

Upload Dataset

```
1 #Obesity Risk Prediction
2 from google.colab import files
3 uploaded = files.upload()
```

Choose Files 3 files

**sample\_submission.csv**(text/csv) - 276818 bytes, last modified: 9/28/2025 - 100% done

**test.csv**(text/csv) - 606882 bytes, last modified: 9/28/2025 - 100% done

**train.csv**(text/csv) - 2036684 bytes, last modified: 9/28/2025 - 100% done

Saving sample\_submission.csv to sample\_submission.csv

Saving test.csv to test.csv

Saving train.csv to train.csv

Load Libraries

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import LabelEncoder
7 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
8 from xgboost import XGBClassifier
9 import warnings
10 warnings.filterwarnings('ignore')
```

Load Data

```
1 train = pd.read_csv("train.csv")# Reading training and test data
2 test = pd.read_csv("test.csv")
3 test_ids = test["id"]# Keeping test IDs for final submission
4 # Removing ID column from train and test
5 train = train.drop(columns=['id'])
6 test = test.drop(columns=['id'])
```

DATA shape and Basic Overview of data

```
1 print("Train shape:", train.shape)
2 print("Test shape:", test.shape)
3 print("\nColumns:", train.columns.tolist())
4 print("\nMissing values:\n", train.isnull().sum())
```

Train shape: (15533, 18)

Test shape: (5225, 17)

Columns: ['id', 'Gender', 'Age', 'Height', 'Weight', 'family\_history\_with\_overweight', 'FAVC', 'FCVC', 'NCP', 'CAEC', 'SMOKE', 'CH2O', 'SCC', 'FAF', 'TUE', 'CALC', 'MTRANS', 'WeightCategory', 'type']

Missing values:

id	0
Gender	0
Age	0
Height	0
Weight	0
family_history_with_overweight	0
FAVC	0
FCVC	0
NCP	0
CAEC	0
SMOKE	0
CH2O	0
SCC	0
FAF	0
TUE	0
CALC	0
MTRANS	0
WeightCategory	0

dtype: int64

```
1 train.head()
```

	id	Gender	Age	Height	Weight	family_history_with_overweight	FAVC	FCVC	NCP	CAEC	SMOKE	CH2O	SCC	
0	0	Male	24.443011	1.699998	81.669950	yes	yes	2.000000	2.983297	Sometimes	no	2.763573	no	0.0
1	1	Female	18.000000	1.560000	57.000000	yes	yes	2.000000	3.000000	Frequently	no	2.000000	no	1.0
2	2	Female	18.000000	1.711460	50.165754	yes	yes	1.880534	1.411685	Sometimes	no	1.910378	no	0.8
3	3	Female	20.952737	1.710730	131.274851	yes	yes	3.000000	3.000000	Sometimes	no	1.674061	no	1.4
4	4	Male	31.641081	1.914186	93.798055	yes	yes	2.679664	1.971472	Sometimes	no	1.979848	no	1.9

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

Exploratory Data Analysis(EDA)

Data Preprocessing

```
1 # Converting Yes/No columns to 1/0
2 yes_no_cols = ['family_history_with_overweight', 'FAVC', 'SMOKE', 'SCC']
3 for col in yes_no_cols:
4     train[col] = train[col].map({'yes': 1, 'no': 0})
5     test[col] = test[col].map({'yes': 1, 'no': 0})
```

```
1 # Frequency-type categorical mapping
2 mapping = {'no': 0, 'Sometimes': 1, 'Frequently': 2, 'Always': 3}
3 for col in ['CAEC', 'CALC']:
4     train[col] = train[col].map(mapping)
5     test[col] = test[col].map(mapping)
```

```
1 # Encode transport column (MTRANS)
2 le_mtrans = LabelEncoder()
3 combined_transport = pd.concat([train['MTRANS'], test['MTRANS']], axis=0)
4 le_mtrans.fit(combined_transport)
5 train['MTRANS'] = le_mtrans.transform(train['MTRANS'])
6 test['MTRANS'] = le_mtrans.transform(test['MTRANS'])
```

```
1 # Encode target
2 le_weight = LabelEncoder()
3 train['WeightCategory'] = le_weight.fit_transform(train['WeightCategory'])
4 print("\ncomplete!")
```

complete!

Train-Validation Split

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import LabelEncoder
3 from xgboost import XGBClassifier
4 from sklearn.metrics import accuracy_score
5 import matplotlib.pyplot as plt
6
```

```
1 # Separating input and target
2 X = train.drop(columns=['WeightCategory'])
3 y = train['WeightCategory']
4
5 # Splitting into train and validation sets
6 X_train, X_val, y_train, y_val = train_test_split(
7     X, y, test_size=0.2, random_state=42, stratify=y
8 )
9 print(f"Train size: {X_train.shape}, Validation size: {X_val.shape}")
10
```

Train size: (12426, 17), Validation size: (3107, 17)

```
1 # Encode categorical columns
2 for col in X_train.select_dtypes(include='object').columns:
3     le = LabelEncoder()
4     X_train[col] = le.fit_transform(X_train[col])
5     X_val[col] = le.transform(X_val[col])
```

XGBoost Parameters

```
1 # Defining refined XGBoost parameters
2
3 refined_params = {
4     'objective': 'multi:softmax', # multi-class classification
5     'num_class': len(np.unique(y)), # number of output categories
```

```

6 'eval_metric': 'mlogloss',      # loss function to evaluate
7 'use_label_encoder': False,    # avoid old warnings
8 'tree_method': 'hist',        # fast histogram-based algorithm
9 'grow_policy': 'lossguide',    # builds deeper trees efficiently
10 'random_state': 42,           # for reproducibility
11
12 # Fine-tuned hyperparameters
13 'learning_rate': 0.021,        # how fast model learns (lower = slower, more stable)
14 'max_depth': 6,               # max depth of each tree
15 'min_child_weight': 2,        # minimum data in a leaf
16 'subsample': 0.71,           # fraction of training data per tree
17 'colsample_bytree': 0.74,     # fraction of features per tree
18 'gamma': 0.8,                # controls overfitting
19 'reg_lambda': 2.4,           # L2 regularization
20 'reg_alpha': 0.22,           # L1 regularization
21 'max_bin': 290               # how fine-grained histogram splits are
22 }

```

```

1 from xgboost import XGBClassifier
2 from sklearn.metrics import accuracy_score
3 best_acc = 0.0 # Manual early stopping
4 best_n = None
5 best_model = None
6 accuracy_scores = []

```

```

1 refined_params.pop('n_estimators', None)
2
3 n_estimator_values = [900, 1000, 1100, 1200, 1220, 1230, 1250, 1270, 1300, 1400]
4 # Loop through different n_estimators values
5 for n in n_estimator_values:
6     temp_model = XGBClassifier(**refined_params, n_estimators=n)
7     temp_model.fit(X_train, y_train)
8
9     preds = temp_model.predict(X_val)
10    acc = accuracy_score(y_val, preds)
11    accuracy_scores.append(acc)
12
13    print(f"n_estimators={n} → Validation Accuracy: {acc*100:.3f}%")
14
15    # Keep track of best performing model
16    if acc > best_acc:
17        best_acc = acc
18        best_n = n
19        best_model = temp_model
20
21 print(f"\n=== Best Model Summary ===")
22 print(f"Best Validation Accuracy: {best_acc*100:.3f}%")
23 print(f"Best n_estimators: {best_n}")
24

```

```

n_estimators=900 → Validation Accuracy: 90.441%
n_estimators=1000 → Validation Accuracy: 90.409%
n_estimators=1100 → Validation Accuracy: 90.505%
n_estimators=1200 → Validation Accuracy: 90.473%
n_estimators=1220 → Validation Accuracy: 90.441%
n_estimators=1230 → Validation Accuracy: 90.473%
n_estimators=1250 → Validation Accuracy: 90.505%
n_estimators=1270 → Validation Accuracy: 90.441%
n_estimators=1300 → Validation Accuracy: 90.409%
n_estimators=1400 → Validation Accuracy: 90.473%

```

```

=== Best Model Summary ===
Best Validation Accuracy: 90.505%
Best n_estimators: 1100

```

## Retrain Final Model on Full Data

```

1 from sklearn.preprocessing import LabelEncoder
2 from xgboost import XGBClassifier
3
4 # Make a copy of X to avoid modifying original dataframe
5 X_full = X.copy()
6 y_full = y.copy()
7
8 # Encode categorical columns
9 for col in X_full.select_dtypes(include='object').columns:
10     le = LabelEncoder()
11     X_full[col] = le.fit_transform(X_full[col])
12
13 # Add best n_estimators to parameters
14 refined_params['n_estimators'] = best_n
15
16 # Train final model on full dataset
17 final_model = XGBClassifier(**refined_params)
18 final_model.fit(X_full, y_full)

```

```

XGBClassifier
XGBClassifier(base_score=None, booster=None, callbacks=None,
               colsample_bylevel=None, colsample_bynode=None,
               colsample_bytree=0.74, device=None, early_stopping_rounds=None,
               enable_categorical=False, eval_metric='mlogloss',
               feature_types=None, feature_weights=None, gamma=0.8,
               grow_policy='lossguide', importance_type=None,
               interaction_constraints=None, learning_rate=0.021, max_bin=290,
               max_cat_threshold=None, max_cat_to_onehot=None,
               max_delta_step=None, max_depth=6, max_leaves=None,
               min_child_weight=2, missing=nan, monotone_constraints=None,
               multi_strategy=None, n_estimators=1100, n_jobs=None, num_class=7, ...)

```

## Predictions on Test Set

```

1 # Make a copy of the test set
2 X_test = test.copy()
3
4 # Encode categorical columns in test set
5 for col in X_test.select_dtypes(include='object').columns:
6     le = LabelEncoder()
7     # Fit-transform on training data, transform on test
8     # Here, we assume you have stored the mapping from training
9     # If not, simplest is to fit on test too (may cause mismatch)
10    X_test[col] = le.fit_transform(X_test[col])
11
12 # Predict on test data
13 test_preds = final_model.predict(X_test)
14
15 # Decode labels back to original
16 pred_labels = le_weight.inverse_transform(test_preds)
17
18 # Prepare submission
19 submission = pd.DataFrame({
20     "id": test_ids,
21     "WeightCategory": pred_labels
22 })
23
24 submission.to_csv("submission_final.csv", index=False)
25 print("Submission saved successfully")
26

```

Submission saved successfully

```

1 from google.colab import files
2 files.download("submission_final.csv")

```

```

1 sns.set(style='whitegrid')
2
3 print("\n===== BASIC INFO =====")
4 print("Train shape:", train.shape)
5 print("Test shape:", test.shape)
6 display(train.info())

```

```

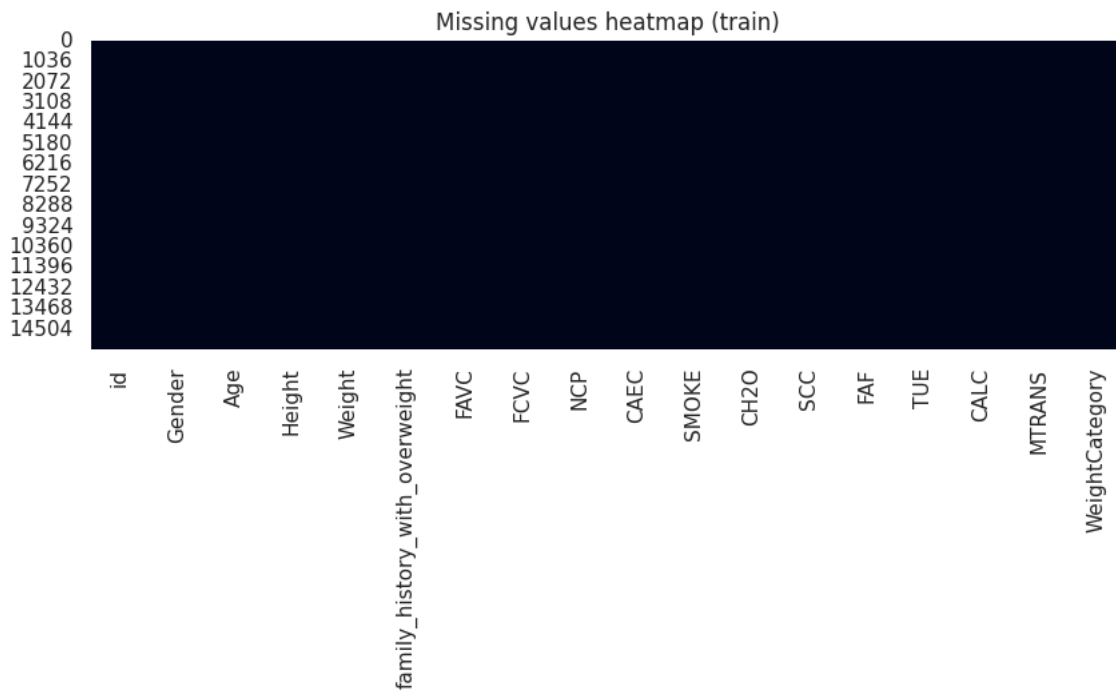
===== BASIC INFO =====
Train shape: (15533, 18)
Test shape: (5225, 17)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 15533 entries, 0 to 15532
Data columns (total 18 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   id                                    15533 non-null  int64
1   Gender                              15533 non-null  object
2   Age                                  15533 non-null  float64
3   Height                              15533 non-null  float64
4   Weight                              15533 non-null  float64
5   family_history_with_overweight      15533 non-null  int64
6   FAVC                                15533 non-null  int64
7   FCVC                                15533 non-null  float64
8   NCP                                  15533 non-null  float64
9   CAEC                                15533 non-null  int64
10  SMOKE                               15533 non-null  int64
11  CH2O                                15533 non-null  float64
12  SCC                                  15533 non-null  int64
13  FAF                                  15533 non-null  float64
14  TUE                                  15533 non-null  float64
15  CALC                                15533 non-null  int64
16  MTRANS                              15533 non-null  int64
17  WeightCategory                      15533 non-null  int64
dtypes: float64(8), int64(9), object(1)
memory usage: 2.1+ MB
None

```

```

1 # Missing values
2 plt.figure(figsize=(10,3))
3 sns.heatmap(train.isnull(), cbar=False)
4 plt.title("Missing values heatmap (train)")
5 plt.show()

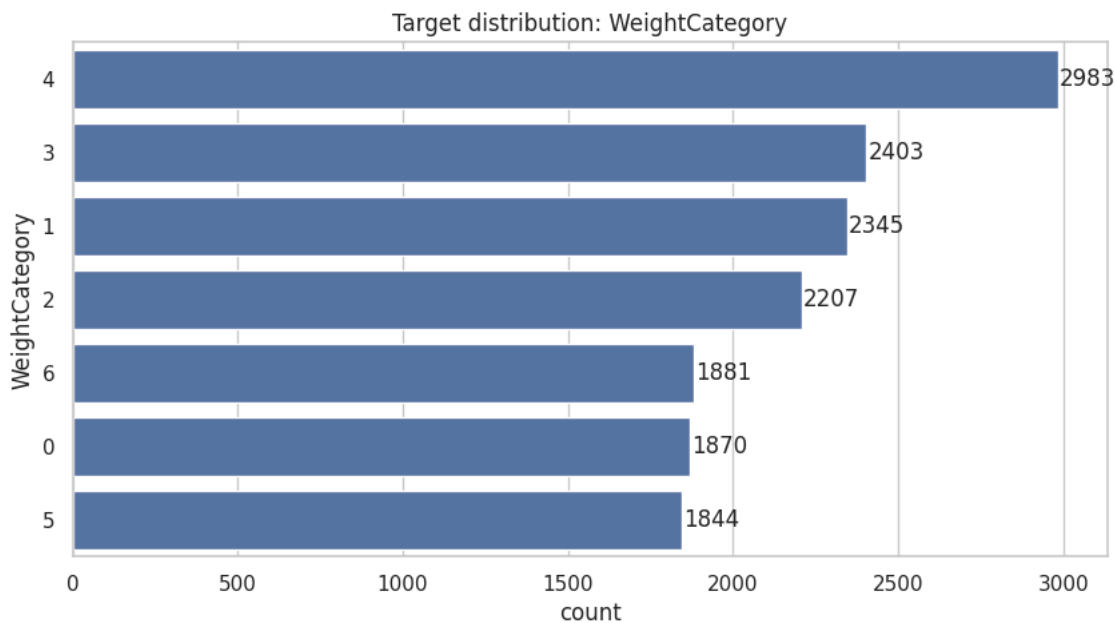
```



```

1 # Target distribution
2 target_col = 'WeightCategory'
3 plt.figure(figsize=(10,5))
4 order = train[target_col].value_counts().index
5 ax = sns.countplot(y=target_col, data=train, order=order)
6 plt.title("Target distribution: WeightCategory")
7 for p in ax.patches:
8     ax.text(p.get_width()+5, p.get_y()+p.get_height()/2,
9            int(p.get_width()), va='center')
10 plt.show()
11

```



```

1
2 print("\nTarget Class Percentages:")
3 print((train[target_col].value_counts(normalize=True)*100).round(2))
4
5 # Numeric summary
6 numeric_cols = train.select_dtypes(include=[np.number]).columns.tolist()
7 display(train[numeric_cols].describe().T)

```

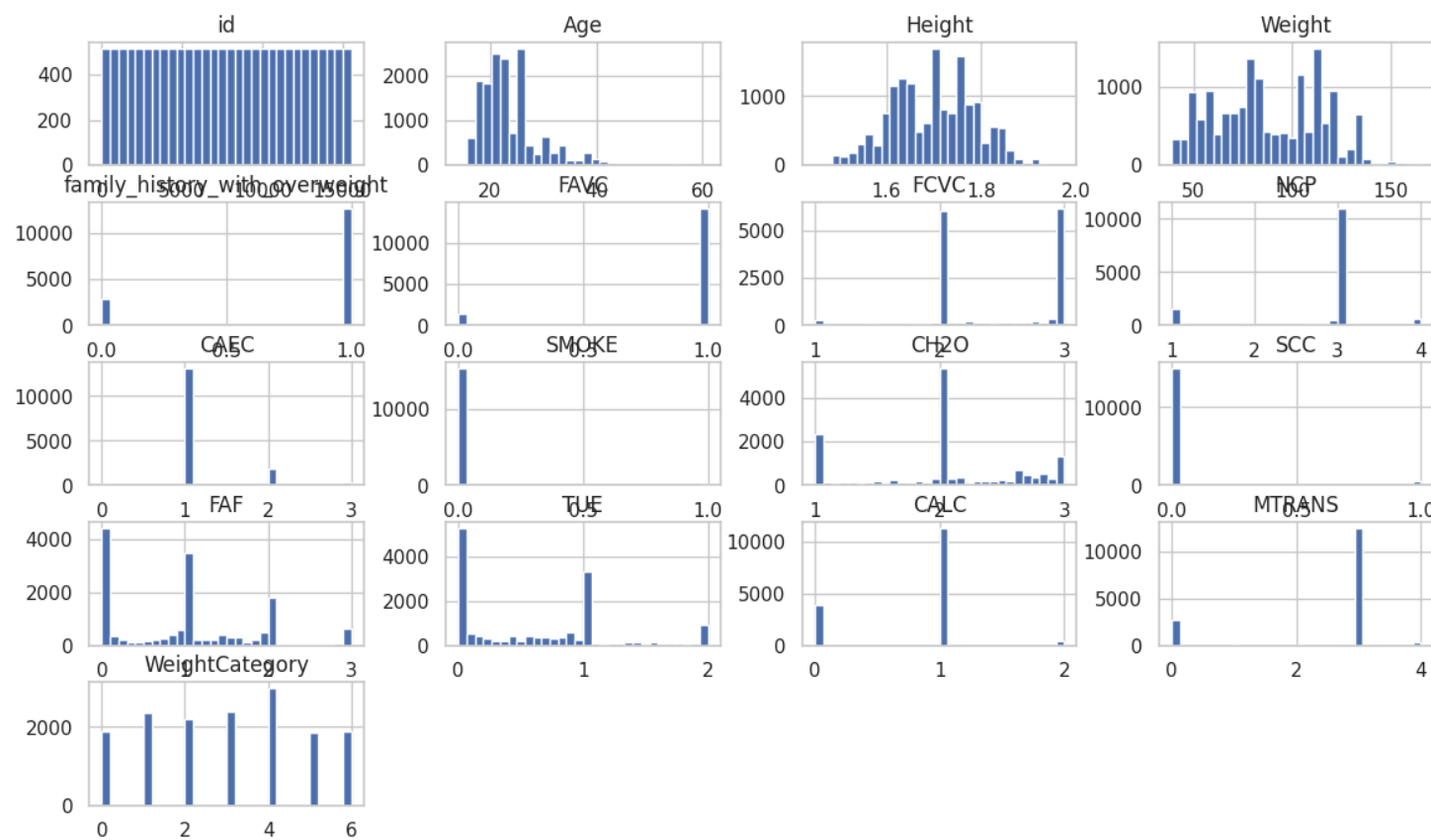
Target Class Percentages:  
WeightCategory  
4 19.20  
3 15.47  
1 15.10  
2 14.21  
6 12.11  
0 12.04  
5 11.87  
Name: proportion, dtype: float64

	count	mean	std	min	25%	50%	75%	max
id	15533.0	7766.000000	4484.135201	0.00	3883.000000	7766.000000	11649.000000	15532.000000
Age	15533.0	23.816308	5.663167	14.00	20.000000	22.771612	26.000000	61.000000
Height	15533.0	1.699918	0.087670	1.45	1.630927	1.700000	1.762921	1.975663
Weight	15533.0	87.785225	26.369144	39.00	66.000000	84.000000	111.600553	165.057269
family_history_with_overweight	15533.0	0.817357	0.386386	0.00	1.000000	1.000000	1.000000	1.000000
FAVC	15533.0	0.913153	0.281620	0.00	1.000000	1.000000	1.000000	1.000000
FCVC	15533.0	2.442917	0.530895	1.00	2.000000	2.342220	3.000000	3.000000
NCP	15533.0	2.760425	0.706463	1.00	3.000000	3.000000	3.000000	4.000000
CAEC	15533.0	1.151098	0.446058	0.00	1.000000	1.000000	1.000000	3.000000
SMOKE	15533.0	0.011395	0.106141	0.00	0.000000	0.000000	0.000000	1.000000
CH2O	15533.0	2.027626	0.607733	1.00	1.796257	2.000000	2.531456	3.000000
SCC	15533.0	0.033091	0.178880	0.00	0.000000	0.000000	0.000000	1.000000
FAF	15533.0	0.976968	0.836841	0.00	0.007050	1.000000	1.582675	3.000000
TUE	15533.0	0.613813	0.602223	0.00	0.000000	0.566353	1.000000	2.000000
CALC	15533.0	0.778922	0.473942	0.00	1.000000	1.000000	1.000000	2.000000
MTRANS	15533.0	2.501384	1.152353	0.00	3.000000	3.000000	3.000000	4.000000
WeightCategory	15533.0	2.987575	1.893756	0.00	1.000000	3.000000	4.000000	6.000000



```
1 # Histograms
2 train[numeric_cols].hist(bins=30, figsize=(14,8))
3 plt.suptitle("Histograms of Numeric Features")
4 plt.show()
```

Histograms of Numeric Features



```

1 # Categorical feature distributions
2 cat_cols = train.select_dtypes(include=['object']).columns.tolist()
3 if cat_cols:
4     plt.figure(figsize=(15,4*len(cat_cols)//3))
5     for i, c in enumerate(cat_cols, 1):
6         plt.subplot((len(cat_cols)+2)//3, 3, i)
7         sns.countplot(x=c, data=train, order=train[c].value_counts().index)
8         plt.title(c)
9         plt.xticks(rotation=45)
10    plt.tight_layout()
11    plt.show()

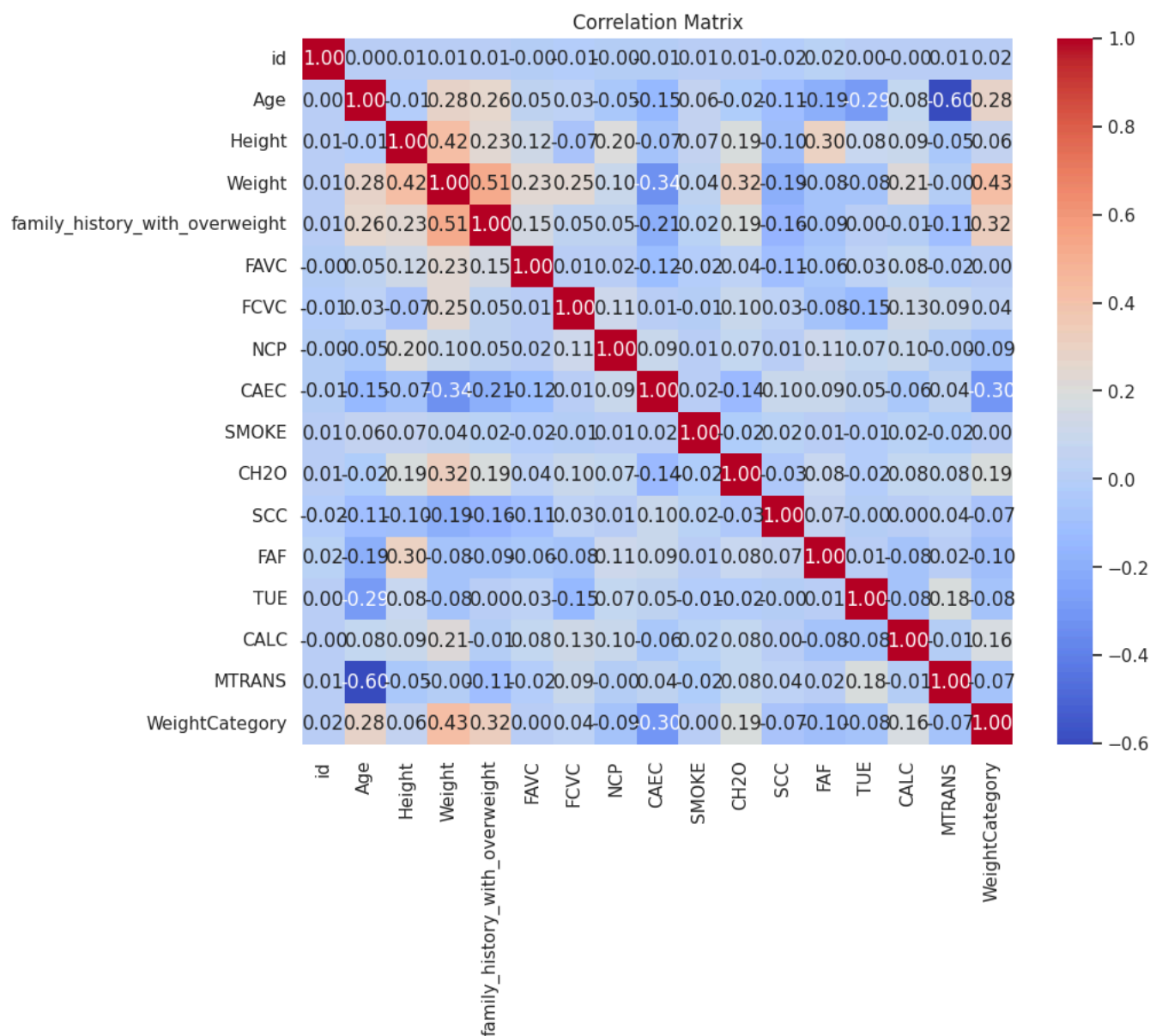
```



```

1 # Correlation matrix
2 plt.figure(figsize=(10,8))
3 sns.heatmap(train[numeric_cols].corr(), annot=True, fmt=".2f", cmap="coolwarm")
4 plt.title("Correlation Matrix")
5 plt.show()
6

```

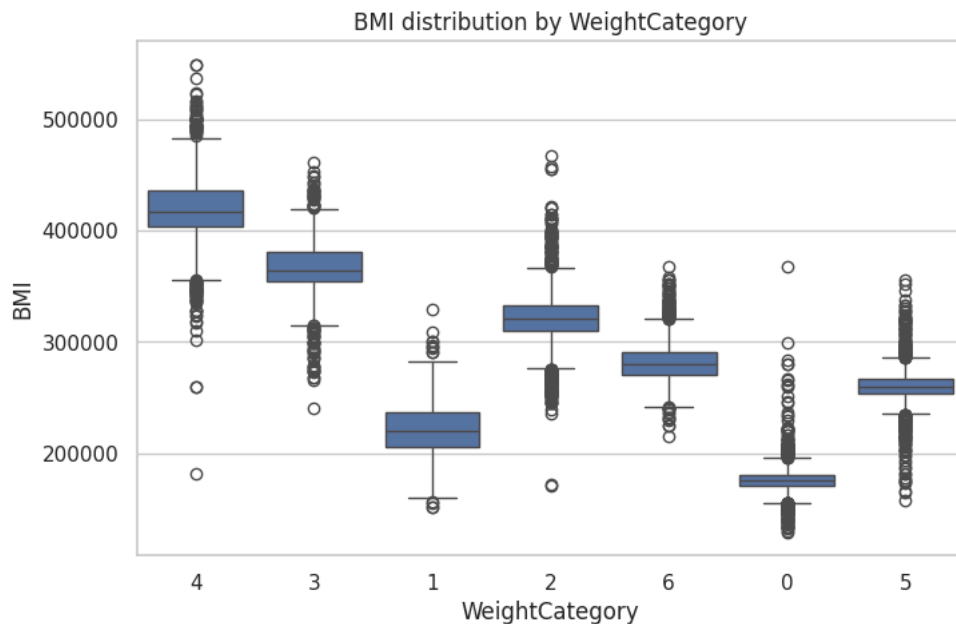


```

1 # BMI analysis if Height/Weight present
2 if set(['Height', 'Weight']).issubset(train.columns):
3     train['Height_m'] = train['Height'] / 100.0
4     train['BMI'] = train['Weight'] / (train['Height_m']**2)
5     plt.figure(figsize=(8,5))
6     sns.boxplot(x=target_col, y='BMI', data=train, order=order)
7     plt.title("BMI distribution by WeightCategory")

```

```
8 plt.show()
9 display(train.groupby(target_col)['BMI'].describe().T)
```



WeightCategory	0	1	2	3	4	5	6
count	1870.000000	2345.000000	2207.000000	2403.000000	2983.000000	1844.000000	1881.000000
mean	176208.069429	220227.886205	321409.896045	365718.139940	418170.010105	260918.382493	282031.880946
std	14360.967040	21992.403195	25631.576473	21144.410780	26903.022934	18002.078553	17973.725467
min	128685.407075	150947.953146	170992.784105	240484.601993	181786.703601	157618.802836	215138.585105
25%	170373.404091	205693.296602	310204.081633	354453.927155	404162.831736	253906.250000	270992.510068
50%	175324.670676	220385.674931	320799.694079	363906.286912	417700.837480	259819.784973	280557.057869
75%	180620.379165	237118.446310	333205.177496	380716.248163	436348.210102	266727.632983	290914.331827
max	367781.150966	328824.141519	468051.876636	462224.831931	549979.913613	355555.555556	367414.556033

### Model Evaluation (Validation Performance)

```
1 y_pred = best_model.predict(X_val)
2 cm = confusion_matrix(y_val, y_pred)
3 acc = accuracy_score(y_val, y_pred)
4
5 print(f"\nBest Validation Accuracy: {acc*100:.3f}%")
6 print("\nClassification Report:\n")
7 print(classification_report(y_val, y_pred))
8
9 plt.figure(figsize=(6,5))
10 sns.heatmap(cm, annot=True, fmt='d', cmap='coolwarm', cbar=False)
11 plt.title("Confusion Matrix for Best Model")
12 plt.xlabel("Predicted")
13 plt.ylabel("True")
14 plt.show()
```



Best Validation Accuracy: 90.505%

Classification Report:

	precision	recall	f1-score	support
0	0.91	0.94	0.93	374
1	0.89	0.89	0.89	469
2	0.89	0.88	0.88	441
3	0.96	0.97	0.97	481
4	0.99	1.00	0.99	597
5	0.81	0.75	0.78	369
6	0.81	0.84	0.82	376
accuracy			0.91	3107
macro avg	0.90	0.90	0.90	3107
weighted avg	0.90	0.91	0.90	3107

Accuracy vs n\_estimatorsPlotConfusion Matrix for Best Model

```
1 accuracy_list = [90.827, 90.795, 90.795, 90.859, 90.795, 90.763, 90.795, 90.795, 90.763, 90.763]
2
3 plt.figure(figsize=(8,5))
4 plt.plot(n_estimator_values, accuracy_list, marker='o', color='teal', linewidth=2)
5 plt.title("Validation Accuracy vs n_estimators")
6 plt.xlabel("n_estimators")
7 plt.ylabel("Validation Accuracy (%)")
8 plt.grid(True, linestyle='--', alpha=0.6)
9
10 best_idx = accuracy_list.index(max(accuracy_list))
11 plt.scatter(n_estimator_values[best_idx], accuracy_list[best_idx],
12            color='red', s=100, label=f"Best ({n_estimator_values[best_idx]})")
13 plt.legend()
14 plt.show()
15
```

