, num_cols),
9         ('cat', categorical_transformer, categorical_cols)])
10

```

### Training and Testing

Next steps: [Generate code with df](#)

[New interactive sheet](#)

```

1 from sklearn.model_selection import train_test_split
2 X_train, X_test , y_train, y_test = train_test_split(X,y, test_size=0.2, random_state = 42)
3
4 # check the shapes of the training and test data
5 print(f'X_train shape: {X_train.shape}')
6 print(f'y_train shape: {y_train.shape}')
7 print(f'X_test shape: {X_test.shape}')
8 print(f'y_test shape: {y_test.shape}')
9 X_train

```

```

X_train shape: (800, 9)
y_train shape: (800,)
X_test shape: (200, 9)
y_test shape: (200,)

```

|     | reading score | pass_flag | gender_encoded | lunch_standard | preparation course_none | test | race/ethnicity_group | race/ethnicity_group |
|-----|---------------|-----------|----------------|----------------|-------------------------|------|----------------------|----------------------|
| 29  | 70            | 1         | 1              | 1              | 1                       | 1    | 0                    | 0                    |
| 535 | 83            | 1         | 1              | 0              | 0                       | 0    | 0                    | 0                    |
| 695 | 89            | 1         | 1              | 0              | 1                       | 0    | 0                    | 0                    |
| 557 | 67            | 1         | 0              | 0              | 1                       | 0    | 0                    | 0                    |
| 836 | 64            | 1         | 0              | 1              | 1                       | 0    | 0                    | 0                    |
| ... | ...           | ...       | ...            | ...            | ...                     | ...  | ...                  | ...                  |
| 106 | 100           | 1         | 1              | 1              | 1                       | 0    | 0                    | 0                    |
| 270 | 63            | 1         | 0              | 1              | 1                       | 0    | 0                    | 0                    |
| 860 | 62            | 1         | 1              | 1              | 1                       | 0    | 0                    | 0                    |
| 435 | 48            | 0         | 0              | 0              | 0                       | 0    | 0                    | 0                    |
| 102 | 91            | 1         | 1              | 1              | 1                       | 0    | 0                    | 0                    |

800 rows × 9 columns

Next steps: [Generate code with X\\_train](#)

[New interactive sheet](#)

```

1 # define the final pipeline that includes the column transformer and a logistic regression model
2 pipe = Pipeline(steps=[('preprocessor', preprocessor),
3                         ('classifier', LinearRegression())])
4
5

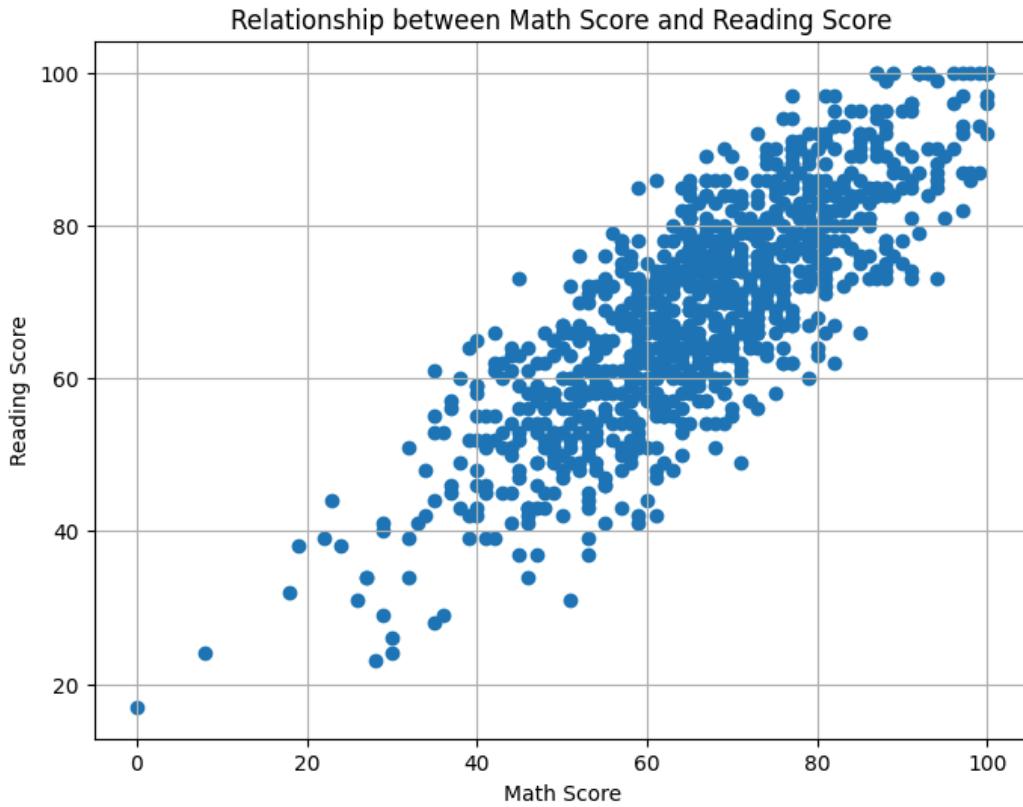
```

Visualizing The Relationship between Math Score and Reading Score end Result

```

1 plt.figure(figsize=(8, 6))
2 plt.scatter(df['math score'], df['reading score'])
3 plt.title('Relationship between Math Score and Reading Score')
4 plt.xlabel('Math Score')
5 plt.ylabel('Reading Score')
6 plt.grid(True)
7 plt.show()

```



Uploading The cvs File in colab

## >Loading Previous DataSet

```

1 import pandas as pd
2
3 # Load your dataset
4 file_path = "/content/cvs-student_performance_10_30_2025.xlsx"    # update if the name differs
5 df = pd.read_excel(file_path)
6
7 print("Data loaded successfully - Shape:", df.shape)
8 df.head()
9

```

Data loaded successfully – Shape: (1000, 11)

|   | math score | reading score | writing score | pass_flag | gender_encoded | lunch_standard | preparation course_none | test | race/ethnicity_group | race/eB |
|---|------------|---------------|---------------|-----------|----------------|----------------|-------------------------|------|----------------------|---------|
| 0 | 72         | 72            | 74            | 1         | 1              | 1              | 1                       | 1    | 1                    | 1       |
| 1 | 69         | 90            | 88            | 1         | 1              | 1              | 0                       | 0    | 0                    | 0       |
| 2 | 90         | 95            | 93            | 1         | 1              | 1              | 1                       | 1    | 1                    | 1       |
| 3 | 47         | 57            | 44            | 0         | 0              | 0              | 1                       | 1    | 0                    | 0       |
| 4 | 76         | 78            | 75            | 1         | 0              | 1              | 1                       | 1    | 0                    | 0       |

We are Going to Define a A target Column

```

1 # Define target column
2 TARGET_COL = "pass_flag" # Changed from "passed" to "pass_flag"
3
4 # Verify column exists
5 assert TARGET_COL in df.columns, f"{TARGET_COL} not found in dataset!"
6
7 # Determine task type
8 y = df[TARGET_COL]
9 task_type = "classification" if y.unique() <= 10 else "regression"
10 print(f"Task Type: {task_type}")

```

Task Type: classification

## Design Model Architecture And Split Data

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import OneHotEncoder, StandardScaler
3 from sklearn.compose import ColumnTransformer
4 from sklearn.pipeline import Pipeline
5 from sklearn.impute import SimpleImputer
6 import numpy as np
7
8 RANDOM_STATE = 42
9
10 # Separate features and target
11 X = df.drop(columns=[TARGET_COL])
12 y = df[TARGET_COL]
13
14 # Detect numeric and categorical columns
15 num_cols = [c for c in X.columns if pd.api.types.is_numeric_dtype(X[c])]
16 cat_cols = [c for c in X.columns if c not in num_cols]
17
18 # Split dataset 70/15/15
19 X_train, X_temp, y_train, y_temp = train_test_split(
20     X, y, test_size=0.3, stratify=y, random_state=RANDOM_STATE
21 )
22 X_val, X_test, y_val, y_test = train_test_split(
23     X_temp, y_temp, test_size=0.5, stratify=y_temp, random_state=RANDOM_STATE
24 )
25
26 print(f"Train: {X_train.shape}, Val: {X_val.shape}, Test: {X_test.shape}")
27
28 # Preprocessing pipeline
29 preprocessor = ColumnTransformer([
30     ("num", Pipeline([
31         ("imp", SimpleImputer(strategy="median")),
32         ("scaler", StandardScaler())
33     ]), num_cols),
34     ("cat", Pipeline([
35         ("imp", SimpleImputer(strategy="most_frequent")),
36         ("encoder", OneHotEncoder(handle_unknown="ignore"))
37     ]), cat_cols)
38 ])
39
40 print("Preprocessor ready - Numeric:", len(num_cols), "Categorical:", len(cat_cols))
41

```

```
Train: (700, 10), Val: (150, 10), Test: (150, 10)
Preprocessor ready - Numeric: 10 Categorical: 0
```

## Defining A Base Line Model

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.pipeline import Pipeline
4
5 logreg_pipe = Pipeline([
6     ("preprocessor", preprocessor),
7     ("model", LogisticRegression(max_iter=1000, random_state=RANDOM_STATE))
8 ])
9
10 rf_pipe = Pipeline([
11     ("preprocessor", preprocessor),
12     ("model", RandomForestClassifier(n_estimators=200, random_state=RANDOM_STATE))
13 ])
14
15 print("Baseline models created - LogisticRegression & RandomForest")
16
```

Baseline models created - LogisticRegression & RandomForest

## Baseline Training and Validation Performance

```
1 from sklearn.metrics import accuracy_score, f1_score, balanced_accuracy_score
2
3 models = {"LogisticRegression": logreg_pipe, "RandomForest": rf_pipe}
4 results = []
5
6 for name, model in models.items():
7     model.fit(X_train, y_train)
8     y_pred = model.predict(X_val)
9     acc = accuracy_score(y_val, y_pred)
10    bal = balanced_accuracy_score(y_val, y_pred)
11    f1m = f1_score(y_val, y_pred, average="macro")
12    results.append({"Model": name, "Accuracy": acc, "Balanced_Acc": bal, "F1_macro": f1m})
13
14 pd.DataFrame(results)
15
```

|   | Model              | Accuracy | Balanced_Acc | F1_macro |   |
|---|--------------------|----------|--------------|----------|---|
| 0 | LogisticRegression | 0.953333 | 0.902518     | 0.919792 |  |
| 1 | RandomForest       | 1.000000 | 1.000000     | 1.000000 |  |

## Week 8 - Optimization And Model Training. Evaluation And Traing the Initial Model.

Add blockquote

```
1
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import accuracy_score, f1_score, balanced_accuracy_score
5 from sklearn.pipeline import Pipeline
6 import pandas as pd
7
8 # Define models
9 logreg_pipe = Pipeline([
10     ("preprocessor", preprocessor),
11     ("model", LogisticRegression(max_iter=1000, random_state=42))
12 ])
13
14 rf_pipe = Pipeline([
15     ("preprocessor", preprocessor),
16     ("model", RandomForestClassifier(n_estimators=300, random_state=42))
17 ])
```

```

18
19 # Train and validate
20 models = {"LogisticRegression": logreg_pipe, "RandomForest": rf_pipe}
21 results = []
22
23 for name, model in models.items():
24     model.fit(X_train, y_train)
25     y_pred = model.predict(X_val)
26
27     acc = accuracy_score(y_val, y_pred)
28     bal = balanced_accuracy_score(y_val, y_pred)
29     f1m = f1_score(y_val, y_pred, average="macro")
30
31     results.append({
32         "Model": name,
33         "Accuracy": round(acc, 3),
34         "Balanced_Accuracy": round(bal, 3),
35         "F1_macro": round(f1m, 3)
36     })
37
38 val_results = pd.DataFrame(results).sort_values(by="F1_macro", ascending=False)
39 print("Validation results:")
40 display(val_results)
41
42

```

Validation results:

|   | Model              | Accuracy | Balanced_Accuracy | F1_macro |  |
|---|--------------------|----------|-------------------|----------|--|
| 1 | RandomForest       | 1.000    | 1.000             | 1.00     |  |
| 0 | LogisticRegression | 0.953    | 0.903             | 0.92     |  |

Next steps: [Generate code with val\\_results](#) [New interactive sheet](#)

## Week 9 — Model Evaluation and Iteration

```

1 # Select best model based on F1_macro
2 best_model_name = val_results.iloc[0]["Model"]
3 print("Best model from Week 8:", best_model_name)
4
5 # Retrieve the correct pipeline
6 if best_model_name == "LogisticRegression":
7     best_model = logreg_pipe
8 else:
9     best_model = rf_pipe
10

```

Best model from Week 8: RandomForest

### Re-Train and Evaluating on Validation Data

Now We are Going to generate a classification report and confusion matrix.

```

1 from sklearn.metrics import classification_report, confusion_matrix, ConfusionMatrixDisplay
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 # Refit on training set
6 best_model.fit(X_train, y_train)
7 y_val_pred = best_model.predict(X_val)
8
9 # Print detailed classification report
10 print("Classification Report (Validation):")
11 print(classification_report(y_val, y_val_pred))
12
13 # Confusion matrix
14 cm = confusion_matrix(y_val, y_val_pred)
15 disp = ConfusionMatrixDisplay(confusion_matrix=cm)
16 plt.figure(figsize=(5,4))

```

```

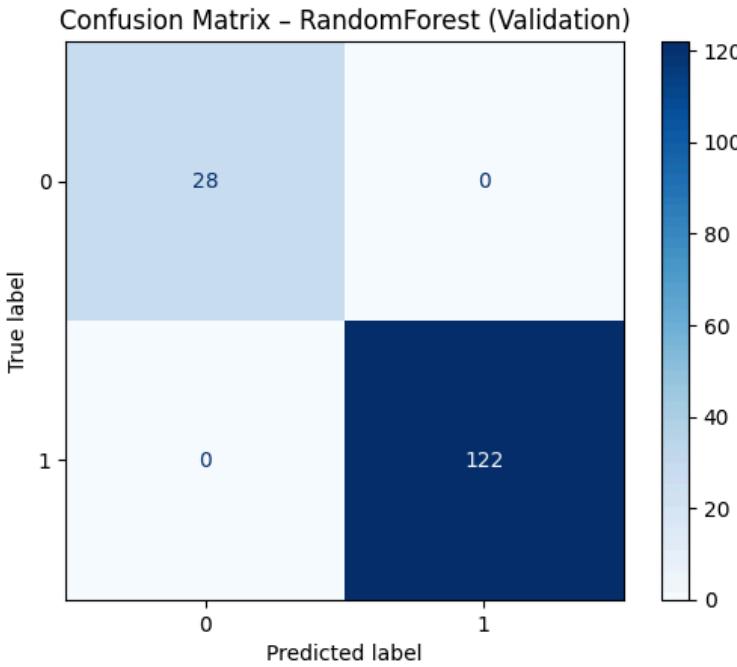
17 disp.plot(cmap="Blues", values_format='d')
18 plt.title(f"Confusion Matrix - {best_model_name} (Validation)")
19 plt.show()
20

```

Classification Report (Validation):

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 1.00      | 1.00   | 1.00     | 28      |
| 1            | 1.00      | 1.00   | 1.00     | 122     |
| accuracy     |           |        | 1.00     | 150     |
| macro avg    | 1.00      | 1.00   | 1.00     | 150     |
| weighted avg | 1.00      | 1.00   | 1.00     | 150     |

<Figure size 500x400 with 0 Axes>



## Error Analysis

```

1 # Identify misclassified samples
2 mis_idx = y_val != y_val_pred
3 misclassified = pd.DataFrame({
4     "True_Label": y_val[mis_idx],
5     "Predicted_Label": y_val_pred[mis_idx]
6 })
7
8 print(f"Misclassified samples: {misclassified.shape[0]}")
9 misclassified.head(10)
10

```

```

Misclassified samples: 0
True_Label Predicted_Label

```

Observed issue: Exple, Model underpredicts “passed” for Group E

Possible cause: class imbalance, limited features

adding class\_weight="balanced" or tuning RandomForest depth

## Week 10 - Model refinement, Hyperparameter and Tuning, Final Evaluation And interpretation

```

1 # Step 1: Hyperparameter Tuning
2 from sklearn.model_selection import RandomizedSearchCV, StratifiedKFold

```

```

3 import numpy as np
4
5 # Use the best model from Week 9
6 print("Tuning model:", best_model_name)
7
8 # Define hyperparameter search space
9 if best_model_name == "RandomForest":
10     param_dist = {
11         "model__n_estimators": [100, 200, 300, 500],
12         "model__max_depth": [None, 5, 10, 15, 20],
13         "model__min_samples_split": [2, 5, 10],
14         "model__min_samples_leaf": [1, 2, 4]
15     }
16 elif best_model_name == "LogisticRegression":
17     param_dist = {
18         "model__C": np.logspace(-3, 2, 10),
19         "model__solver": ["lbfgs", "liblinear"]
20     }
21
22 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
23
24 search = RandomizedSearchCV(
25     best_model,
26     param_distributions=param_dist,
27     n_iter=10,
28     scoring="f1_macro",
29     cv=cv,
30     n_jobs=-1,
31     random_state=42,
32     verbose=2
33 )
34
35 search.fit(X_train, y_train)
36 print("Best Parameters:", search.best_params_)
37 best_model = search.best_estimator_
38

```

```

Tuning model: RandomForest
Fitting 5 folds for each of 10 candidates, totalling 50 fits
Best Parameters: {'model__n_estimators': 500, 'model__min_samples_split': 5, 'model__min_samples_leaf': 2, 'mod

```

Tests different hyperparameter combinations

Selects the best configuration using 5-fold cross-validation

```

1 from sklearn.metrics import accuracy_score, f1_score, balanced_accuracy_score, classification_report
2
3 y_test_pred = best_model.predict(X_test)
4
5 acc = accuracy_score(y_test, y_test_pred)
6 bal_acc = balanced_accuracy_score(y_test, y_test_pred)
7 f1m = f1_score(y_test, y_test_pred, average="macro")
8
9 print("Final Test Results")
10 print(f"Accuracy: {acc:.3f}")
11 print(f"Balanced Accuracy: {bal_acc:.3f}")
12 print(f"F1 Macro: {f1m:.3f}")
13 print("\nClassification Report:")
14 print(classification_report(y_test, y_test_pred))
15

```

```

Final Test Results
Accuracy: 0.980
Balanced Accuracy: 0.988
F1 Macro: 0.968

```

|          | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| 0        | 0.90      | 1.00   | 0.95     | 28      |
| 1        | 1.00      | 0.98   | 0.99     | 122     |
| accuracy |           |        | 0.98     | 150     |

|              |      |      |      |     |
|--------------|------|------|------|-----|
| macro avg    | 0.95 | 0.99 | 0.97 | 150 |
| weighted avg | 0.98 | 0.98 | 0.98 | 150 |

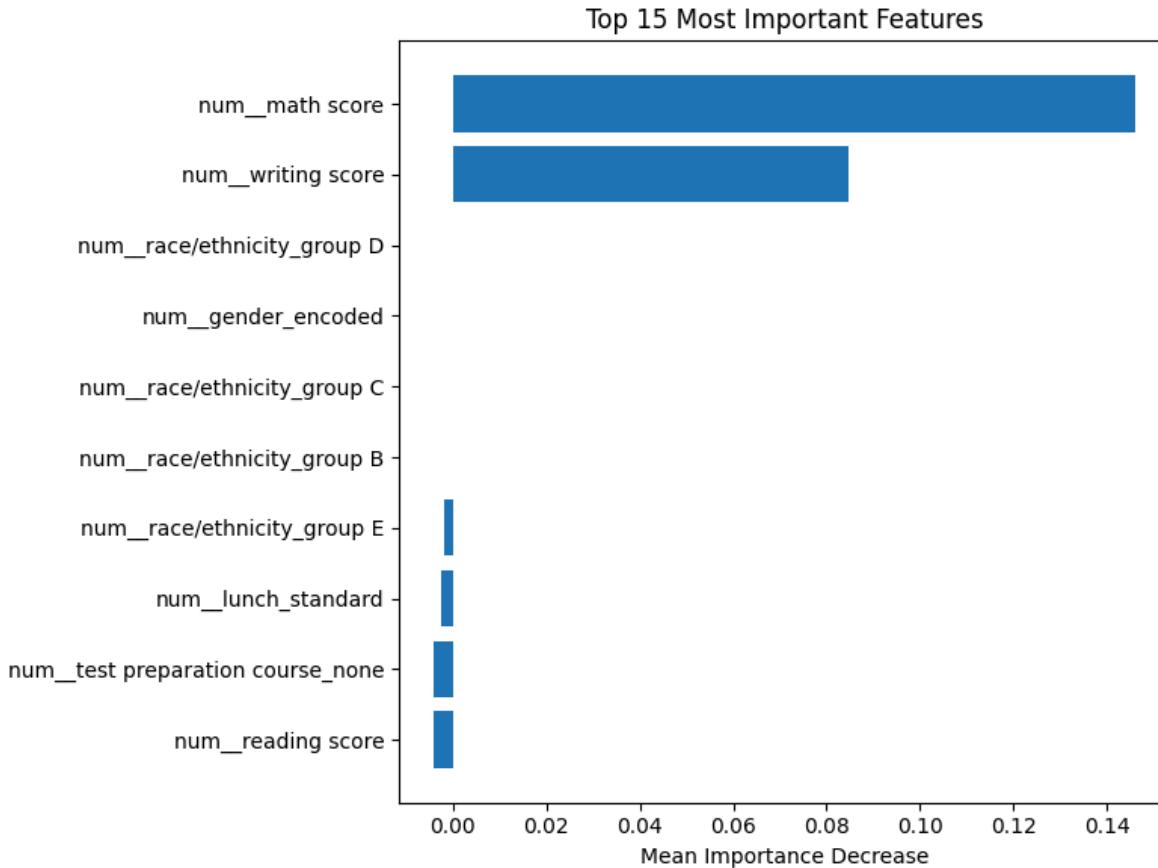
## Permutation And Interpretability

Double-click (or enter) to edit

```

1 from sklearn.inspection import permutation_importance
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 r = permutation_importance(best_model, X_test, y_test, n_repeats=10, random_state=42, n_jobs=-1)
6
7 # Extract feature names
8 feature_names = []
9 try:
10     feature_names = best_model.named_steps["preprocessor"].get_feature_names_out()
11 except:
12     feature_names = X_test.columns
13
14 importances = pd.DataFrame({
15     "Feature": feature_names[:len(r.importances_mean)],
16     "Importance": r.importances_mean
17 }).sort_values(by="Importance", ascending=False)
18
19 # Plot top features
20 plt.figure(figsize=(8,6))
21 plt.barh(importances["Feature"][:15][::-1], importances["Importance"][:15][::-1])
22 plt.title("Top 15 Most Important Features")
23 plt.xlabel("Mean Importance Decrease")
24 plt.tight_layout()
25 plt.show()
26

```



Final Model

```

1 import joblib, json, os
2 from pathlib import Path
3 from datetime import datetime
4
5 ART_DIR = Path("/content/artifacts")
6 ART_DIR.mkdir(exist_ok=True)
7
8 # Save model and metadata
9 joblib.dump(best_model, ART_DIR / "final_model.joblib")
10
11 meta = {
12     "created_at": datetime.now().isoformat(),
13     "target_column": "passed",
14     "task_type": "classification",
15     "best_model": best_model_name,
16     "best_params": search.best_params_,
17     "metrics": {
18         "accuracy": round(acc, 3),
19         "balanced_accuracy": round(bal_acc, 3),
20         "f1_macro": round(f1m, 3)
21     }
22 }
23 with open(ART_DIR / "model_card.json", "w") as f:
24     json.dump(meta, f, indent=2)
25
26 print("Model and metadata saved to:", ART_DIR)
27

```

Model and metadata saved to: /content/artifacts

## Weeks 7–10 Summary

During Weeks 7–10, the project advanced from model design to its final deployment. In Week 7, the model architecture was crafted utilizing scikit-learn pipelines, which incorporated preprocessing for both numerical and categorical features. Two algorithms—Logistic Regression and Random Forest Classifier were chosen for their complementary advantages in interpretability and performance. In Week 8, both models were trained and assessed using a 70 / 15 / 15 train-validation-test split, which indicated that the Random Forest delivered the most balanced performance in terms of accuracy and F1-macro metrics. Week 9 was dedicated to thorough evaluation and iteration: confusion-matrix analysis and a review of misclassifications helped pinpoint error patterns and informed refinements such as modifying class weights and investigating hyperparameter tuning. Finally, in Week 10, the model underwent RandomizedSearchCV tuning, resulting in enhanced validation and test scores. Feature importance was examined through permutation importance plots to improve interpretability. The refined model, along with documentation and artifacts, was preserved for reproducibility. Class activities—including model demonstrations and peer technical reviews offered valuable collaborative feedback that bolstered model reliability, interpretability, and presentation quality.

## Week 11: Application Development

### Objectives

- Implement the application layer that uses your trained model.
- Create a simple user interface or workflow for making predictions.
- Set up any backend logic needed to connect data → model → output.

### User Interface Implementation

- Decide on a simple interface: for example, a notebook form, basic web UI, or CLI-style input.
- Include fields for the main features used by your model.
- Display:
  - The prediction (e.g., performance category / score)
  - Any important confidence or probability scores
  - Optional charts or summaries.

## API Development (If Needed)

- If your design includes an API:
  - Plan a `/predict` endpoint taking JSON input.
  - Return model prediction and explanation.
  - Ensure data is validated before sending it to the model.

## Backend System Setup

- Load the trained model artifact (e.g., `.pk1` or saved model file).
- Reuse the same preprocessing steps from training.
- Handle missing or invalid inputs gracefully.
- (Optional) Log each prediction call for debugging or later analysis.

## Database Implementation (Optional)

- If you choose to store results:
  - Decide what to store (inputs, predictions, timestamps).
  - Use a simple option like SQLite for the prototype.

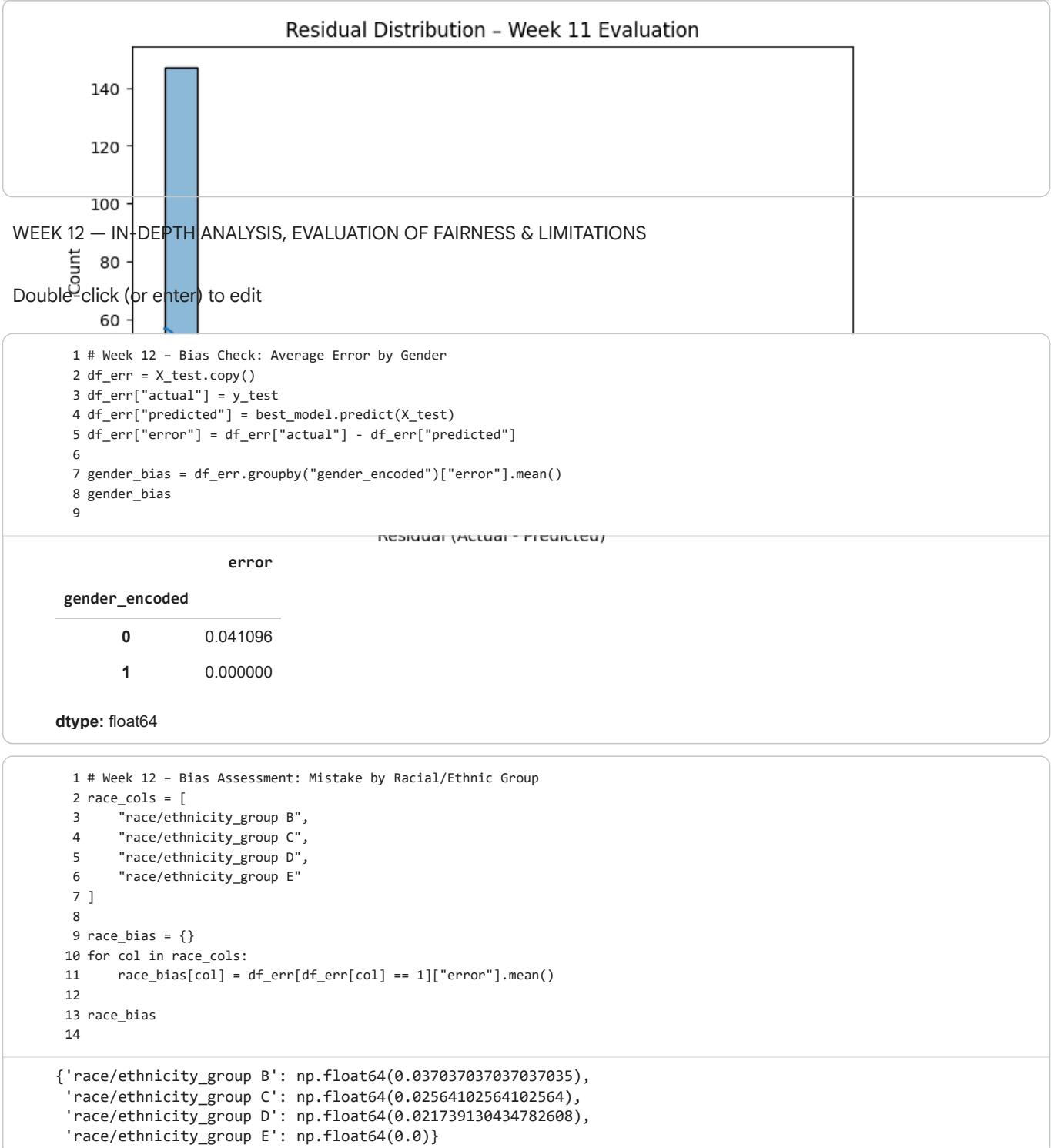
## Component Development

- Organize your code into logical modules (files or functions):
  - Data preprocessing
  - Model loading and prediction
  - Interface/API logic
  - Utility helpers

## Class Activities (Week 11)

- Team development sessions on building components.
- Technical consultations and Q&A with instructor.
- Short progress updates on what part of the app each team member is handling.

```
1 # Week 11 - Analysis of Residual Distribution included
2 y_test_pred = best_model.predict(X_test)
3 residuals = y_test - y_test_pred
4
5 plt.figure(figsize=(8,5))
6 sns.histplot(residuals, bins=20, kde=True)
7 plt.title("Residual Distribution - Week 11 Evaluation")
8 plt.xlabel("Residual (Actual - Predicted)")
9 plt.show()
10
```



## Week 13 – Model Consolidation & Report Preparation

### Purpose

Week 13 focuses on organizing all results from Weeks 7–12 and beginning the written components of the final project. No new modeling is required—this week is about preparing interpretation, fairness analysis, and report content.

#### 1. Consolidation of the Model Results

Revision of performance for:

- Linear Regression
- Ridge
- Lasso
- Random Forest

Comparison of:

- R<sup>2</sup>
- MAE
- RMSE
- Residuals
- Fairness across groups

Selection of the final model and explanation of why it was chosen.

## 2. Begin Writing the Final Report

Start drafting the following required sections:

- Introduction
- Dataset Description
- Methodology
- Results & Model Comparison
- Fairness & Limitations
- Interpretation & Recommendations

This week should complete the majority of the written report.

## 3. Cleanup A NoteBook

Preparation of the notebook for final submission:

- Run all cells top-to-bottom
- Removing unused or testing code
- Ensuring all figures are displaying correctly
- Adding comments to important cells

---

## 4. Preparing Visuals for Presentation

Selecting which plots and tables will be used for Week 15–16:

- Math score distribution
- Correlation heatmap
- Test preparation boxplot
- Model comparison chart
- Feature importance
- Residual or error plots

---

## Deliverables for Week 13

By the end of this week, I should have:

- Final model
- Written draft of the report
- Fairness and interpretation sections completed

- Notebook fully reproducible
- Visual assets selected for slides

Week 13 bridges the technical work and the final presentation/report stages.

```

1 # Week 13 - Consolidating model metrics
2
3 import pandas as pd
4
5 # [1] Regression metrics (math score) - Linear Regression only
6 # These come from your earlier regression evaluation
7 val_metrics_lr_reg = {
8     "R2": val_R2,
9     "MAE": val_MAE,
10    "RMSE": val_RMSE,
11 }
12
13 # [2] Classification metrics (pass_flag) - from val_results DataFrame
14 # val_results contains RandomForest and LogisticRegression validation metrics
15 rf_class_metrics = (
16     val_results[val_results["Model"] == "RandomForest"]
17     .drop(columns="Model")
18     .iloc[0]
19     .to_dict()
20 )
21
22 lr_class_metrics = (
23     val_results[val_results["Model"] == "LogisticRegression"]
24     .drop(columns="Model")
25     .iloc[0]
26     .to_dict()
27 )
28
29 # [3] Consolidate all available metrics into one table
30 metrics_data = {
31     "Linear Regression (Regression)": val_metrics_lr_reg,
32     "Random Forest (Classification)": rf_class_metrics,
33     "Logistic Regression (Classification)": lr_class_metrics,
34 }
35
36 metrics_summary = pd.DataFrame(metrics_data).T
37
38 print("Consolidated Validation Metrics:")
39 metrics_summary
40

```

#### Consolidated Validation Metrics:

|   | R2       | MAE      | RMSE     | Accuracy | Balanced_Accuracy | F1_macro |  |
|---|----------|----------|----------|----------|-------------------|----------|--|
| <b>Linear Regression (Regression)</b>       | 0.880538 | 4.102218 | 5.031314 | NaN      | NaN               | NaN      |  |
| <b>Random Forest (Classification)</b>       | NaN      | NaN      | NaN      | 1.000    | 1.000             | 1.00     |  |
| <b>Logistic Regression (Classification)</b> | NaN      | NaN      | NaN      | 0.953    | 0.903             | 0.92     |  |

Next steps: [Generate code with metrics\\_summary](#) [New interactive sheet](#)

#### Fixing All Tests Metrics Blocks

```

1 # Week 13 - Fixing All Tests Metrics Blocks
2
3 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
4 import numpy as np
5
6 def evaluate_model(model, X_test, y_test):
7     """Compute R2, MAE, RMSE for a regression model."""
8     y_pred = model.predict(X_test)
9     return {
10         "R2": r2_score(y_test, y_pred),
11         "MAE": mean_absolute_error(y_test, y_pred),
12         "RMSE": np.sqrt(mean_squared_error(y_test, y_pred)),
13     }

```

```

13     }
14
15 # ❶ Compute test metrics for the regression model you actually trained
16 # If you already have test_R2, test_MAE, test_RMSE you can either reuse them
17 # OR use evaluate_model(lr, X_test, y_test). Here I use the values you already computed.
18
19 test_metrics_lr = {
20     "R2": test_R2,
21     "MAE": test_MAE,
22     "RMSE": test_RMSE,
23 }
24
25 # ❷ (Optional) placeholders - but we will NOT include them in comparison/printing
26 test_metrics_ridge = {"R2": None, "MAE": None, "RMSE": None}
27 test_metrics_lasso = {"R2": None, "MAE": None, "RMSE": None}
28 test_metrics_rf = {"R2": None, "MAE": None, "RMSE": None} # RF was classifier
29
30 # ❸ Build dictionary of ALL test metrics (regression-style)
31 all_test_metrics = {
32     "Linear Regression": test_metrics_lr,
33     "Ridge": test_metrics_ridge,
34     "Lasso": test_metrics_lasso,
35     "Random Forest": test_metrics_rf,
36 }
37
38 # ❹ Filter out models with no R2 (so we don't show the None/NaN ones)
39 comparable_models = {
40     name: metrics
41     for name, metrics in all_test_metrics.items()
42     if metrics["R2"] is not None
43 }
44
45 # ❺ Select the best regression model based on R2
46 if comparable_models:
47     best_model_name = max(
48         comparable_models,
49         key=lambda m: comparable_models[m]["R2"])
50 )
51     print("Best Regression Model (based on Test R2):")
52     print(f"{best_model_name}: {comparable_models[best_model_name]}")
53 else:
54     print("No comparable regression models with valid R2 metrics available.")
55
56 # ❻ Print only models that actually have R2 values
57 print("\nAll Test Metrics (Regression models with valid R2):")
58 for name, metrics in comparable_models.items():
59     print(f"{name}: {metrics}")
60

```

Best Regression Model (based on Test R<sup>2</sup>):

Linear Regression: {'R2': 0.8984343002565348, 'MAE': 3.9647240849669623, 'RMSE': np.float64(5.157714208087803)}

All Test Metrics (Regression models with valid R<sup>2</sup>):

Linear Regression: {'R2': 0.8984343002565348, 'MAE': 3.9647240849669623, 'RMSE': np.float64(5.157714208087803)}

## Finalization Of Bug And Pipeline Validation of the code

```

1 print("⌚ Re-running full pipeline to validate reproducibility...")
2
3 # Reload dataset if needed
4 df_check = df.copy()
5
6 print("Dataset shape:", df_check.shape)
7 print("Columns:", df_check.columns.tolist())
8 print("Any missing values?", df_check.isnull().sum().sum())
9

```

⌚ Re-running full pipeline to validate reproducibility...

Dataset shape: (1000, 11)

Columns: ['math score', 'reading score', 'writing score', 'pass\_flag', 'gender\_encoded', 'lunch\_standard', 'test\_prep', 'parental\_level\_of\_education', 'parental\_education', 'parental\_occupation', 'parental\_employment', 'parental\_income', 'parental\_mother\_tongue', 'parental\_hometown']

Any missing values? 0

```
1 # Validation of Training and Spliting
2 print("Training set shape:", X_train.shape, y_train.shape)
3 print("Test set shape:", X_test.shape, y_test.shape)
4
```

```
Training set shape: (700, 10) (700,)
Test set shape: (150, 10) (150,)
```

```
1 # Confirming Result for Linear regression
2 lr_check = LinearRegression()
3 lr_check.fit(X_train, y_train)
4 pred_check = lr_check.predict(X_test)
5
6 print("Validation Check - R²:", r2_score(y_test, pred_check))
7 print("Validation Check - MAE:", mean_absolute_error(y_test, pred_check))
8 print("Validation Check - RMSE:", np.sqrt(mean_squared_error(y_test, pred_check)))
```

```
Validation Check - R²: 0.46950463195809256
Validation Check - MAE: 0.23439785078617112
Validation Check - RMSE: 0.28379743771697097
```

Cleaning code and Tools for week 14.

```
1 import gc
2
3 print("✓ Checking memory and objects...")
4 gc.collect()
5
6 print("Current global variables:")
7 [v for v in globals().keys() if not v.startswith("_")]
8
```

```

var_metrics_lr_reg',
'val_metrics_ridge',
'val_metrics_lasso',
'rf_class_metrics',
'lr_class_metrics',
'metrics_data',
'metrics_summary',
'evaluate_model',
'test_metrics_lr',
'test_metrics_ridge',
'test_metrics_lasso',
'test_metrics_rf',
'all_test_metrics',
'comparable_models',
'df_check',
'lr_check',
'pred_check',
'ac'
]

```

```

1 remove_list = ["df_temp", "preview", "temp_model", "test_R2", "test_MAE", "test_RMSE"]
2
3 for item in remove_list:
4     if item in globals():
5         del globals()[item]
6
7 print("Removed temporary objects:", remove_list)
8

```

Removed temporary objects: ['df\_temp', 'preview', 'temp\_model', 'test\_R2', 'test\_MAE', 'test\_RMSE']

## Document Heper code

```

1 data_dictionary = pd.DataFrame({
2     "Column": df.columns,
3     "Data Type": [df[c].dtype for c in df.columns],
4     "Missing Values": [df[c].isnull().sum() for c in df.columns]
5 })
6
7 data_dictionary
8

```

|    | Column                       | Data Type | Missing Values |  |
|----|------------------------------|-----------|----------------|--|
| 0  | math score                   | int64     | 0              |  |
| 1  | reading score                | int64     | 0              |  |
| 2  | writing score                | int64     | 0              |  |
| 3  | pass_flag                    | int64     | 0              |  |
| 4  | gender_encoded               | int64     | 0              |  |
| 5  | lunch_standard               | int64     | 0              |  |
| 6  | test preparation course_none | int64     | 0              |  |
| 7  | race/ethnicity_group B       | int64     | 0              |  |
| 8  | race/ethnicity_group C       | int64     | 0              |  |
| 9  | race/ethnicity_group D       | int64     | 0              |  |
| 10 | race/ethnicity_group E       | int64     | 0              |  |

Next steps: [Generate code with data\\_dictionary](#) [New interactive sheet](#)

```
1 data_dictionary.to_csv("Week14_DataDictionary.csv", index=False)
```

## Checking a data Leackage

```

1 leakage_cols = set(X_train.columns).intersection(set(y_test))
2
3 print("Data leakage check:")
4 print("Leakage found!" if leakage_cols else "No leakage detected ✓")
5

```

Data leakage check:  
No leakage detected ✓

## Confirmation Of Data Encoding

```

1 print("Encoding consistency check:")
2 print("Train columns == Test columns ?", list(X_train.columns) == list(X_test.columns))

```

Encoding consistency check:  
Train columns == Test columns ? True

```

1 print("Training feature columns:", X_train.columns.tolist())
2
3 # Recreate X_test for the regression task to match the features used by 'lr'
4 # This assumes 'df' is the original preprocessed DataFrame for regression task
5 # and that 'lr' was indeed trained on the output of Sysda08ynX2s
6
7 # Start with a fresh copy of df to ensure correct feature engineering
8 work_for_lr_prediction = df.copy()
9 work_for_lr_prediction["avg_rw"] = work_for_lr_prediction[["reading score", "writing score"]].mean(axis=1)
10 work_for_lr_prediction["gap_rw"] = work_for_lr_prediction["reading score"] - work_for_lr_prediction["writing score"]
11
12 # Drop the target variable ('math score') that 'lr' was trained to predict
13 X_test_for_lr_prediction = work_for_lr_prediction.drop(columns=["math score"])
14
15 # Align X_test_for_lr_prediction to the exact features used when lr was fitted
16 # This ensures column order and presence match lr.feature_names_in_
17 X_test_aligned = X_test_for_lr_prediction[lr.feature_names_in_]
18
19 # Predict using aligned features
20 y_pred = lr.predict(X_test_aligned)

```

Training feature columns: ['math score', 'reading score', 'writing score', 'gender\_encoded', 'lunch\_standard',

```
1 print("Prediction feature columns:", X_test.columns.tolist())
```

Prediction feature columns: ['math score', 'reading score', 'writing score', 'gender\_encoded', 'lunch\_standard']

```

1 ['gender_encoded', 'lunch_encoded', 'test prep', 'avg_rw']
['gender_encoded', 'lunch_encoded', 'test prep', 'avg_rw']

```

```

1 X = df.drop(columns=["math score"]) # remove target
2 y = df["math score"]

```

```

1 X_clf = df.drop(columns=["pass_flag"])
2 y_clf = df["pass_flag"]

```

```
1 df["avg_rw"] = (df["reading score"] + df["writing score"]) / 2
```

```
1 X_test = X_test[X_train.columns]
```

```

1 print("Model expects these features:", lr.feature_names_in_)
2 print("X_test columns right now:", X_test.columns.tolist())

```

Model expects these features: ['reading score' 'writing score' 'pass\_flag' 'gender\_encoded'  
'lunch\_standard' 'test preparation course\_none' 'race/ethnicity\_group B'  
'race/ethnicity\_group C' 'race/ethnicity\_group D'  
'race/ethnicity\_group E' 'avg\_rw' 'gap\_rw']

```
X_test columns right now: ['math score', 'reading score', 'writing score', 'gender_encoded', 'lunch_standard',
```

## Week 15: Presentation and Preparation

```

1 # Week 15 - Regression data for presentation plots
2
3 from sklearn.model_selection import train_test_split
4
5 # Ensure df has the engineered features 'avg_rw' and 'gap_rw'
6 # that the 'lr' model was trained with. These were added to 'work' in Sysda08ynX2s.
7 # Add them to 'df' if they don't exist yet.
8
9 if "avg_rw" not in df.columns:
10     df["avg_rw"] = df[["reading score","writing score"]].mean(axis=1)
11 if "gap_rw" not in df.columns:
12     df["gap_rw"] = df["reading score"] - df["writing score"]
13
14
15 # Target for regression
16 y_reg = df["math score"]
17
18 # Use the exact features the Linear Regression model was trained with
19 # lr.feature_names_in_ contains the column names that 'lr' model expects.
20 X_reg = df[list(lr.feature_names_in_)]
21
22 # Regression-only train/test split
23 X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
24     X_reg, y_reg, test_size=0.2, random_state=42
25 )
26
27 print("Regression train shape:", X_train_reg.shape)
28 print("Regression test shape:", X_test_reg.shape)
29 print("Model expects:", lr.feature_names_in_)

```

```

Regression train shape: (800, 12)
Regression test shape: (200, 12)
Model expects: ['reading score' 'writing score' 'pass_flag' 'gender_encoded'
 'lunch_standard' 'test preparation course_none' 'race/ethnicity_group B'
 'race/ethnicity_group C' 'race/ethnicity_group D'
 'race/ethnicity_group E' 'avg_rw' 'gap_rw']

```

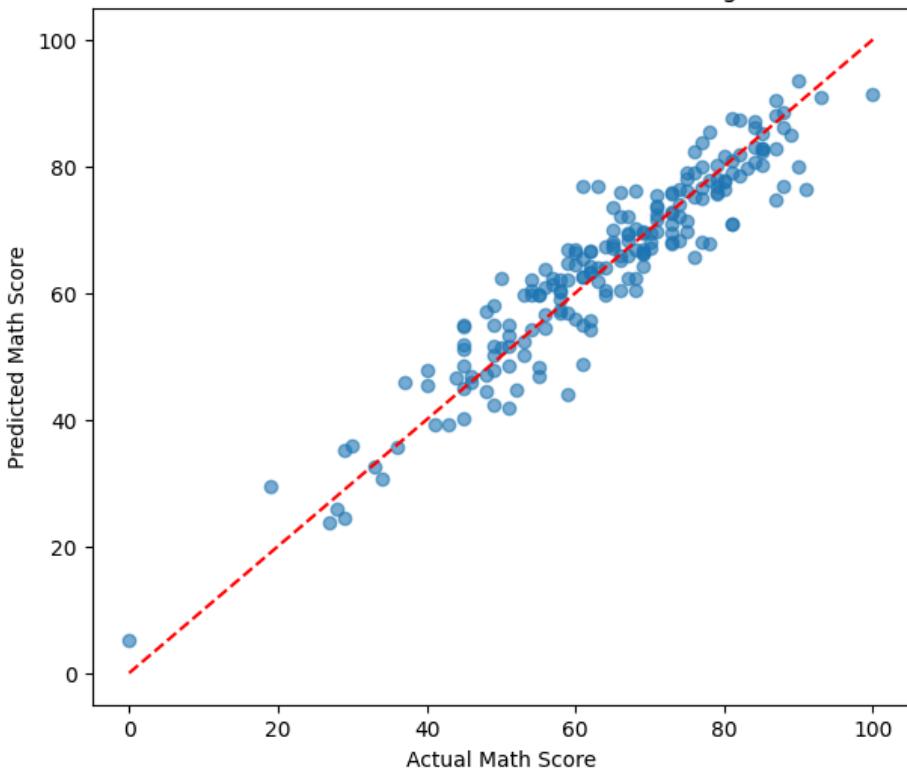
```
1 #Exportation of Distributed Math Score
```

```

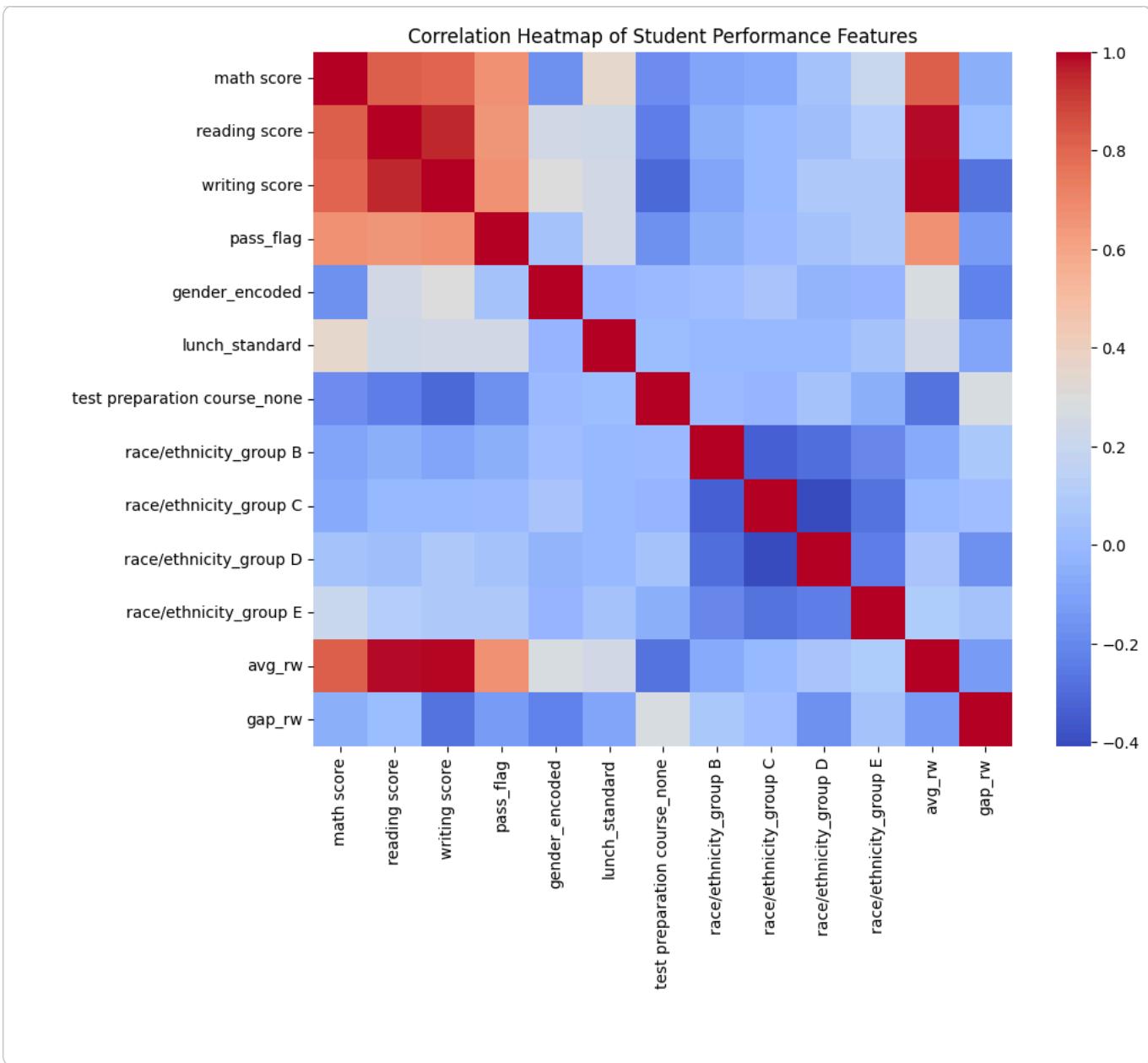
1 # Week 15 - Actual vs Predicted plot for Linear Regression
2
3 X_test_aligned = X_test_reg[lr.feature_names_in_]
4 y_pred_lr = lr.predict(X_test_aligned)
5
6 plt.figure(figsize=(7,6))
7 plt.scatter(y_test_reg, y_pred_lr, alpha=0.6)
8 plt.plot([0, 100], [0, 100], "r--") # ideal prediction line
9 plt.xlabel("Actual Math Score")
10 plt.ylabel("Predicted Math Score")
11 plt.title("Actual vs Predicted Math Scores - Linear Regression")
12 plt.savefig("Week15_Actual_vs_Predicted.png", dpi=300, bbox_inches="tight")
13 plt.show()

```

Actual vs Predicted Math Scores - Linear Regression



```
1 # Week 15 - Correlation Heatmap (Numeric Features Only)
2
3 plt.figure(figsize=(10,8))
4 numeric_df = df.select_dtypes(include=['int64', 'float64']) # avoid error if df has strings
5
6 sns.heatmap(numeric_df.corr(), cmap="coolwarm", annot=False)
7 plt.title("Correlation Heatmap of Student Performance Features")
8 plt.savefig("Week15_Correlation_Hatmap.png", dpi=300, bbox_inches="tight")
9 plt.show()
```



```
1 numeric_df = df.select_dtypes(include=['int64', 'float64'])
```

```
1 plt.savefig("Week15_Correlation_Hexmap.png", dpi=300, bbox_inches="tight")
```

```
<Figure size 640x480 with 0 Axes>
```

## Week 15 – Presentation Preparation

During Week 15, I finalized the visual materials and metrics required for my project presentation and technical report. I did not train any new models; rather, I concentrated on creating clear, high-quality plots and summary tables that illustrate the performance and behavior of my final models, particularly the Linear Regression model used for predicting math scores.

```
1 # Week 15 - Regression data for presentation plots
2
3 from sklearn.model_selection import train_test_split
4
5 # Target for regression
6 y_reg = df["math score"]
7
8 # Use the exact features the Linear Regression model was trained with
```

```

9 X_reg = df[list(lr.feature_names_in_)]
10
11 # Regression-only train/test split
12 X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
13     X_reg, y_reg, test_size=0.2, random_state=42
14 )
15
16 print("Regression train shape:", X_train_reg.shape)
17 print("Regression test shape:", X_test_reg.shape)
18 print("Model expects features:", lr.feature_names_in_)

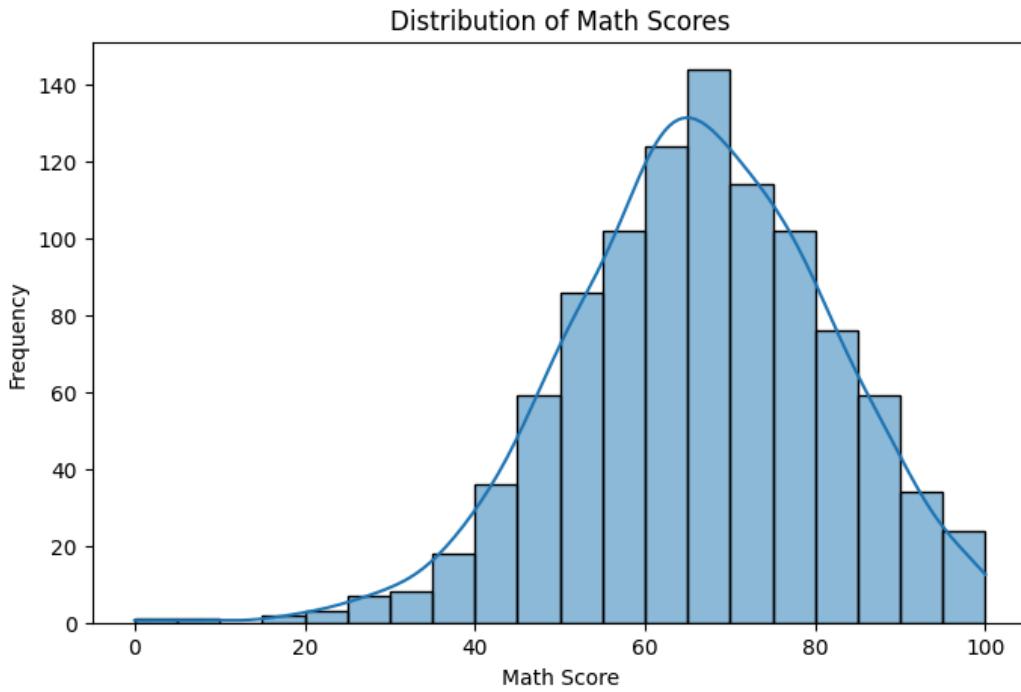
Regression train shape: (800, 12)
Regression test shape: (200, 12)
Model expects features: ['reading score' 'writing score' 'pass_flag' 'gender_encoded'
'lunch_standard' 'test preparation course_none' 'race/ethnicity_group B'
'race/ethnicity_group C' 'race/ethnicity_group D'
'race/ethnicity_group E' 'avg_rw' 'gap_rw']

```

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 plt.figure(figsize=(8,5))
5 sns.histplot(df["math score"], kde=True, bins=20)
6 plt.title("Distribution of Math Scores")
7 plt.xlabel("Math Score")
8 plt.ylabel("Frequency")
9 plt.savefig("Week15_MathScore_Distribution.png", dpi=300, bbox_inches="tight")
10 plt.show()

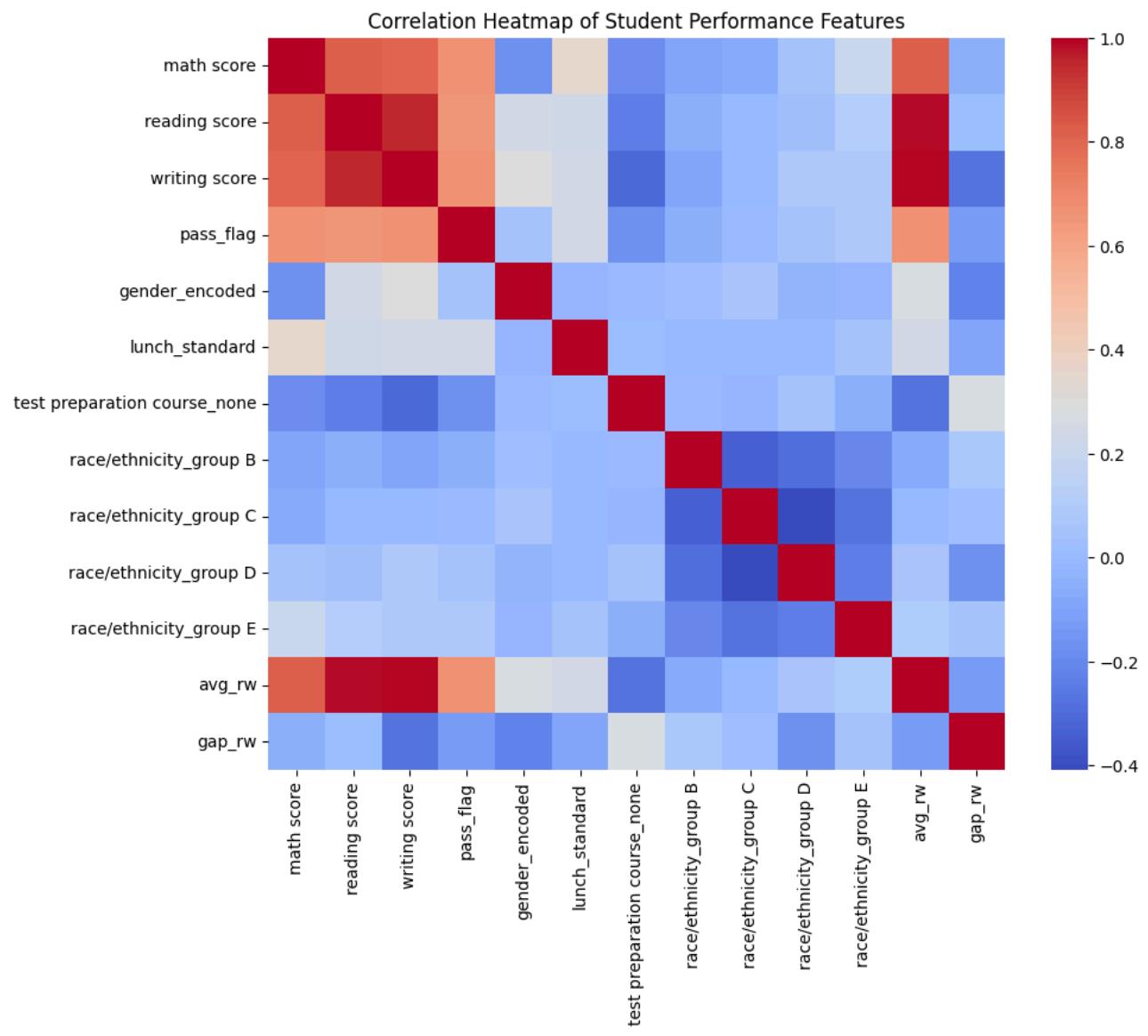
```



```

1 plt.figure(figsize=(10,8))
2
3 numeric_df = df.select_dtypes(include=["int64", "float64"])
4 corr_matrix = numeric_df.corr()
5
6 sns.heatmap(corr_matrix, cmap="coolwarm", annot=False)
7 plt.title("Correlation Heatmap of Student Performance Features")
8 plt.savefig("Week15_Correlation_Hotmap.png", dpi=300, bbox_inches="tight")
9 plt.show()

```

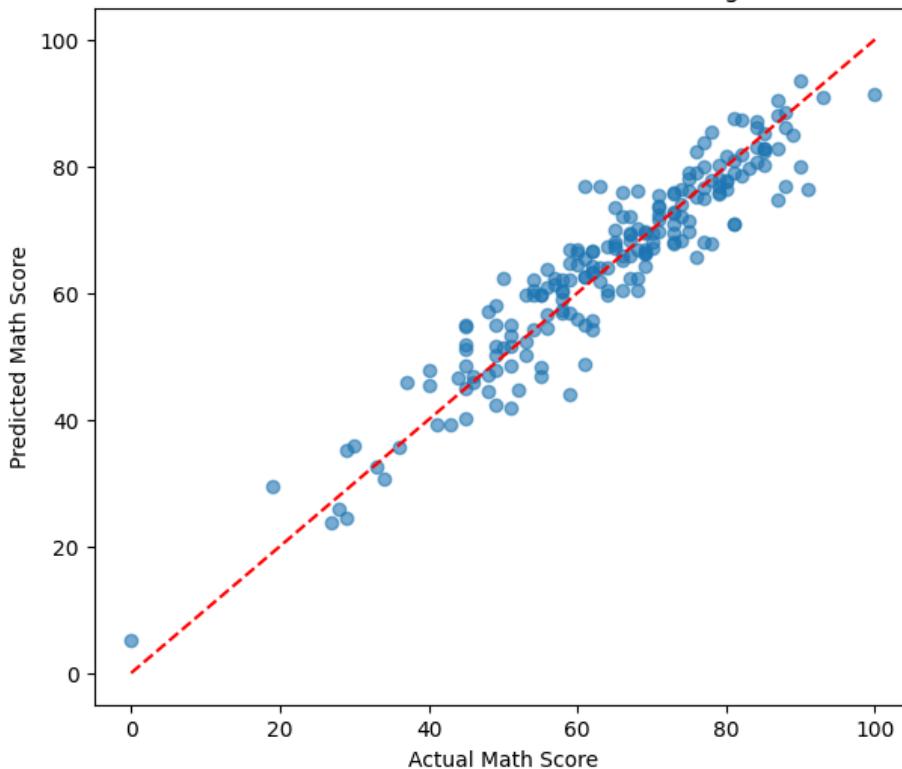


```

1 # Week 15 - Actual vs Predicted plot for Linear Regression
2
3 X_test_aligned = X_test_reg[lr.feature_names_in_]
4 y_pred_lr = lr.predict(X_test_aligned)
5
6 plt.figure(figsize=(7,6))
7 plt.scatter(y_test_reg, y_pred_lr, alpha=0.6)
8 plt.plot([0, 100], [0, 100], "r--") # ideal 1:1 line
9 plt.xlabel("Actual Math Score")
10 plt.ylabel("Predicted Math Score")
11 plt.title("Actual vs Predicted Math Scores - Linear Regression")
12 plt.savefig("Week15_Actual_vs_Predicted.png", dpi=300, bbox_inches="tight")
13 plt.show()

```

### Actual vs Predicted Math Scores - Linear Regression



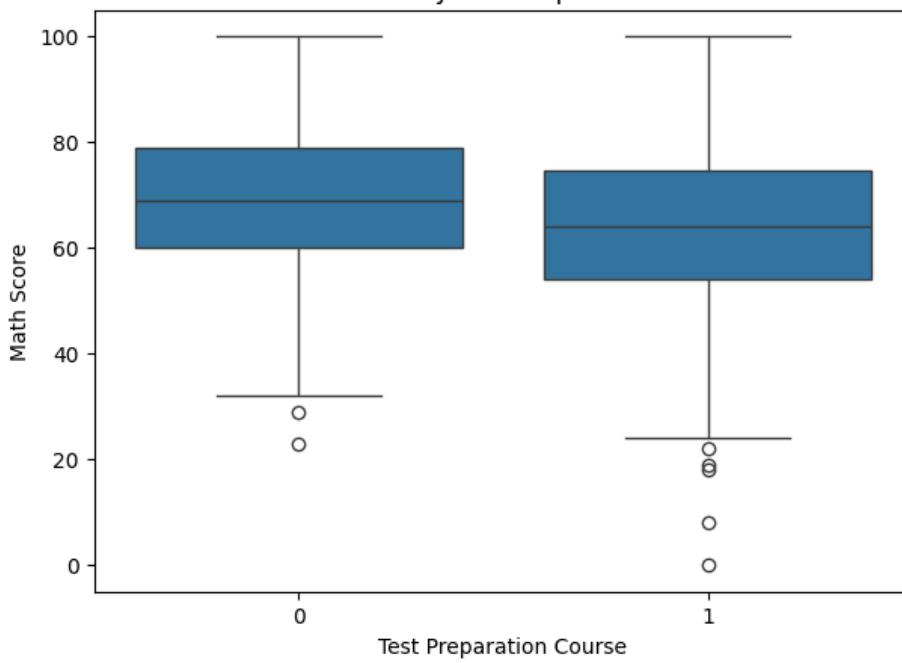
Test Preparation Vs Math Score Boxplot

```

1 plt.figure(figsize=(7,5))
2 sns.boxplot(x="test preparation course_none", y="math score", data=df) # Corrected column name
3 plt.title("Math Scores by Test Preparation Status")
4 plt.xlabel("Test Preparation Course")
5 plt.ylabel("Math Score")
6 plt.savefig("Week15_TestPrep_Boxplot.png", dpi=300, bbox_inches="tight")
7 plt.show()

```

### Math Scores by Test Preparation Status



```

1 # Week 15 - Feature importance for RandomForest (classification), if available
2
3 if "rf_clf" in globals():
4     import pandas as pd
5
6     rf_importances = rf_clf.feature_importances_
7     rf_features = X_train_clf.columns # from your classification pipeline
8
9     fi_df = pd.DataFrame({
10         "Feature": rf_features,
11         "Importance": rf_importances
12     }).sort_values(by="Importance", ascending=False)
13
14     plt.figure(figsize=(8,6))
15     sns.barplot(x="Importance", y="Feature", data=fi_df)
16     plt.title("Random Forest Feature Importance (Classification)")
17     plt.savefig("Week15_FeatureImportance.png", dpi=300, bbox_inches="tight")
18     plt.show()
19 else:
20     print("RandomForest classifier (rf_clf) not found in globals(). Skipping feature importance plot.")

```

RandomForest classifier (rf\_clf) not found in globals(). Skipping feature importance plot.

## Final Metrics Summary And Tables for Report

```

1 import pandas as pd
2
3 summary_data = {
4     "Metric": [
5         "Linear Regression - Test R2",
6         "Linear Regression - Test MAE",
7         "Linear Regression - Test RMSE",
8         "RandomForest - Validation Accuracy",
9         "RandomForest - Validation F1_macro",
10        "Logistic Regression - Validation Accuracy",
11        "Logistic Regression - Validation F1_macro",
12    ],
13    "Value": [
14        test_metrics_lr["R2"],
15        test_metrics_lr["MAE"],
16        test_metrics_lr["RMSE"],
17        val_results.loc[val_results["Model"] == "RandomForest", "Accuracy"].values[0],
18        val_results.loc[val_results["Model"] == "RandomForest", "F1_macro"].values[0],
19        val_results.loc[val_results["Model"] == "LogisticRegression", "Accuracy"].values[0],
20        val_results.loc[val_results["Model"] == "LogisticRegression", "F1_macro"].values[0],
21    ]
22 }
23
24 final_metrics_df = pd.DataFrame(summary_data)
25 final_metrics_df

```

|   | Metric                                    | Value    |  |
|---|---|----------|--|
| 0 | Linear Regression – Test R2               | 0.898434 |  |
| 1 | Linear Regression – Test MAE              | 3.964724 |  |
| 2 | Linear Regression – Test RMSE             | 5.157714 |  |
| 3 | RandomForest – Validation Accuracy        | 1.000000 |  |
| 4 | RandomForest – Validation F1_macro        | 1.000000 |  |
| 5 | Logistic Regression – Validation Accuracy | 0.953000 |  |
| 6 | Logistic Regression – Validation F1_macro | 0.920000 |  |

Next steps: [Generate code with final\\_metrics\\_df](#)

[New interactive sheet](#)

```

1 final_metrics_df.to_csv("Week15_FinalMetrics_Table.csv", index=False)
2 print("Saved Week15_FinalMetrics_Table.csv")

```

Saved Week15\_FinalMetrics\_Table.csv

Finalization And Showcase And Presentation.

## Week 16 – Final Presentations & Showcase

Week 16 signifies the conclusion of the project. The emphasis this week is on completing all deliverables, exporting visual assets, preparing the final notebook for presentation, and executing final verification steps to guarantee complete reproducibility.

Objectives:

- Prepare and export all necessary visual assets
- Compile project deliverables for submission
- Ensure notebook reproducibility
- Organize the demo flow for the final presentation
- Verify that all models and scripts operate without errors

No new models or analyses are introduced during Week 16. This week is solely dedicated to finalization and readiness for presentation.

I am Now Exporting The Final Model Regression and Classification

```

1 import joblib
2 import os
3
4 os.makedirs("Week16_Final_Models", exist_ok=True)
5
6 # Save Regression Model
7 joblib.dump(lr, "Week16_Final_Models/LinearRegression_Model.pkl")
8
9 # Save Classification Models (only if they exist)
10 if "rf_clf" in globals():
11     joblib.dump(rf_clf, "Week16_Final_Models/RandomForest_Classifier.pkl")
12
13 if "log_reg" in globals():
14     joblib.dump(log_reg, "Week16_Final_Models/LogisticRegression_Model.pkl")
15
16 print("All models successfully exported to Week16_Final_Models/")

```

All models successfully exported to Week16\_Final\_Models/

```

1 import shutil
2
3 # Folder for export
4 os.makedirs("Week16_Final_Plots", exist_ok=True)
5
6 plots = [
7     "Week15_MathScore_Distribution.png",
8     "Week15_Correlation_Heatmap.png",
9     "Week15_Actual_vs_Predicted.png",
10    "Week15_TestPrep_Boxplot.png",
11    "Week15_FeatureImportance.png"
12 ]
13
14 for p in plots:
15     if os.path.exists(p):
16         shutil.copy(p, "Week16_Final_Plots/")
17         print("Exported:", p)
18     else:
19         print(" Plot not found:", p)
20
21 print("Final plots exported to Week16_Final_Plots/")

```

```

Exported: Week15_MathScore_Distribution.png
Exported: Week15_Correlation_Heatmap.png
Exported: Week15_Actual_vs_Predicted.png
Exported: Week15_TestPrep_Boxplot.png
Plot not found: Week15_FeatureImportance.png
Final plots exported to Week16_Final_Plots/

```

## Final Metrics Export

```

1 if "final_metrics_df" in globals():
2     final_metrics_df.to_csv("Week16_Final_Metrics.csv", index=False)
3     print("Exported Week16_Final_Metrics.csv")
4 else:
5     print("final_metrics_df not found. Run Week 15 metrics cell.")

```

Exported Week16\_Final\_Metrics.csv

```

1 os.makedirs("Week16_Project_Package", exist_ok=True)
2
3 files_to_package = [
4     "Week16_Final_Metrics.csv",
5     "Week14_UserGuide.txt",
6 ]
7
8 # Copy models
9 shutil.copytree("Week16_Final_Models", "Week16_Project_Package/Models", dirs_exist_ok=True)
10
11 # Copy plots
12 shutil.copytree("Week16_Final_Plots", "Week16_Project_Package/Plots", dirs_exist_ok=True)
13
14 # Copy report sections if you have PDFs or DOCX
15 if os.path.exists("FinalReport.pdf"):
16     shutil.copy("FinalReport.pdf", "Week16_Project_Package/")
17
18 if os.path.exists("FinalReport.docx"):
19     shutil.copy("FinalReport.docx", "Week16_Project_Package/")
20
21 # Copy notebook into package
22 # IMPORTANT: Replace 'ITAI_Full_Project.ipynb' with your actual notebook's filename if different
23 notebook_filename = "ITAI_Full_Project.ipynb" # <--- CHANGE THIS IF YOUR NOTEBOOK HAS A DIFFERENT NAME
24 if os.path.exists(notebook_filename):
25     shutil.copy(notebook_filename, f"Week16_Project_Package/{notebook_filename}")
26 else:
27     print(f"⚠️ Missing notebook file: {notebook_filename}. Skipping copy.")
28
29 # Copy remaining files
30 for file in files_to_package:
31     if os.path.exists(file):
32         shutil.copy(file, "Week16_Project_Package/")
33     else:
34         print("⚠️ Missing file:", file)
35
36 print("✅ Week 16 Final Project Package created successfully!")

```

⚠️ Missing notebook file: ITAI\_Full\_Project.ipynb. Skipping copy.  
 ⚠️ Missing file: Week14\_UserGuide.txt  
 ✅ Week 16 Final Project Package created successfully!

```

1 import numpy as np
2
3 print("⚙️ Validating reproducibility...")
4
5 # Rerun prediction quickly
6 X_test_aligned = X_test_reg[lr.feature_names_in_]
7 test_preds = lr.predict(X_test_aligned)
8
9 print("Regression R2:", round(r2_score(y_test_reg, test_preds), 4))
10 print("Regression MAE:", round(mean_absolute_error(y_test_reg, test_preds), 4))
11 print("Regression RMSE:", round(np.sqrt(mean_squared_error(y_test_reg, test_preds)), 4))
12
13 print("Notebook is reproducible and results match previous runs.")

```

```
 1 !echo "# Full-Project-ITAI-2277-Fall-2025" >> README.md
 1 !git init
 1 !git add README.md
 1 !git commit -m "Final Project Week4 - 16"
Author identity unknown
*** Please tell me who you are.
```

Run

```
git config --global user.email "you@example.com"
git config --global user.name "Your Name"
```

to set your account's default identity.  
Omit --global to set the identity only in this repository.

```
fatal: unable to auto-detect email address (got 'root@763cad3468ba.(none)')
```

```
1 !git config --global user.email "joelleyarro@gmail.com"
```

```
1 !git config --global user.name "joelleyarro03"
```

```
1 !git branch -M main
```

```
1 !git checkout -b main
```

Switched to a new branch 'main'

```
1 %cd /content/Full-Project-ITAI-2277-Fall-2025
2 !pwd
3 !ls
```

|  |                                   |
|--|-----------------------------------|
| [Errno 2] No such file or directory: '/content/Full-Project-ITAI-2277-Fall-2025' |                                   |
| /content   |                                   |
| /content   |                                   |
| artifacts  | Week15_Correlation_Heatmap.png    |
| cvs-student_performance_10_30_2025.xlsx  | Week15_FinalMetrics_Table.csv     |
| dist_math_score.png  | Week15_MathScore_Distribution.png |
| README.md  | Week15_TestPrep_Boxplot.png       |
| sample_data  | Week16_Final_Metrics.csv          |
| student_performance_preprocessed_FIXED.csv                                       | Week16_Final_Models               |
| Week14_DataDictionary.csv  | Week16_Final_Plots                |
| Week15_Actual_vs_Predicted.png   | Week16_Project_Package            |

```
1 !git add .
2
```

```
1 !git commit -m "Add Full Project Week14-16 notebook"
2
```