

✓ *Importing necessary Libraries*

```
1 # Import NumPy, which can deal with multi-dimensional arrays such as matrix
  intuitively.
2 import numpy as np # A useful package for dealing with mathematical processes
3 import pandas as pd # A common package for viewing tabular data
4 import matplotlib.pyplot as plt # We will be using Matplotlib for our graphs
5 import seaborn as sns; sns.set() # for plot styling
6 from google.colab import drive
7 drive.mount('/content/drive')
8 import sklearn.datasets # We want to be able to access the sklearn datasets
  again
9 from sklearn.metrics import accuracy_score, precision_score, recall_score,
  f1_score, balanced_accuracy_score # required for evaluating classification
  models
10 from sklearn.preprocessing import StandardScaler # We will be using the inbuilt
  preprocessing functions sklearn provides
11 from sklearn.model_selection import train_test_split # A library that can
  automatically perform data splitting for us
```

➡ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", forc

✓ Data Exploration

```
1 # Load the dataset  
2 data = pd.read_csv('/content/drive/MyDrive/ML-Sem2/Coursework/COMP1801_Coursework_Datas
```

```
1 # Display the first few rows of the dataset
2 print(data.head())
```

```

⇒ Lifespan partType microstructure coolingRate quenchTime forgeTime \
0 1469.17 Nozzle equiGrain 13 3.84 6.47
1 1793.64 Block singleGrain 19 2.62 3.48
2 700.60 Blade equiGrain 28 0.76 1.34
3 1082.10 Nozzle colGrain 9 2.01 2.19
4 1838.83 Blade colGrain 16 4.13 3.87

HeatTreatTime Nickel% Iron% Cobalt% Chromium% smallDefects \
0 46.87 65.73 16.52 16.82 0.93 10
1 44.70 54.22 35.38 6.14 4.26 19
2 9.54 51.83 35.95 8.81 3.41 35
3 20.29 57.03 23.33 16.86 2.78 0
4 16.13 59.62 27.37 11.45 1.56 10

largeDefects sliverDefects seedLocation castType
0 0 0 Bottom Die
1 0 0 Bottom Investment
2 3 0 Bottom Investment
3 1 0 Top Continuous
4 0 0 Top Die

```

```
1 # Check the structure and data types
2 print(data.info())
```

```
↵ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Lifespan               1000 non-null   float64
1   partType               1000 non-null   object
2   microstructure         1000 non-null   object
3   coolingRate             1000 non-null   int64
4   quenchTime              1000 non-null   float64
5   forgeTime              1000 non-null   float64
6   HeatTreatTime          1000 non-null   float64
7   Nickel%                1000 non-null   float64
8   Iron%                  1000 non-null   float64
9   Cobalt%                1000 non-null   float64
10  Chromium%              1000 non-null   float64
11  smallDefects            1000 non-null   int64
12  largeDefects            1000 non-null   int64
13  sliverDefects           1000 non-null   int64
14  seedLocation            1000 non-null   object
15  castType                1000 non-null   object
dtypes: float64(8), int64(4), object(4)
memory usage: 125.1+ KB
None
```

```
1 # Get summary statistics for numerical columns
2 print(data.describe())
```

```

↳
count    Lifespan    coolingRate    quenchTime    forgeTime    HeatTreatTime  \
mean    1298.556320    17.639000    2.764230    5.464600    30.194510
std     340.071434    7.491783    1.316979    2.604513    16.889415
min     417.990000    5.000000    0.500000    1.030000    1.030000
25%    1047.257500    11.000000    1.640000    3.170000    16.185000
50%    1266.040000    18.000000    2.755000    5.475000    29.365000
75%    1563.050000    24.000000    3.970000    7.740000    44.955000
max    2134.530000    30.000000    4.990000    10.000000    59.910000

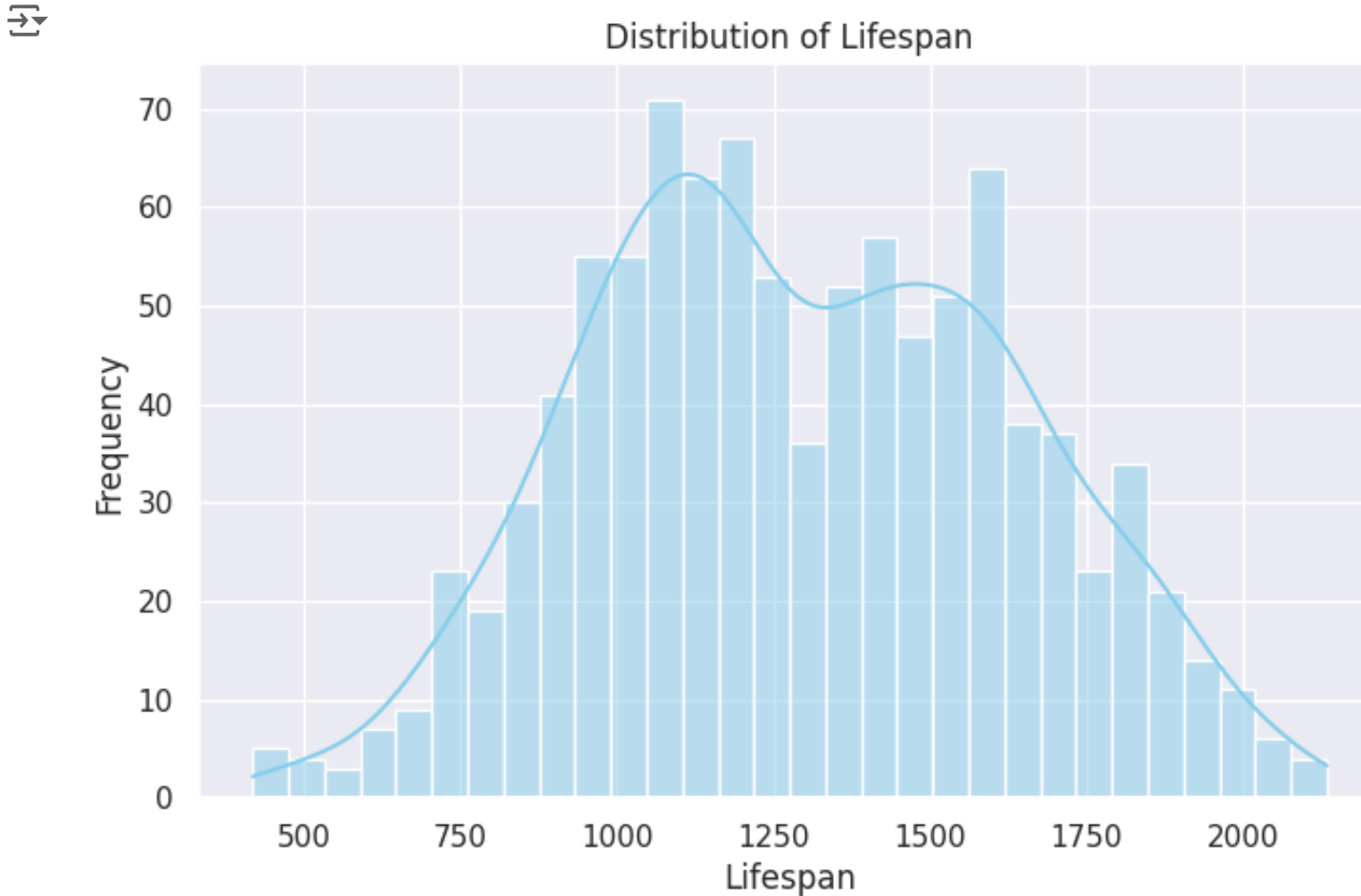
count    Nickel%    Iron%    Cobalt%    Chromium%    smallDefects  \
mean     60.243080    24.553580    12.434690    2.768650    17.311000
std      5.790475    7.371737    4.333197    1.326496    12.268365
min      50.020000    6.660000    5.020000    0.510000    0.000000
25%      55.287500    19.387500    8.597500    1.590000    7.000000
50%      60.615000    24.690000    12.585000    2.865000    18.000000
75%      65.220000    29.882500    16.080000    3.922500    26.000000
max      69.950000    43.650000    19.990000    4.990000    61.000000

count    largeDefects    sliverDefects
mean      0.550000      0.292000
std       1.163982      1.199239
min       0.000000      0.000000
25%       0.000000      0.000000
50%       0.000000      0.000000
75%       0.000000      0.000000
max       4.000000      8.000000

```

EDA

```
1 # Plot the distribution of the target variable (lifespan)
2 plt.figure(figsize=(8, 5))
3 sns.histplot(data['Lifespan'], kde=True, bins=30, color='skyblue')
4 plt.title('Distribution of Lifespan')
5 plt.xlabel('Lifespan')
6 plt.ylabel('Frequency')
7 plt.show()
```

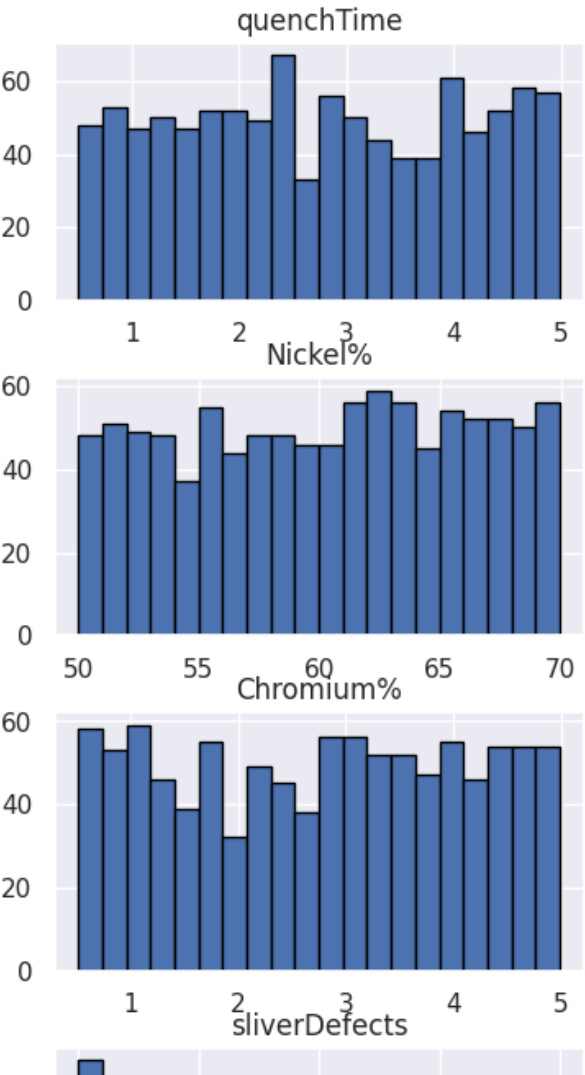
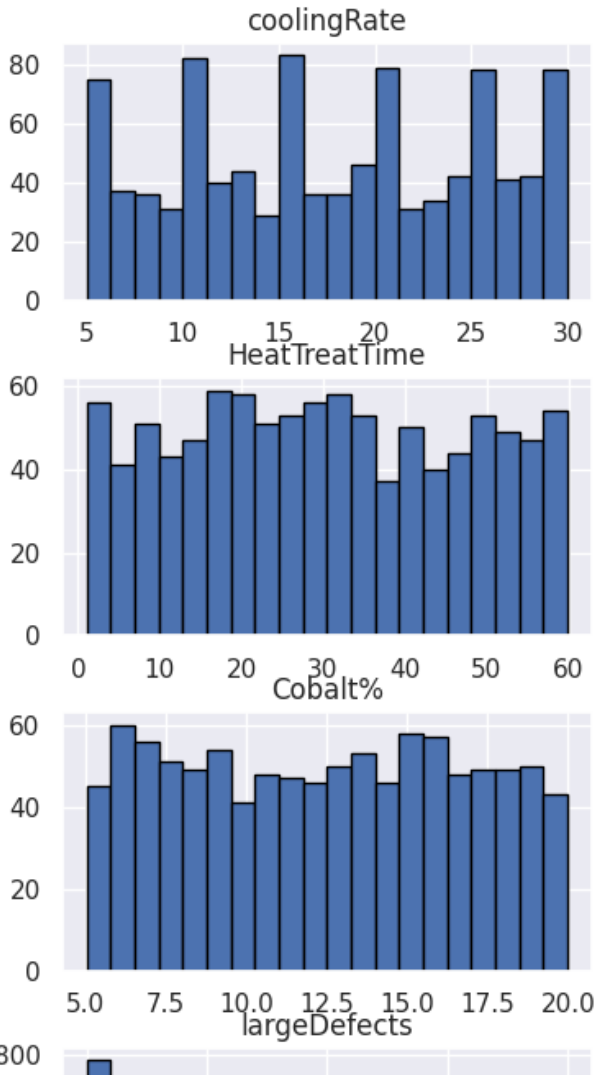
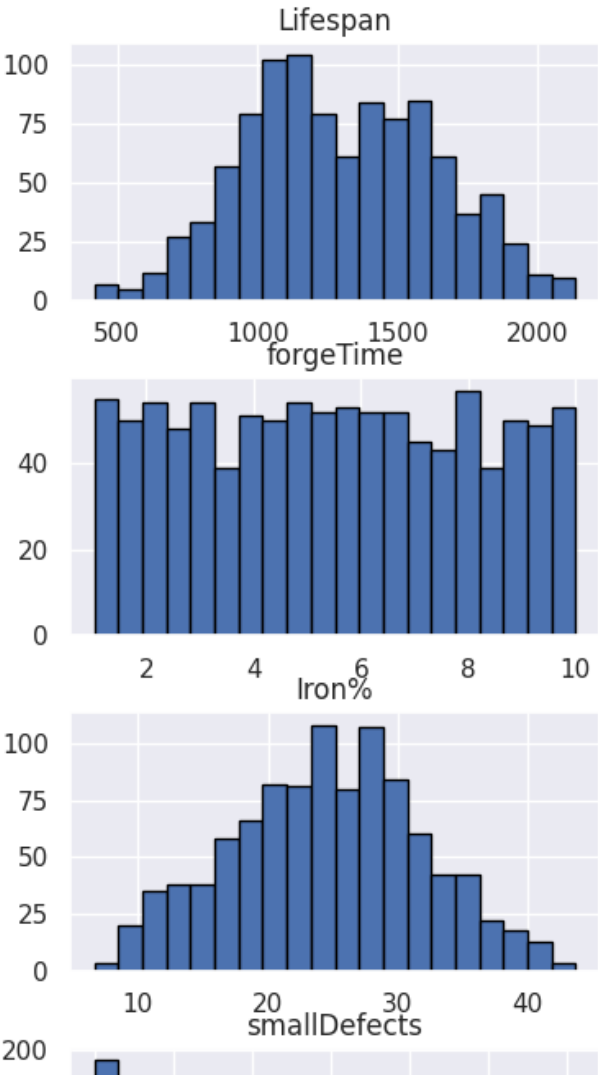


```
1 # Plot histograms for features
2 data.hist(figsize=(15, 10), bins=20, edgecolor='black')
3 plt.suptitle('Histograms of Features')
```

4 plt.show()



Histograms of Features



The histograms provide insights into the distribution of various features affecting product lifespan. The lifespan distribution is right-skewed, indicating a majority of shorter lifespans. Other features like cooling rate, quench time, and material composition seem to have a more uniform distribution. Understanding these distributions can help identify factors influencing product lifespan and guide potential improvements.

0 10 20 30 40 50 60 0 1 2 3 4 0 2 4 6 8

Exploring Categorical Features

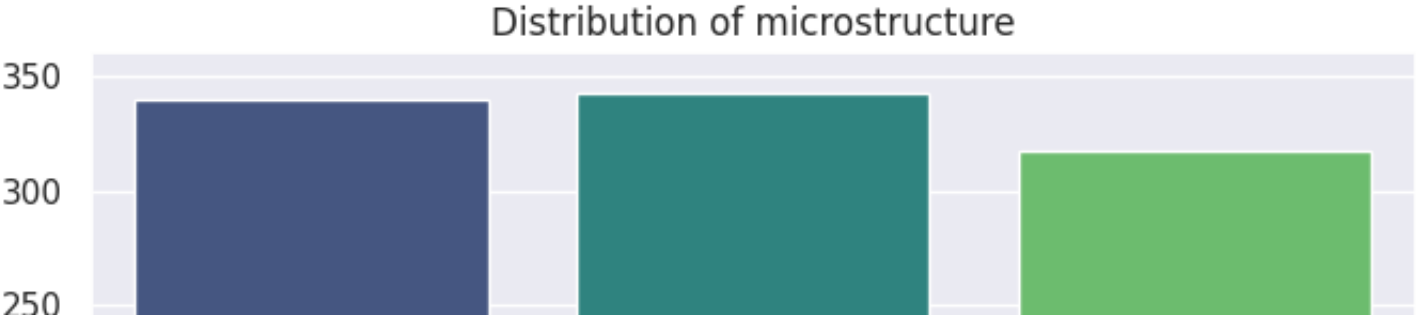
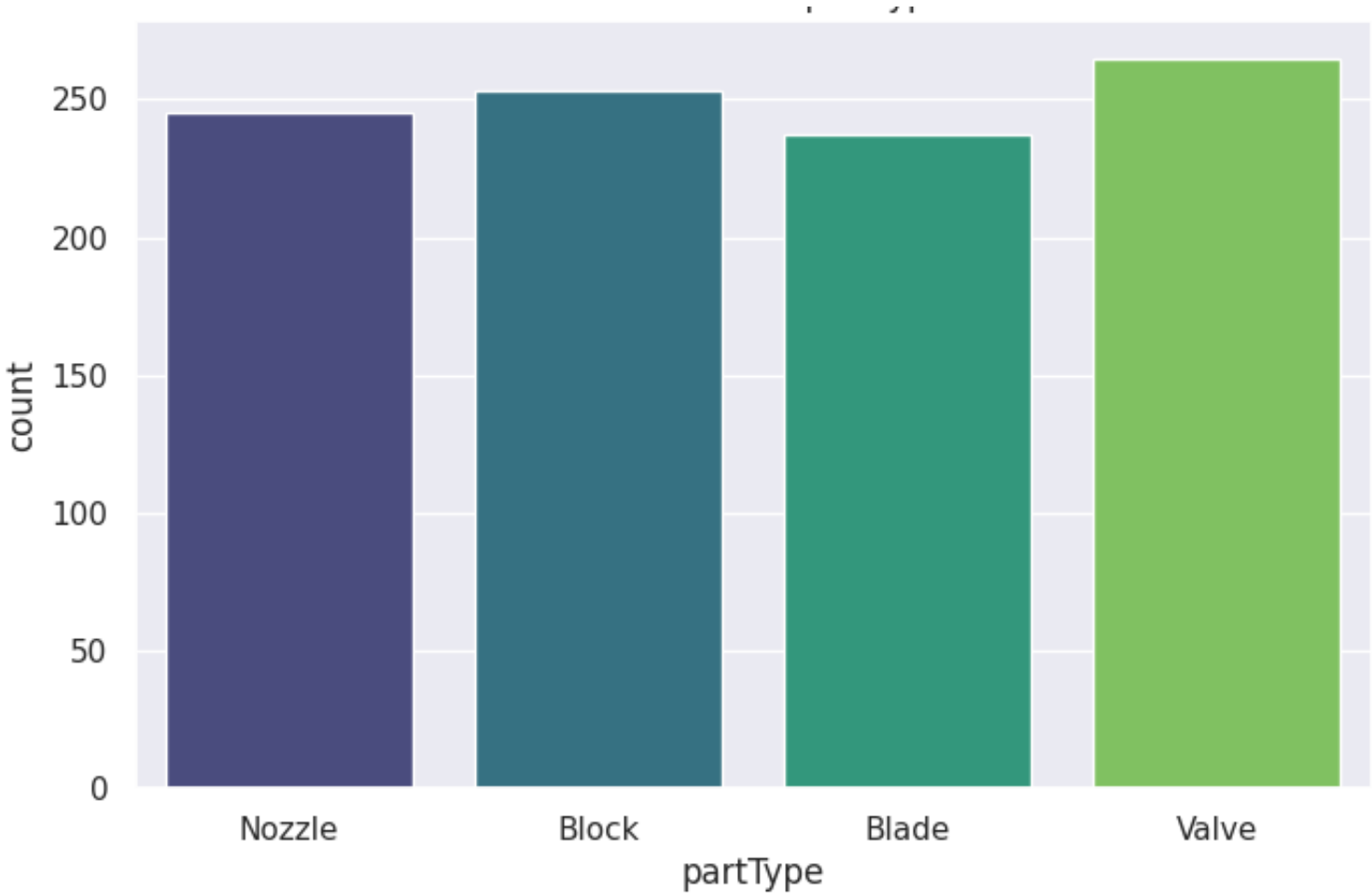
```
1 # Get unique values for categorical features
2 categorical_features = data.select_dtypes(include=['object']).columns.tolist()
3
4 for feature in categorical_features:
5     print(f"{feature}'s Unique Values List: {data[feature].unique()}")
```

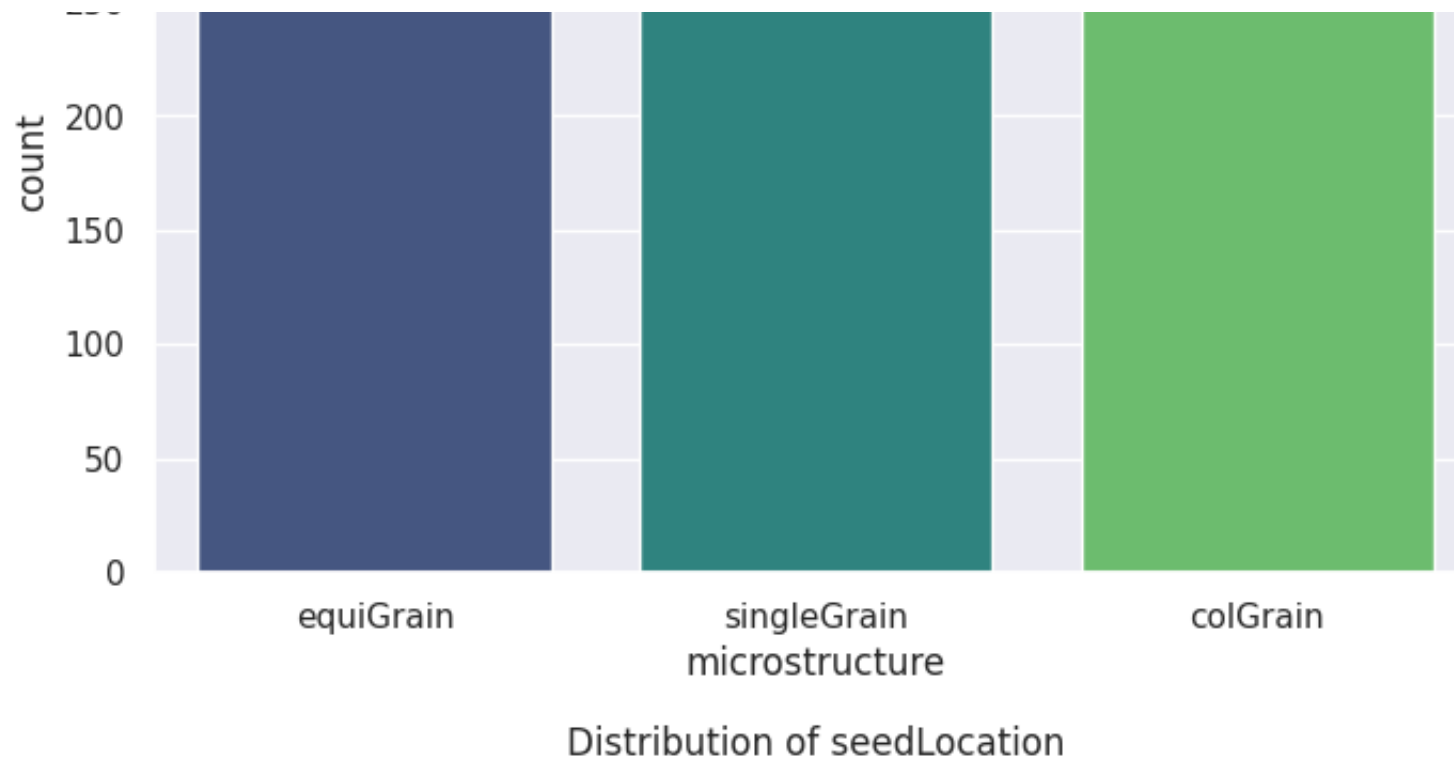
```
➞ partType's Unique Values List: ['Nozzle' 'Block' 'Blade' 'Valve']
   microstructure's Unique Values List: ['equiGrain' 'singleGrain' 'colGrain']
   seedLocation's Unique Values List: ['Bottom' 'Top']
   castType's Unique Values List: ['Die' 'Investment' 'Continuous']
```

```
1 # Visualize distributions of categorical features
2 for feature in categorical_features:
3     plt.figure(figsize=(8, 5))
4     sns.countplot(x=feature, data=data, hue=feature, palette='viridis')
5     plt.title(f'Distribution of {feature}')
6     plt.show()
```

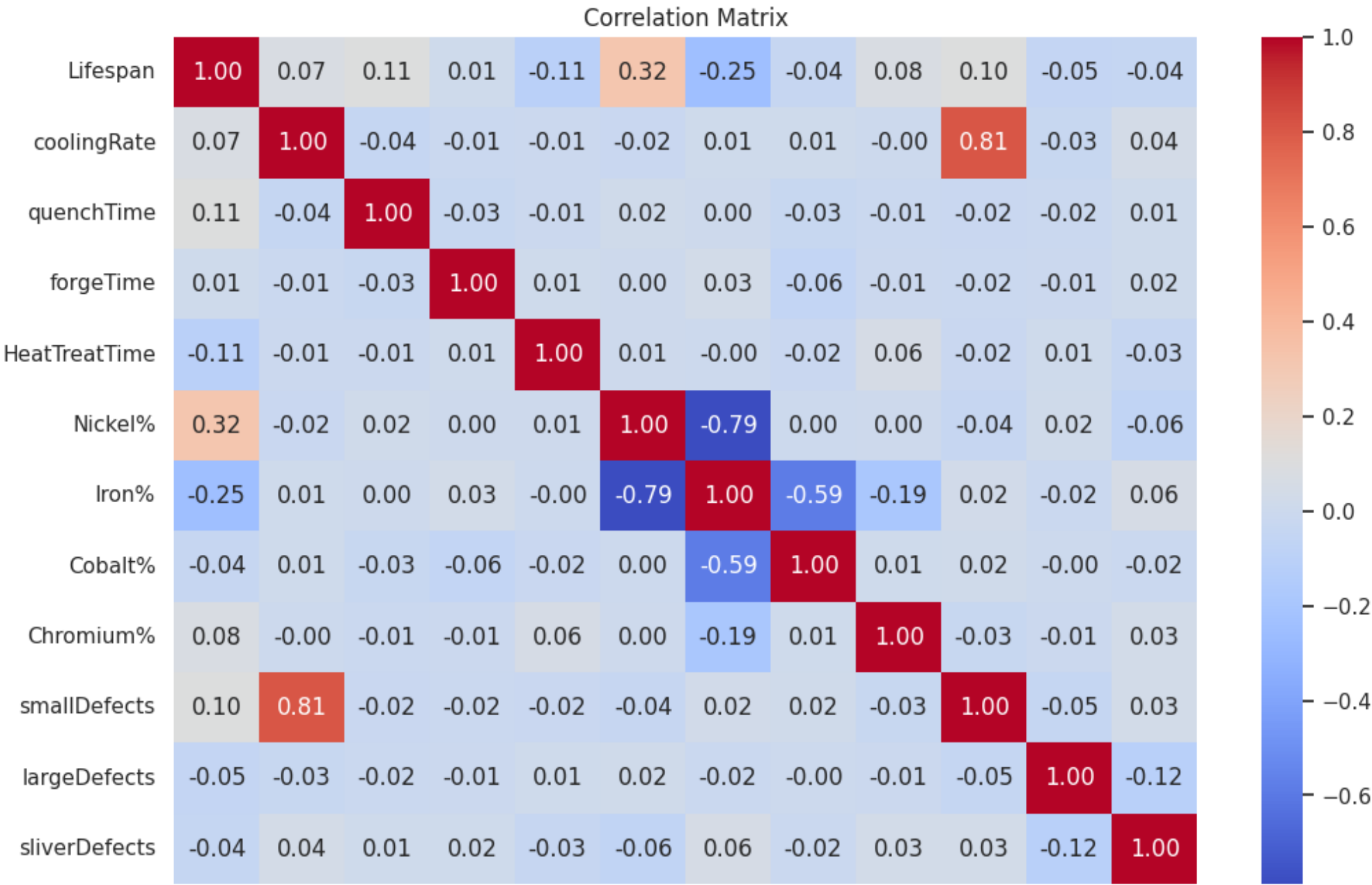


Distribution of partType





```
1 # Correlation heatmap for numerical features
2
3 datanumeric = data.select_dtypes(include=['int64', 'float64'])
4 plt.figure(figsize=(12, 8))
5
6 sns.heatmap(datanumeric.corr(), annot=True, fmt='.2f', cmap='coolwarm')
7 plt.title('Correlation Matrix')
8 plt.show()
9
```



```
1 # Identify categorical and numerical features
2 categorical_features = data.select_dtypes(include=['object']).columns.tolist() # Non-n
3 numerical_features = data.select_dtypes(include=['number']).columns.tolist() # Numeric
4
5 # Display the identified features
6 print("Categorical Features:", categorical_features)
7 print("Numerical Features:", numerical_features)
8
```

⇒ Categorical Features: ['partType', 'microstructure', 'seedLocation', 'castType']
Numerical Features: ['Lifespan', 'coolingRate', 'quenchTime', 'forgeTime', 'HeatTreatTime', 'Nickel%', 'Iron%',

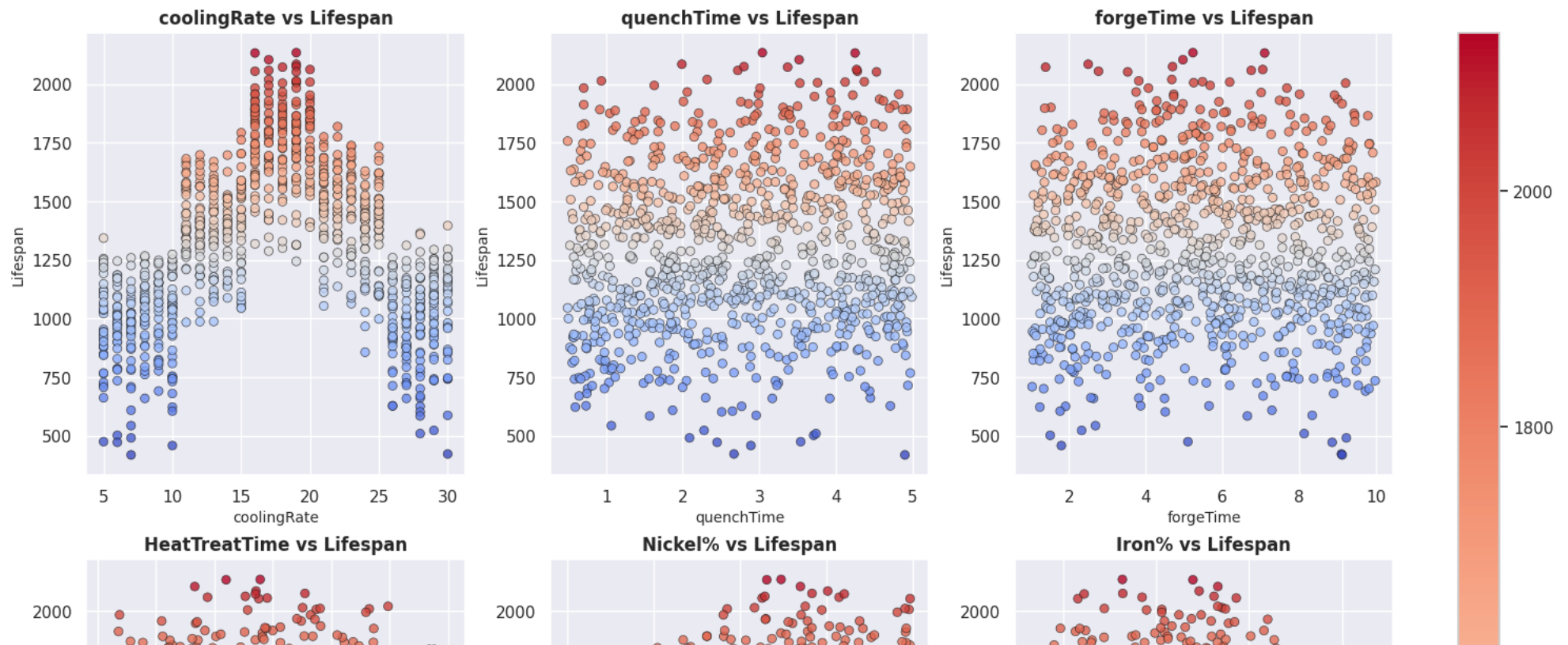
```
1 #Scatter plot numerical features vs lifespan
2 import math
3
4 # Get the numerical features except for the target variable
5 numerical_features = [col for col in data.select_dtypes(include=['float64', 'int64']).c
6
7 # Calculate number of rows and columns for the grid
8 num_features = len(numerical_features)
9 cols = 3 # Number of plots per row
10 rows = math.ceil(num_features / cols) # Number of rows needed
11
12 # Set the figure size and grid layout
13 fig, axes = plt.subplots(rows, cols, figsize=(15, 5 * rows), constrained_layout=True)
14
```

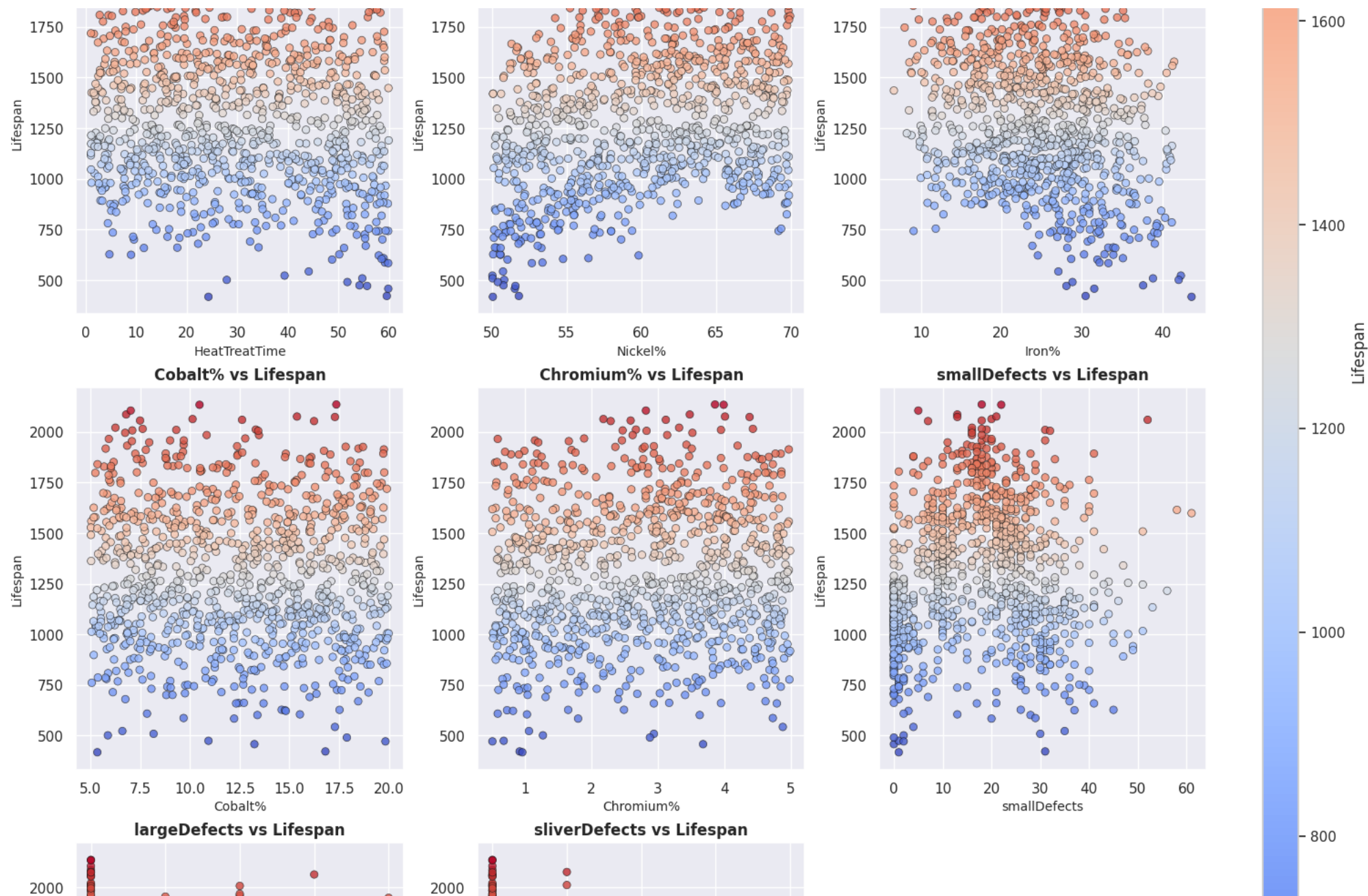
```
15 # Flatten the axes array for easier indexing
16 axes = axes.flatten()
17
18 # Set Seaborn theme
19 sns.set_theme(style="whitegrid")
20
21 # Loop through each feature and plot
22 for idx, feature in enumerate(numerical_features):
23     scatter = sns.scatterplot(
24         ax=axes[idx],
25         x=data[feature],
26         y=data['Lifespan'],
27         hue=data['Lifespan'], # Use 'Lifespan' for coloring
28         palette='coolwarm',
29         alpha=0.8,
30         edgecolor='k',
31         legend=False
32     )
33
34     # Add title and axis labels
35     axes[idx].set_title(f'{feature} vs Lifespan', fontsize=12, weight='bold')
36     axes[idx].set_xlabel(feature, fontsize=10)
37     axes[idx].set_ylabel('Lifespan', fontsize=10)
38
39 # Hide any unused subplots
40 for ax in axes[len(numerical_features):]:
```

```

41     ax.axis('off')
42
43 # Add a colorbar to the figure
44 sm = plt.cm.ScalarMappable(cmap='coolwarm', norm=plt.Normalize(vmin=data['Lifespan'].mi
45 cbar = fig.colorbar(sm, ax=axes, location='right', aspect=50, pad=0.05)
46 cbar.set_label('Lifespan', fontsize=12)
47
48 plt.show()
49

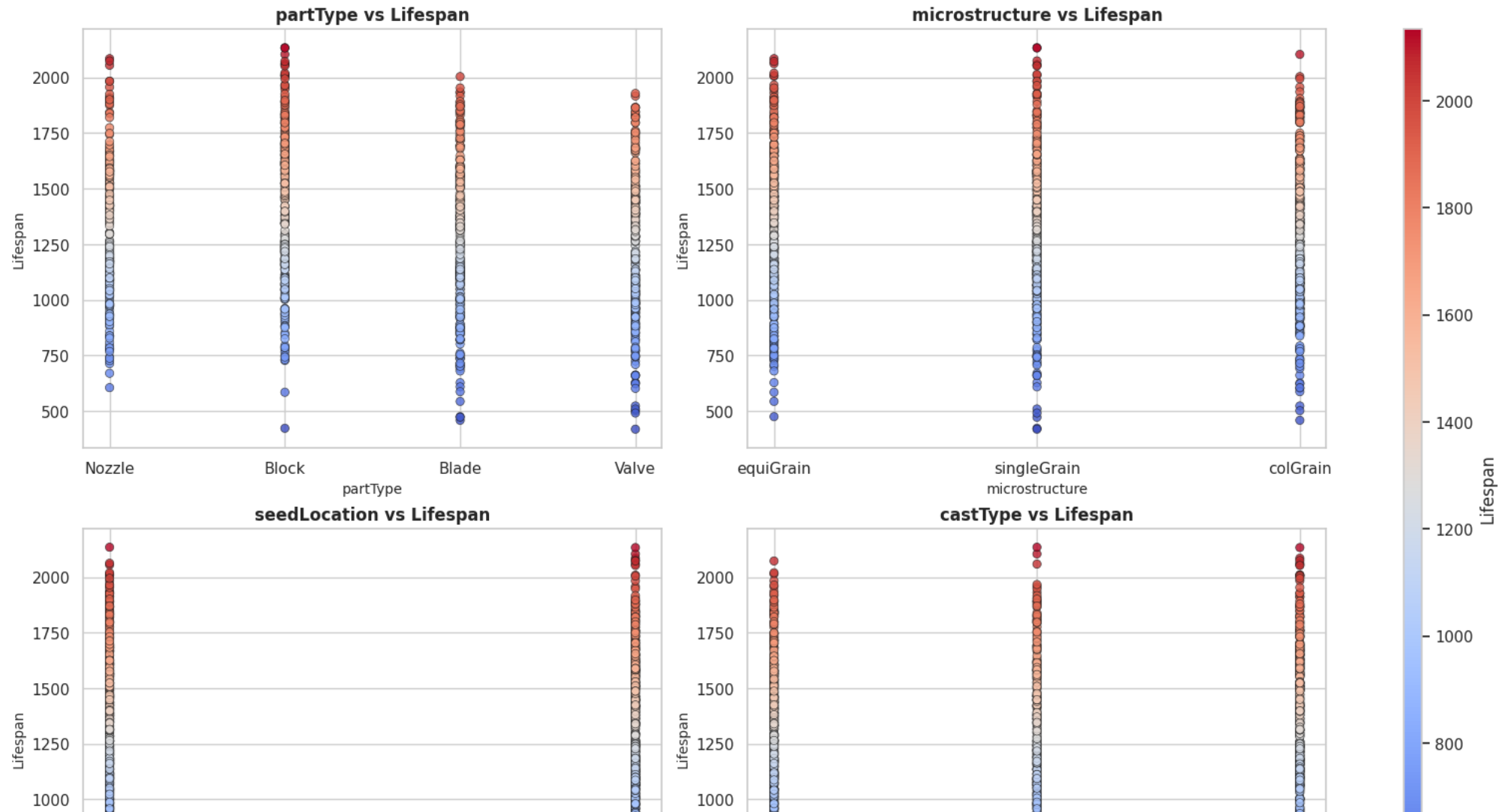
```






```
1 import math
2
3 # Get the numerical features except for the target variable
4 numerical_features = [col for col in data.select_dtypes(include=['float64', 'int64']).c
5
6 # Calculate number of rows and columns for the grid
7 num_features = len(categorical_features)
8 cols = 2 # Number of plots per row
9 rows = math.ceil(num_features / cols) # Number of rows needed
10
11 # Set the figure size and grid layout
12 fig, axes = plt.subplots(rows, cols, figsize=(15, 5 * rows), constrained_layout=True)
13
14 # Flatten the axes array for easier indexing
15 axes = axes.flatten()
16
17 # Set Seaborn theme
18 sns.set_theme(style="whitegrid")
19
20 # Loop through each feature and plot
21 for idx, feature in enumerate(categorical_features):
22     scatter = sns.scatterplot(
23         ax=axes[idx],
24         x=data[feature],
25         y=data['Lifespan'],
26         hue=data['Lifespan'], # Use 'Lifespan' for coloring
```

```
27     palette='coolwarm',
28     alpha=0.8,
29     edgecolor='k',
30     legend=False
31 )
32
33 # Add title and axis labels
34 axes[idx].set_title(f'{feature} vs Lifespan', fontsize=12, weight='bold')
35 axes[idx].set_xlabel(feature, fontsize=10)
36 axes[idx].set_ylabel('Lifespan', fontsize=10)
37
38 # Hide any unused subplots
39 for ax in axes[len(categorical_features):]:
40     ax.axis('off')
41
42 # Add a colorbar to the figure
43 sm = plt.cm.ScalarMappable(cmap='coolwarm', norm=plt.Normalize(vmin=data['Lifespan'].min(), vmax=data['Lifespan'].max()))
44 cbar = fig.colorbar(sm, ax=axes, location='right', aspect=50, pad=0.05)
45 cbar.set_label('Lifespan', fontsize=12)
46
47 plt.show()
48
```



```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler, OneHotEncoder
3 from sklearn.compose import ColumnTransformer
```

```
4
5 # Define the target column and features
6 target_column = 'Lifespan'
7 X = data.drop(columns=[target_column])
8 y = data[target_column]
9
10 # Automatically detect categorical and numerical features
11 categorical_features = X.select_dtypes(include=['object']).columns.tolist()
12 numerical_features = X.select_dtypes(include=['number']).columns.tolist()
13
14 # Print detected features
15 print("Categorical Features:", categorical_features)
16 print("Numerical Features:", numerical_features)
17
18 # Split the dataset into training and testing sets (80% training, 20% testing)
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=4
20
21 # Define a preprocessing pipeline
22 preprocessor = ColumnTransformer(
23     transformers=[
24         ('num', StandardScaler(), numerical_features), # Scale numerical features
25         ('cat', OneHotEncoder(drop='first'), categorical_features) # One-hot encode ca
26     ]
27 )
28
29 # Apply preprocessing to the training and testing sets
```

```
30 X_train = preprocessor.fit_transform(X_train)
31 X_test = preprocessor.transform(X_test)
32
33 # Print the shapes of the processed data
34 print(f"Shape of X_train: {X_train.shape}")
35 print(f"Shape of X_test: {X_test.shape}")
36 print(f"Shape of y_train: {y_train.shape}")
37 print(f"Shape of y_test: {y_test.shape}")
38
```

⇒ Categorical Features: ['partType', 'microstructure', 'seedLocation', 'castType']
Numerical Features: ['coolingRate', 'quenchTime', 'forgeTime', 'HeatTreatTime', 'Nickel%', 'Iron%', 'Cobalt%', '']
Shape of X_train: (800, 19)
Shape of X_test: (200, 19)
Shape of y_train: (800,)
Shape of y_test: (200,)

```
1 from sklearn.linear_model import Ridge
2 from sklearn.model_selection import RandomizedSearchCV
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score, explained_variance_score
5 import numpy as np
6 from scipy.stats import uniform
7
8 # Scaling the data (important for Ridge regression)
9 scaler = StandardScaler()
10 X_train_scaled = scaler.fit_transform(X_train)
```

```
11 X_test_scaled = scaler.transform(X_test)
12
13 # Ridge Regression Model
14 ridge = Ridge()
15
16 # Hyperparameter tuning using RandomizedSearchCV
17 # Define a wider distribution for alpha (log scale would often be more appropriate, but
18 ridge_params = {
19     'alpha': uniform(0.1, 100) # Randomized search over alpha in the range [0.1, 100]
20 }
21
22 # Perform RandomizedSearchCV
23 random_search = RandomizedSearchCV(estimator=ridge, param_distributions=ridge_params, n
24
25 # Fit the model
26 random_search.fit(X_train_scaled, y_train)
27
28 # Best hyperparameters found by RandomizedSearchCV
29 print("Best Parameters:", random_search.best_params_)
30
31 # Best model from RandomizedSearchCV
32 best_ridge_model = random_search.best_estimator_
33
34 # Predictions using the best model
35 y_train_pred = best_ridge_model.predict(X_train_scaled)
36 y_test_pred = best_ridge_model.predict(X_test_scaled)
```

```
37
38 # Evaluate the model on training and testing data
39 train_rmse = np.sqrt(mean_squared_error(y_train, y_train_pred))
40 test_rmse = np.sqrt(mean_squared_error(y_test, y_test_pred))
41 train_mae = mean_absolute_error(y_train, y_train_pred)
42 test_mae = mean_absolute_error(y_test, y_test_pred)
43 train_r2 = r2_score(y_train, y_train_pred)
44 test_r2 = r2_score(y_test, y_test_pred)
45 train_explained_variance = explained_variance_score(y_train, y_train_pred)
46 test_explained_variance = explained_variance_score(y_test, y_test_pred)
47
48 # Print evaluation metrics
49 print(f"Train RMSE: {train_rmse}")
50 print(f"Test RMSE: {test_rmse}")
51 print(f"Train MAE: {train_mae}")
52 print(f"Test MAE: {test_mae}")
53 print(f"Train R2: {train_r2}")
54 print(f"Test R2: {test_r2}")
55 print(f"Train Explained Variance: {train_explained_variance}")
56 print(f"Test Explained Variance: {test_explained_variance}")
57
```

⇒ Best Parameters: {'alpha': 66.35222843539819}
Train RMSE: 306.4802970768733
Test RMSE: 300.01997922693323
Train MAE: 261.4550108061923
Test MAE: 253.21894063540793
Train R²: 0.20718947746556693
Test R²: 0.13081284044543373
Train Explained Variance: 0.20718947746556693
Test Explained Variance: 0.1377299709215305

```
1 ridge_grid = RandomizedSearchCV(ridge, ridge_params, cv=5, scoring='neg_mean_squared_er
2 ridge_grid.fit(X_train_scaled, y_train)
```

⇒

- ▶ **RandomizedSearchCV** ⓘ ?
 - ▶ **best_estimator_**: Ridge
 - ▶ Ridge ?


```
1 # Best Ridge Regression model and parameters
2 best_ridge = ridge_grid.best_estimator_
3 print("Best Ridge Parameters:", ridge_grid.best_params_)
4
5 # Evaluate on test set
6 ridge_preds = best_ridge.predict(X_test_scaled)
7
8 # Calculate evaluation metrics
```

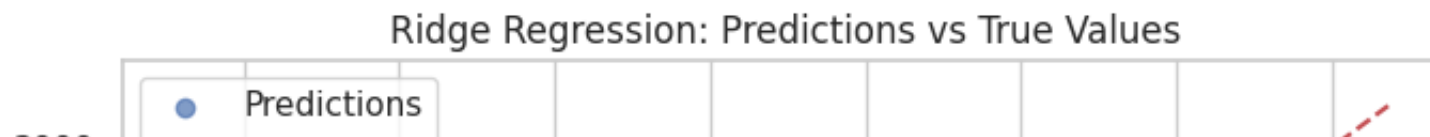


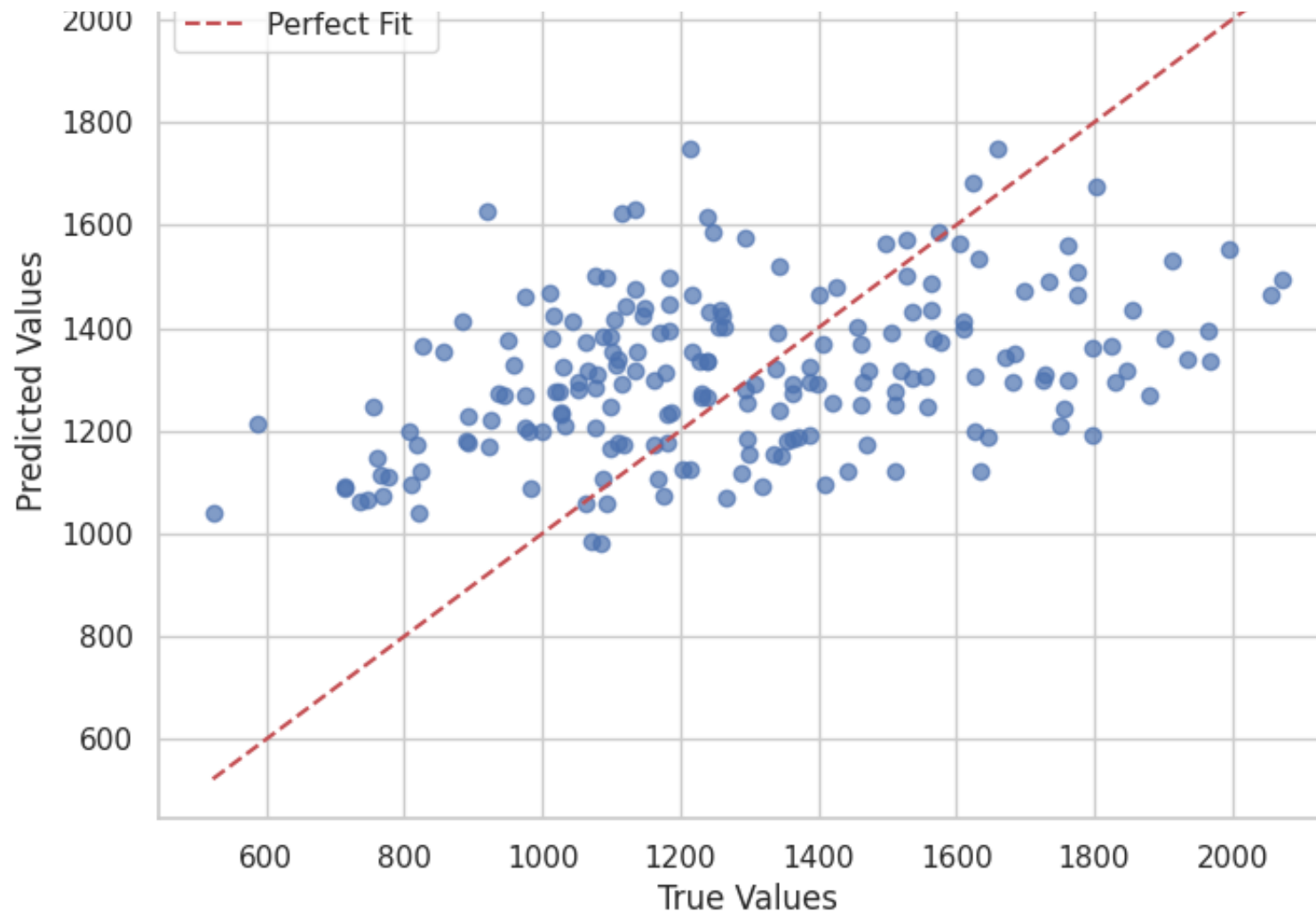
```

9 ridge_rmse = np.sqrt(mean_squared_error(y_test, ridge_preds))
10 ridge_mae = mean_absolute_error(y_test, ridge_preds)
11 ridge_r2 = r2_score(y_test, ridge_preds)
12 ridge_explained_variance = explained_variance_score(y_test, ridge_preds)
13
14 # Print evaluation metrics
15 print(f"Ridge Regression Test RMSE: {ridge_rmse:.2f}")
16 print(f"Ridge Regression Test MAE: {ridge_mae:.2f}")
17 print(f"Ridge Regression Test R²: {ridge_r2:.2f}")
18 print(f"Ridge Regression Test Explained Variance: {ridge_explained_variance:.2f}")
19 # Ridge Regression: Predictions vs True Values
20 plt.figure(figsize=(8, 6))
21 plt.scatter(y_test, ridge_preds, alpha=0.7, label='Predictions')
22 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--', label='Perf')
23 plt.title('Ridge Regression: Predictions vs True Values')
24 plt.xlabel('True Values')
25 plt.ylabel('Predicted Values')
26 plt.legend()
27 plt.show()

```

 Best Ridge Parameters: {'alpha': 62.36876997566669}
 Ridge Regression Test RMSE: 300.09
 Ridge Regression Test MAE: 253.30
 Ridge Regression Test R²: 0.13
 Ridge Regression Test Explained Variance: 0.14





```
1 import xgboost as xgb
2 from sklearn.model_selection import RandomizedSearchCV
3 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score, explain
```

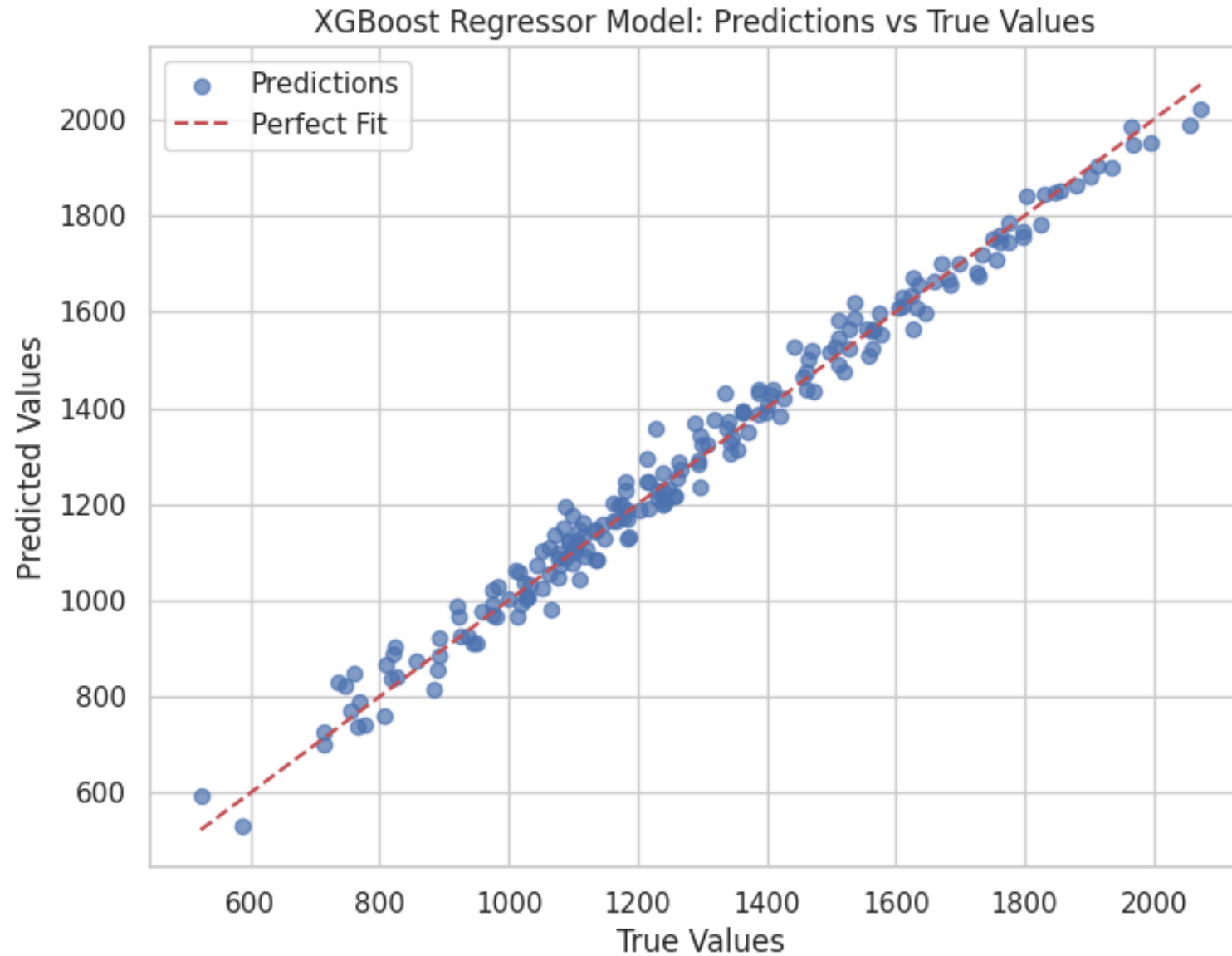
```
4 import numpy as np
5
6 # XGBoost Regressor Model
7 xgb_model = xgb.XGBRegressor(random_state=42)
8
9 # Reduced hyperparameter tuning using RandomizedSearchCV for faster results
10 xgb_params = {
11     'n_estimators': [100, 200, 300], # Limited to a smaller number of trees for faster
12     'learning_rate': [0.01, 0.05, 0.1], # A limited range of learning rates
13     'max_depth': [3, 5, 7], # Exploring only three tree depths
14     'min_child_weight': [1, 3], # Limited range to avoid too many combinations
15     'subsample': [0.7, 0.8, 1.0], # Testing three values for subsample
16     'colsample_bytree': [0.7, 0.8, 1.0], # Testing three values for colsample_bytree
17     'gamma': [0, 0.1], # Reduced range of gamma values for faster testing
18     'reg_alpha': [0, 0.1], # L2 regularization with fewer values
19     'reg_lambda': [0, 0.1] # L1 regularization with fewer values
20 }
21
22 # RandomizedSearchCV for faster hyperparameter tuning
23 xgb_random = RandomizedSearchCV(xgb_model, xgb_params, cv=5, scoring='neg_mean_squared_
24                                 n_iter=10, verbose=2, random_state=42, n_jobs=-1)
25 xgb_random.fit(X_train, y_train)
26
27 # Best model
28 best_xgb = xgb_random.best_estimator_
29 print("Best XGBoost Parameters:", xgb_random.best_params_)
```

```
30
31 # Evaluate on test set
32 xgb_preds = best_xgb.predict(X_test)
33 xgb_rmse = np.sqrt(mean_squared_error(y_test, xgb_preds))
34 xgb_mae = mean_absolute_error(y_test, xgb_preds)
35 xgb_r2 = r2_score(y_test, xgb_preds)
36 xgb_explained_variance = explained_variance_score(y_test, xgb_preds)
37
38 print(f"XGBoost Regressor Test RMSE: {xgb_rmse:.2f}")
39 print(f"XGBoost Regressor Test MAE: {xgb_mae:.2f}")
40 print(f"XGBoost Regressor Test R2: {xgb_r2:.2f}")
41 print(f"XGBoost Regressor Test Explained Variance: {xgb_explained_variance:.2f}")
42
```

➦ Fitting 5 folds for each of 10 candidates, totalling 50 fits
Best XGBoost Parameters: {'subsample': 1.0, 'reg_lambda': 0, 'reg_alpha': 0.1, 'n_estimators': 300, 'min_child_w
XGBoost Regressor Test RMSE: 39.28
XGBoost Regressor Test MAE: 31.09
XGBoost Regressor Test R²: 0.99
XGBoost Regressor Test Explained Variance: 0.99

```
1 # XGBoost Regressor Model: Predictions vs True Values
2 plt.figure(figsize=(8, 6))
3 plt.scatter(y_test, xgb_preds, alpha=0.7, label='Predictions')
4 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--', label='Perf
5 plt.title(' XGBoost Regressor Model: Predictions vs True Values')
6 plt.xlabel('True Values')
```

```
7 plt.ylabel('Predicted Values')  
8 plt.legend()  
9 plt.show()
```



```
1 from sklearn.ensemble import RandomForestRegressor
```

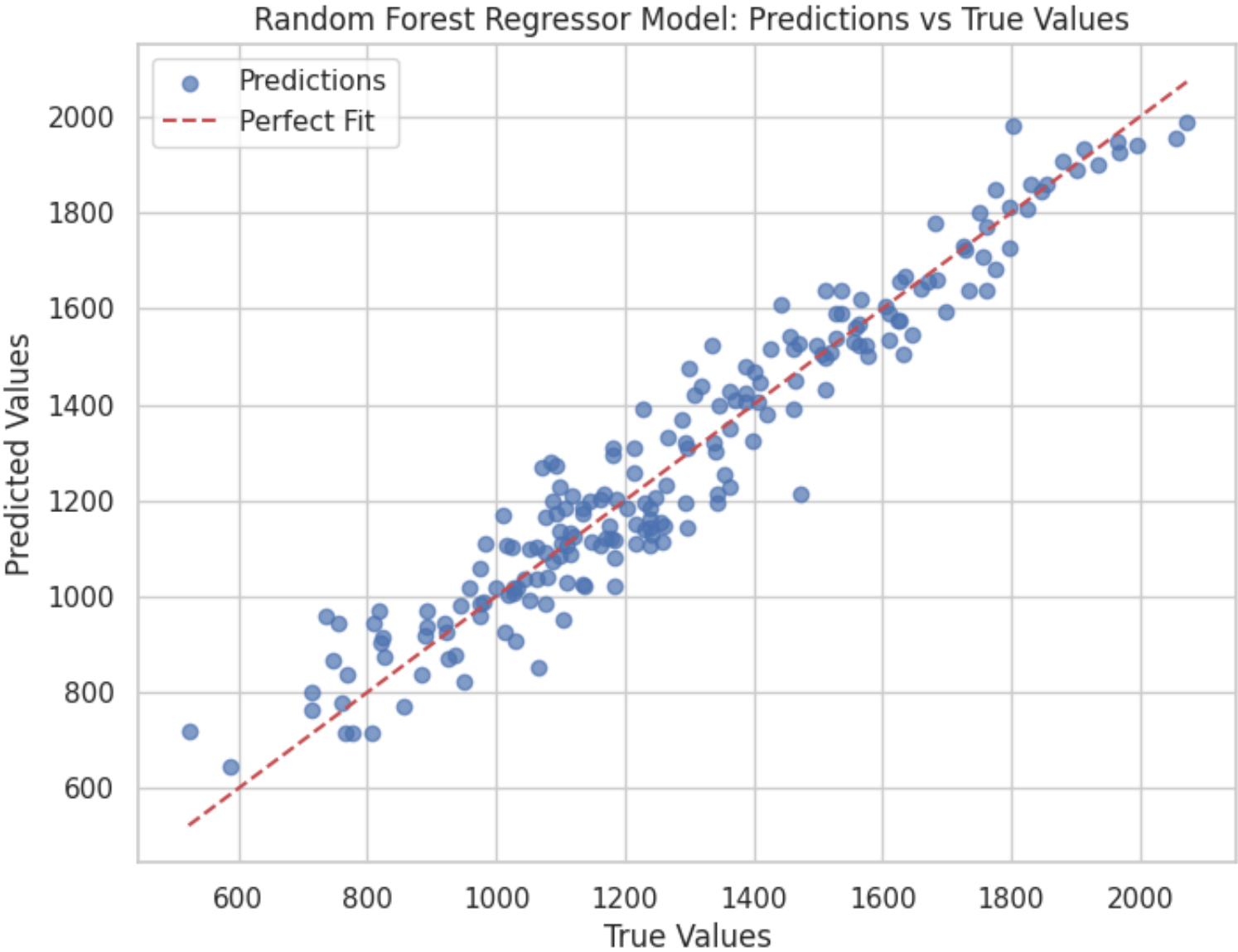
```
1 from sklearn.ensemble import RandomForestRegressor
2 from sklearn.model_selection import RandomizedSearchCV, KFold
3 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score, explained_variance_score
4 import numpy as np
5
6 # Random Forest Regressor Model
7 rf = RandomForestRegressor(random_state=42)
8
9 # Hyperparameter tuning using RandomizedSearchCV
10 rf_params = {
11     'n_estimators': [50, 100, 200, 300, 500], # More options for number of estimators
12     'max_depth': [None, 10, 20, 30, 40, 50], # Expand depth range
13     'min_samples_split': [2, 5, 10, 15, 20], # More options for splitting
14     'min_samples_leaf': [1, 2, 4, 5, 10], # More options for leaf size
15     'max_features': ['sqrt', 'log2', None], # More options for max features
16     'bootstrap': [True, False],
17     'warm_start': [True, False], # Allow reuse of previous trees to increase estimator
18     'random_state': [42] # For reproducibility
19 }
20
21 # Implementing RandomizedSearchCV with KFold to ensure proper cross-validation for regr
22 kf = KFold(n_splits=5, shuffle=True, random_state=42) # Using KFold for regression tas
23
24 rf_random_search = RandomizedSearchCV(
25     estimator=rf,
26     param_distributions=rf_params, # Randomly sample from these parameters
27     n_iter=100, # Number of parameter combinations to test (increased for better explo
```

```
28     cv=kf,          # Cross-validation strategy
29     scoring='neg_mean_squared_error', # Use negative MSE for regression
30     n_jobs=-1,      # Use all available CPUs for faster processing
31     verbose=2,      # Provide more detailed output during fitting
32     random_state=42,
33     error_score='raise'
34 )
35
36 # Fit the model with the training data
37 rf_random_search.fit(X_train, y_train)
38
39 # Best Random Forest model and parameters
40 best_rf = rf_random_search.best_estimator_
41 print("Best Random Forest Parameters:", rf_random_search.best_params_)
42
43 # Evaluate on test set
44 rf_preds = best_rf.predict(X_test)
45
46 # Evaluate using multiple metrics
47 rf_rmse = np.sqrt(mean_squared_error(y_test, rf_preds))
48 rf_mae = mean_absolute_error(y_test, rf_preds)
49 rf_r2 = r2_score(y_test, rf_preds)
50 rf_explained_variance = explained_variance_score(y_test, rf_preds)
51
52 # Print evaluation metrics
53 print(f"Random Forest Regressor Test RMSE: {rf_rmse:.2f}")
```



```
⇒ Fitting 5 folds for each of 100 candidates, totalling 500 fits
Best Random Forest Parameters: {'warm_start': True, 'random_state': 42, 'n_estimators': 500, 'min_samples_split': 10}
Random Forest Regressor Test RMSE: 86.13
Random Forest Regressor Test MAE: 68.48
Random Forest Regressor Test R2: 0.93
Random Forest Regressor Test Explained Variance: 0.93
```


```
1 # Random Forest Regressor Model: Predictions vs True Values
2 plt.figure(figsize=(8, 6))
3 plt.scatter(y_test, rf_preds, alpha=0.7, label='Predictions')
4 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--', label='Perf
5 plt.title(' Random Forest Regressor Model: Predictions vs True Values')
6 plt.xlabel('True Values')
7 plt.ylabel('Predicted Values')
8 plt.legend()
9 plt.show()
```



```

1 # Summarize results in a DataFrame for comparison
2 results = {
3     'Model': ['Ridge Regression', 'Random Forest Regressor', 'XGBoost Regressor'],
4     #'Train Score': [ridge_train_score, rf_train_score, xgb_train_score],
5     #'Test Score': [ridge_test_score, rf_test_score, xgb_test_score],
6     'RMSE': [ridge_rmse, rf_rmse, xgb_rmse],
7     'MAE': [ridge_mae, rf_mae, xgb_mae],
8     'R²': [ridge_r2, rf_r2, xgb_r2],
9     'Explained Variance': [ridge_explained_variance, rf_explained_variance, xgb_explaine
10 }
11
12 results_df = pd.DataFrame(results)
13 results_df
14

```



	Model	RMSE	MAE	R ²	Explained Variance
0	Ridge Regression	300.091482	253.303066	0.130398	0.137331
1	Random Forest Regressor	86.129146	68.478550	0.928367	0.928450
2	XGBoost Regressor	39.278167	31.091480	0.985102	0.985371

```

1 # Assuming 'data' is your DataFrame and 'Lifespan' is the target column
2 X = data.drop(columns=['Lifespan']) # Drop target column to get features
3 y = data['Lifespan'] # Target column
4

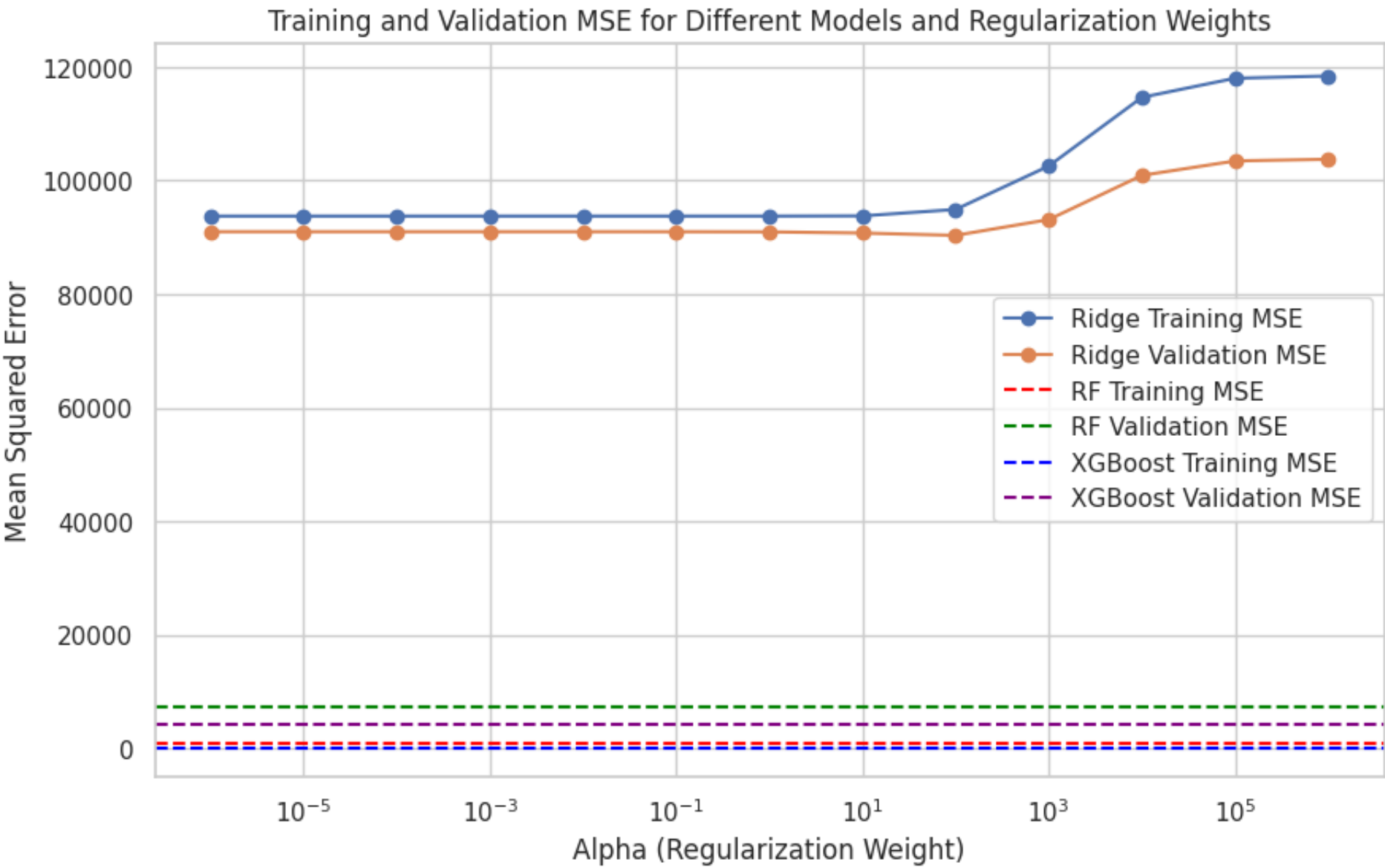
```

```
5 # Automatically detect categorical and numerical features
6 categorical_features = X.select_dtypes(include=['object']).columns.tolist()
7 numerical_features = X.select_dtypes(include=['number']).columns.tolist()
8
9 # Preprocessing for numerical and categorical features
10 preprocessor = ColumnTransformer(
11     transformers=[
12         ('num', StandardScaler(), numerical_features), # Scale numerical features
13         ('cat', OneHotEncoder(drop='first'), categorical_features) # One-hot encode ca
14     ]
15 )
16
17 # Split data into training and validation sets (80% training, 20% validation)
18 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)
19
20 # Apply preprocessing to the features
21 X_train = preprocessor.fit_transform(X_train)
22 X_val = preprocessor.transform(X_val)
23
24 # List of alpha values to test for Ridge Regression (logarithmic scale from 1e-6 to 1e6
25 alphas = np.logspace(-6, 6, 13)
26
27 # Initialize lists for storing errors
28 ridge_train_errors = []
29 ridge_val_errors = []
30 rf_train_errors = []
```

```
31 rf_val_errors = []
32 xgb_train_errors = []
33 xgb_val_errors = []
34
35 # Loop over alpha values for Ridge Regression
36 for alpha in alphas:
37     # Ridge regression model
38     ridge_model = Ridge(alpha=alpha)
39     ridge_model.fit(X_train, y_train)
40
41     # Predictions for training and validation sets
42     ridge_train_pred = ridge_model.predict(X_train)
43     ridge_val_pred = ridge_model.predict(X_val)
44
45     # Calculate MSE for training and validation
46     ridge_train_mse = mean_squared_error(y_train, ridge_train_pred)
47     ridge_val_mse = mean_squared_error(y_val, ridge_val_pred)
48
49     ridge_train_errors.append(ridge_train_mse)
50     ridge_val_errors.append(ridge_val_mse)
51
52 # Now let's also compare Random Forest and XGBoost regression models:
53 # Initialize and train Random Forest Regressor and XGBoost Regressor
54 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
55 xgb_model = xgb.XGBRegressor(objective='reg:squarederror', random_state=42)
56
```

```
57 # Train Random Forest Model
58 rf_model.fit(X_train, y_train)
59 rf_train_pred = rf_model.predict(X_train)
60 rf_val_pred = rf_model.predict(X_val)
61 rf_train_mse = mean_squared_error(y_train, rf_train_pred)
62 rf_val_mse = mean_squared_error(y_val, rf_val_pred)
63
64 rf_train_errors.append(rf_train_mse)
65 rf_val_errors.append(rf_val_mse)
66
67 # Train XGBoost Model
68 xgb_model.fit(X_train, y_train)
69 xgb_train_pred = xgb_model.predict(X_train)
70 xgb_val_pred = xgb_model.predict(X_val)
71 xgb_train_mse = mean_squared_error(y_train, xgb_train_pred)
72 xgb_val_mse = mean_squared_error(y_val, xgb_val_pred)
73
74 xgb_train_errors.append(xgb_train_mse)
75 xgb_val_errors.append(xgb_val_mse)
76
77 # Now plot the MSE results for all models
78 plt.figure(figsize=(10, 6))
79 plt.plot(alphas, ridge_train_errors, label='Ridge Training MSE', marker='o')
80 plt.plot(alphas, ridge_val_errors, label='Ridge Validation MSE', marker='o')
81 plt.axhline(rf_train_errors[0], color='red', linestyle='--', label='RF Training MSE')
82 plt.axhline(rf_val_errors[0], color='green', linestyle='--', label='RF Validation MSE')
```

```
83 plt.axhline(xgb_train_errors[0], color='blue', linestyle='--', label='XGBoost Training')
84 plt.axhline(xgb_val_errors[0], color='purple', linestyle='--', label='XGBoost Validation')
85 plt.xscale('log') # Use logarithmic scale for alpha
86 plt.xlabel('Alpha (Regularization Weight)')
87 plt.ylabel('Mean Squared Error')
88 plt.title('Training and Validation MSE for Different Models and Regularization Weights')
89 plt.legend()
90 plt.grid(True)
91 plt.show()
92
```



Explanation: The plot visualizes the training and validation Mean Squared Error (MSE) for three different machine learning models (Ridge Regression, Random Forest, and XGBoost) across a range of regularization weights (α). The plot illustrates the trade-off between underfitting and overfitting. XGBoost demonstrates the best performance with the lowest validation MSE, indicating a good balance between bias and variance. Ridge Regression also performs relatively well, while Random Forest tends to overfit the training data.

✓ Part 4: Classification Implementation

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder
4 from sklearn.compose import ColumnTransformer
5 from sklearn.cluster import KMeans
6 from sklearn.metrics import silhouette_score
7 from sklearn.model_selection import train_test_split, GridSearchCV
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier
10 from sklearn.metrics import accuracy_score, f1_score, classification_report, confusion_
11 from imblearn.over_sampling import SMOTE
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14
```

```

1 # Load the dataset
2 data = pd.read_csv('/content/COMP1801_Coursework_Dataset.csv')
3
4 # Display the first few rows and dataset shape
5 print("Dataset Shape:", data.shape)
6 print(data.head())
7

```

➞ Dataset Shape: (1000, 16)

	Lifespan	partType	microstructure	coolingRate	quenchTime	forgeTime	\
0	1469.17	Nozzle	equiGrain	13	3.84	6.47	
1	1793.64	Block	singleGrain	19	2.62	3.48	
2	700.60	Blade	equiGrain	28	0.76	1.34	
3	1082.10	Nozzle	colGrain	9	2.01	2.19	
4	1838.83	Blade	colGrain	16	4.13	3.87	

	HeatTreatTime	Nickel%	Iron%	Cobalt%	Chromium%	smallDefects	\
0	46.87	65.73	16.52	16.82	0.93	10	
1	44.70	54.22	35.38	6.14	4.26	19	
2	9.54	51.83	35.95	8.81	3.41	35	
3	20.29	57.03	23.33	16.86	2.78	0	
4	16.13	59.62	27.37	11.45	1.56	10	

	largeDefects	sliverDefects	seedLocation	castType
0	0	0	Bottom	Die
1	0	0	Bottom	Investment
2	3	0	Bottom	Investment
3	1	0	Top	Continuous
4	0	0	Top	Die

```
1 # Display dataset info
2 print("Dataset Shape:", data.shape)
3 print("Columns:", data.columns)
```

```
⇒ Dataset Shape: (1000, 16)
Columns: Index(['Lifespan', 'partType', 'microstructure', 'coolingRate', 'quenchTime',
               'forgeTime', 'HeatTreatTime', 'Nickel%', 'Iron%', 'Cobalt%',
               'Chromium%', 'smallDefects', 'largeDefects', 'sliverDefects',
               'seedLocation', 'castType'],
              dtype='object')
```

```
1 # Create the binary target variable
2 data['usable'] = (data['Lifespan'] >= 1500).astype(int)
3
4 # Check the distribution of the new target
5 print(data['usable'].value_counts())
6
```

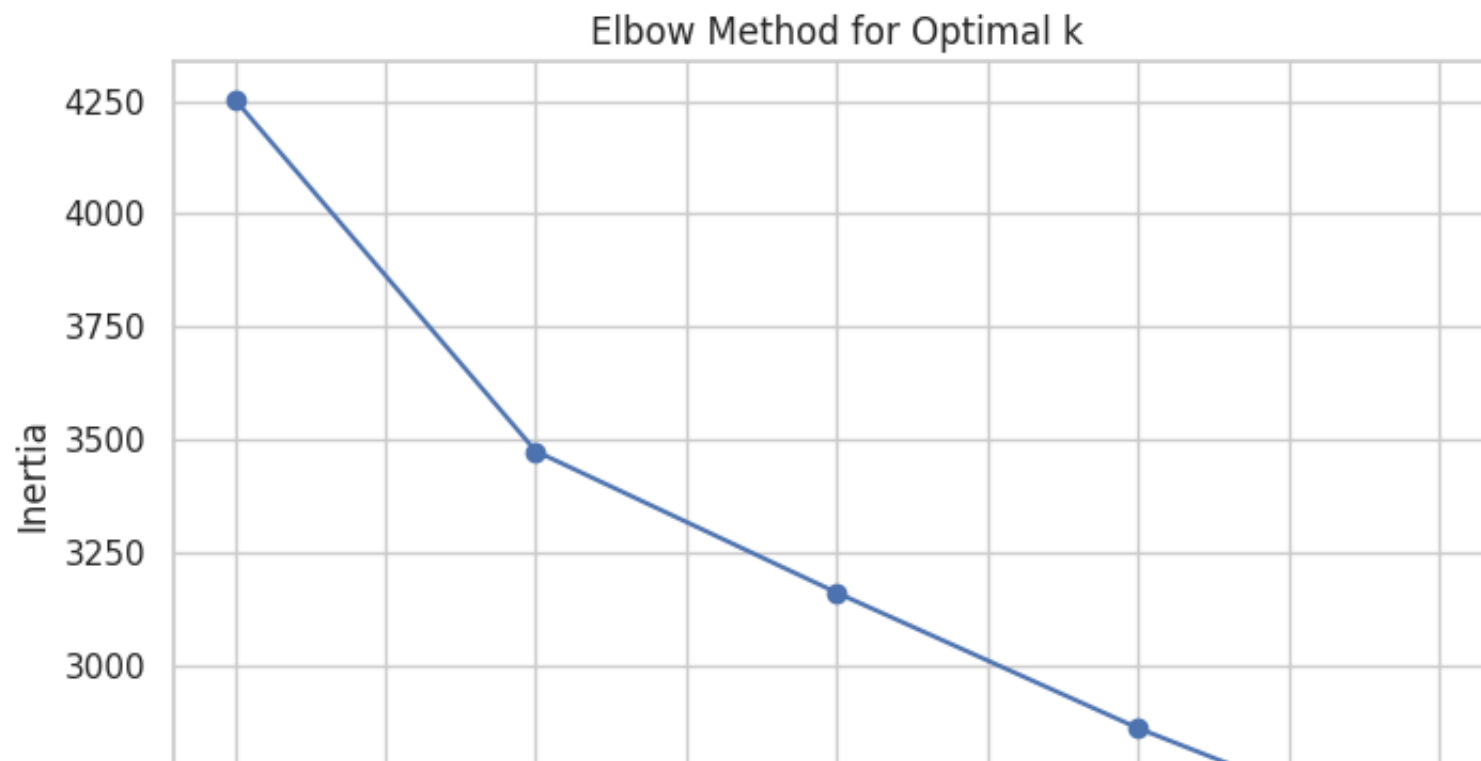
```
⇒ usable
0      694
1      306
Name: count, dtype: int64
```

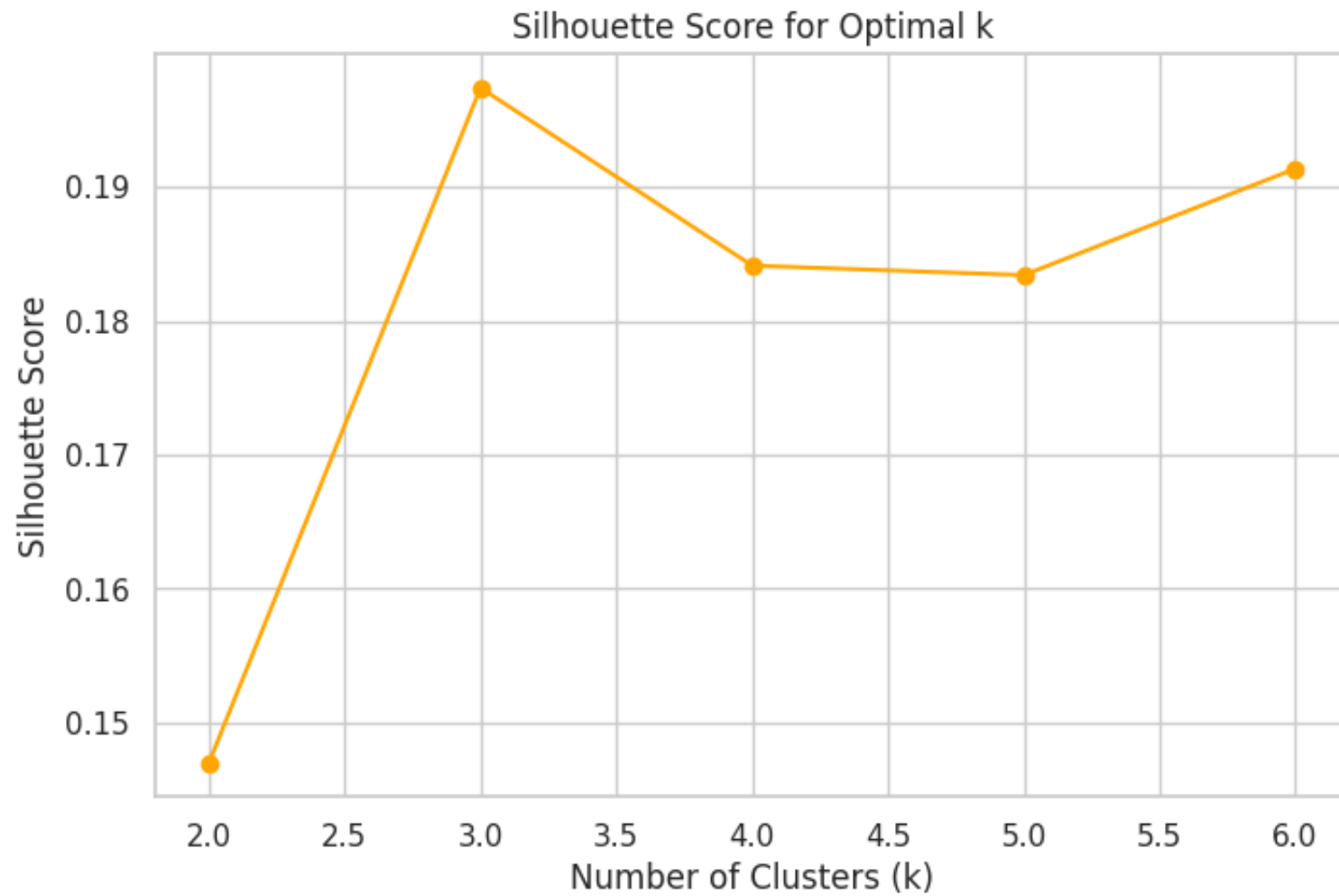
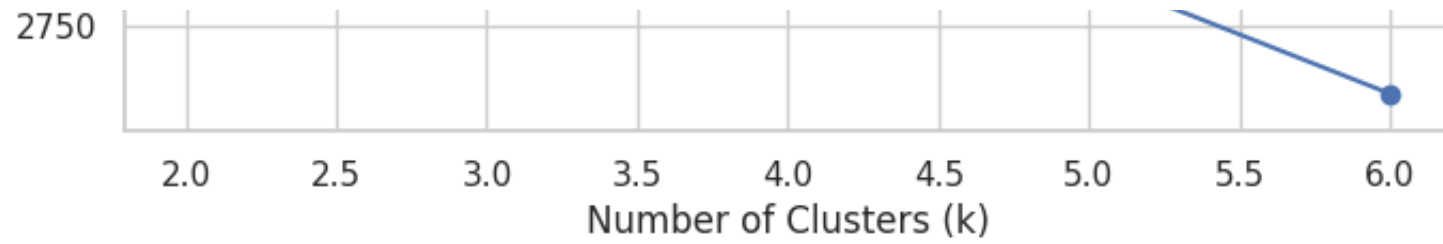
Clustering-Based Feature Crafting

```
1 # Select numerical features related to lifespan prediction
2 clustering_features = data[['Lifespan', 'coolingRate', 'quenchTime', 'forgeTime', 'Heat
3
4 # Standardize the features
5 scaler = StandardScaler()
6 clustering_features_scaled = scaler.fit_transform(clustering_features)
7
```


```
1 # Find optimal number of clusters using Elbow and Silhouette methods
2 inertia = []
3 silhouette_scores = []
4 k_values = range(2, 7)
5
6 for k in k_values:
7     kmeans = KMeans(n_clusters=k, random_state=42)
8     kmeans.fit(clustering_features_scaled)
9     inertia.append(kmeans.inertia_)
10    silhouette_scores.append(silhouette_score(clustering_features_scaled, kmeans.labels
11
12 # Plot Elbow Method
13 plt.figure(figsize=(8, 5))
14 plt.plot(k_values, inertia, marker='o')
15 plt.title('Elbow Method for Optimal k')
16 plt.xlabel('Number of Clusters (k)')
17 plt.ylabel('Inertia')
```

```
18 plt.show()
19
20 # Plot Silhouette Scores
21 plt.figure(figsize=(8, 5))
22 plt.plot(k_values, silhouette_scores, marker='o', color='orange')
23 plt.title('Silhouette Score for Optimal k')
24 plt.xlabel('Number of Clusters (k)')
25 plt.ylabel('Silhouette Score')
26 plt.show()
27
```



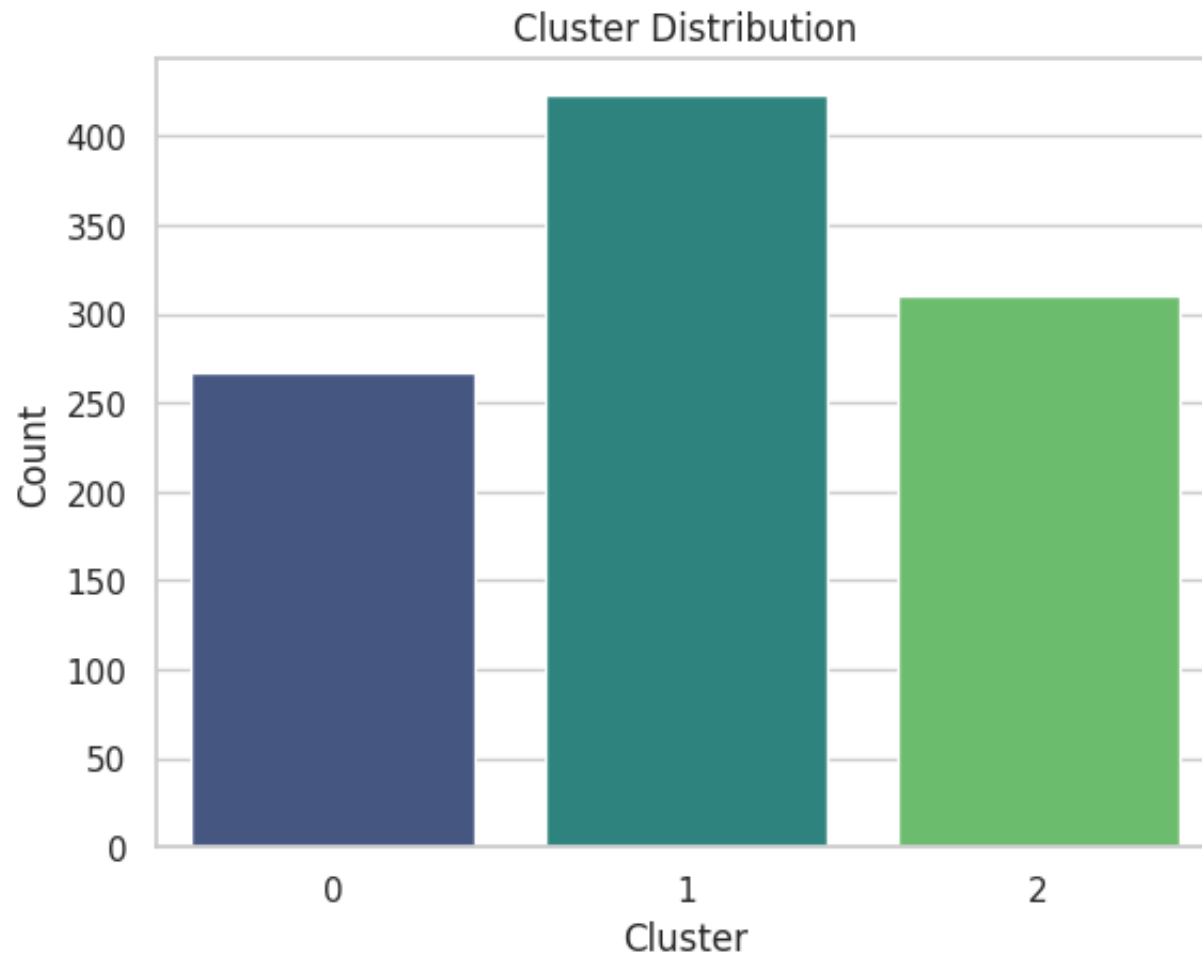


```
1 # Apply K-Means clustering with the chosen k
2 optimal_k = 3
3 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
4 data['Lifespan_clusters'] = kmeans.fit_predict(clustering_features_scaled)
5
6 # Visualize the cluster distribution
7 sns.countplot(x=data['Lifespan_clusters'], palette='viridis')
8 plt.title('Cluster Distribution')
9 plt.xlabel('Cluster')
10 plt.ylabel('Count')
11 plt.show()
12
```

 <ipython-input-38-6abf837243ed>:7: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable

```
sns.countplot(x=data['Lifespan_clusters'], palette='viridis')
```




```
1 # Binary classification target
2 y_binary = data['usable']
3
4 # Multi-class clustering target
5 y_clusters = data['Lifespan_clusters']
6
7 # Drop unnecessary columns
8 X = data.drop(columns=['Lifespan', 'usable', 'Lifespan_clusters'])
9 X = pd.get_dummies(X, drop_first=True) # One-hot encode categorical features
10
```

```
1 # Binary classification split
2 X_train_bin, X_test_bin, y_train_bin, y_test_bin = train_test_split(X, y_binary, test_s
3
4 # Multi-class classification split
5 X_train_clust, X_test_clust, y_train_clust, y_test_clust = train_test_split(X, y_cluste
6
```

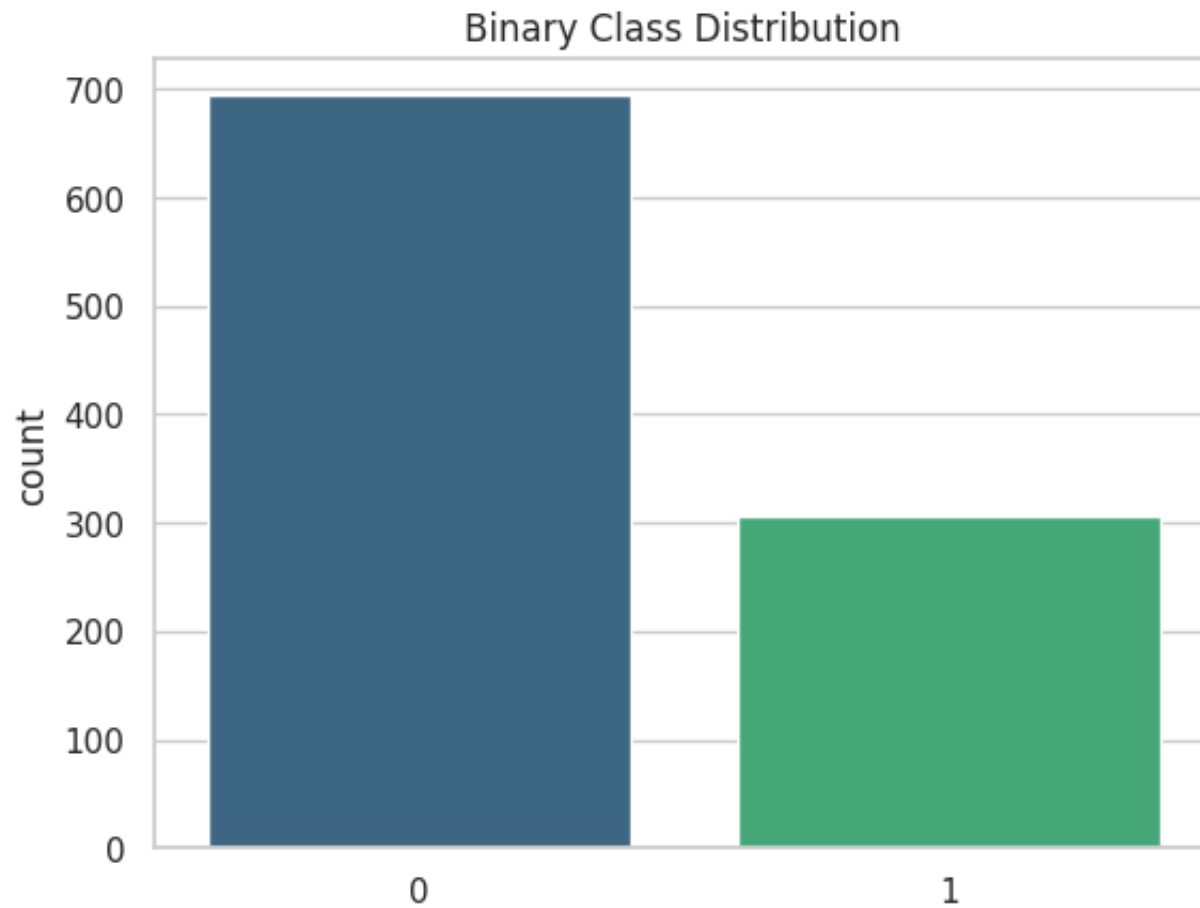
```
1 # Binary class balance
2 sns.countplot(x=y_binary, palette='viridis')
3 plt.title('Binary Class Distribution')
4 plt.show()
5
6 # Multi-class balance
7 sns.countplot(x=y_clusters, palette='viridis')
```

```
8 plt.title('Multi-Class Cluster Distribution')
9 plt.show()
10
```

 <ipython-input-41-b5f36c39fceb>:2: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable

```
sns.countplot(x=y_binary, palette='viridis')
```

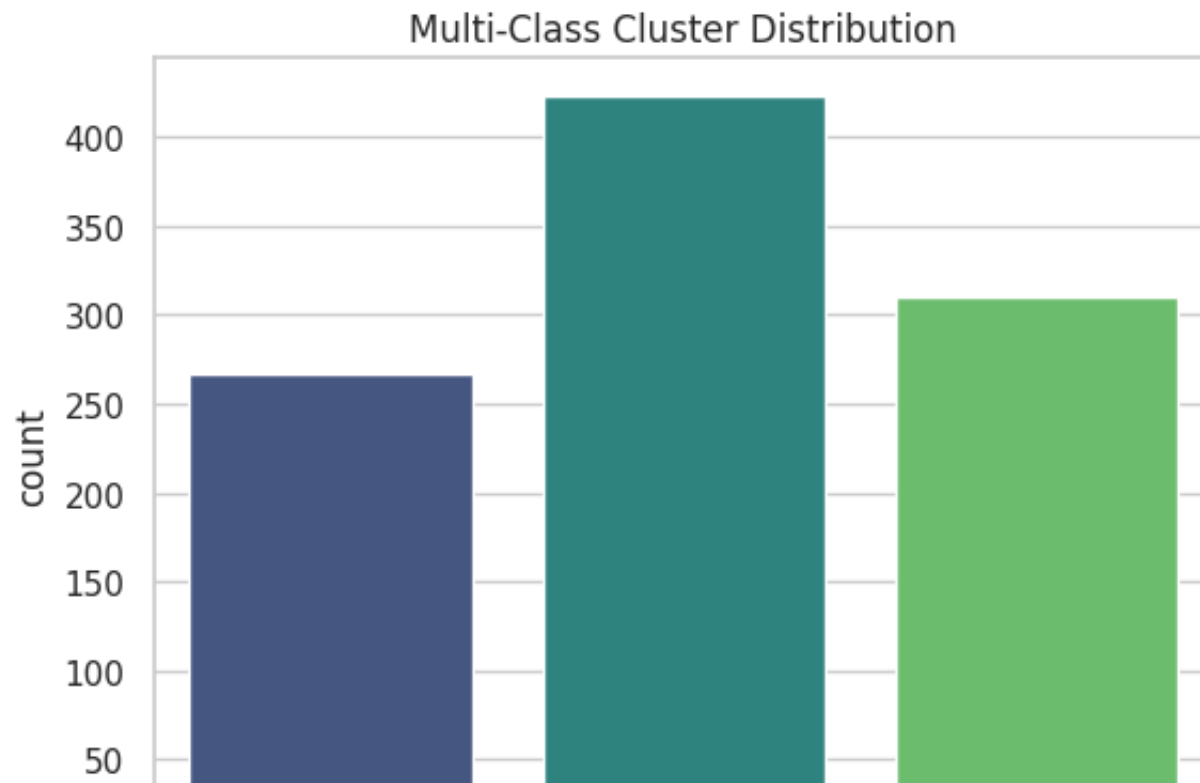


usable

```
<ipython-input-41-b5f36c39fceb>:7: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable

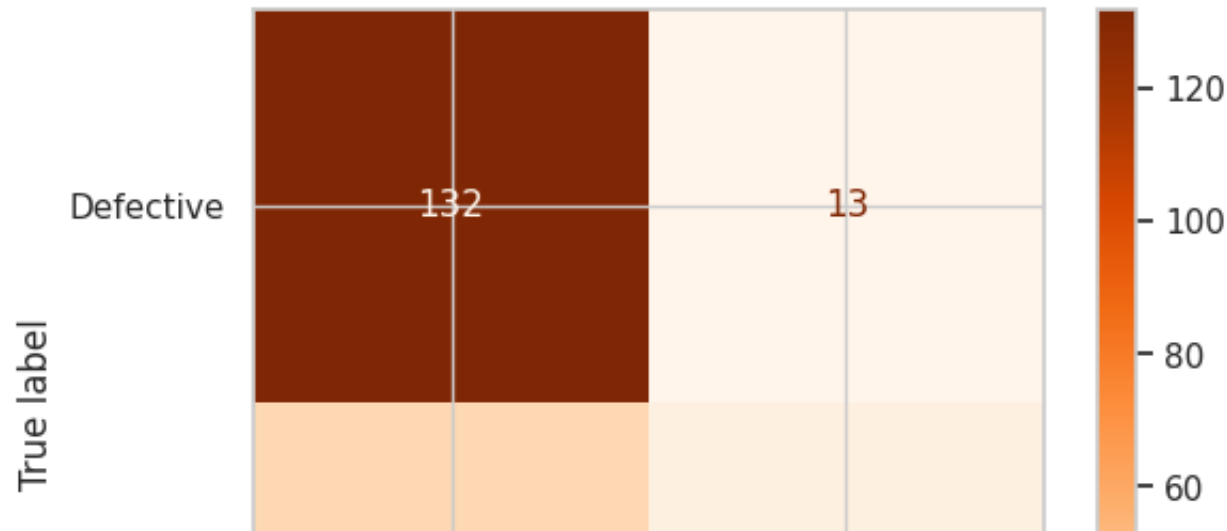
```
sns.countplot(x=y_clusters, palette='viridis')
```

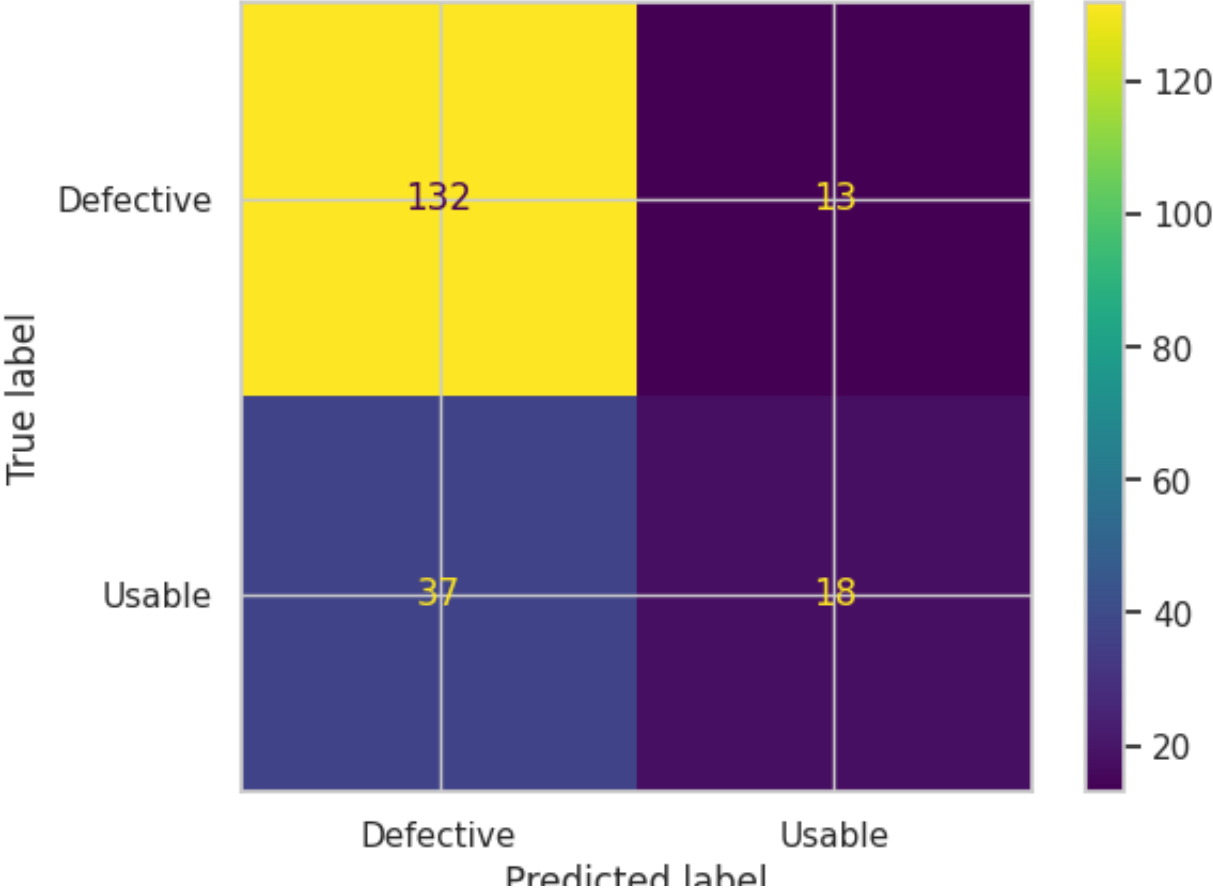
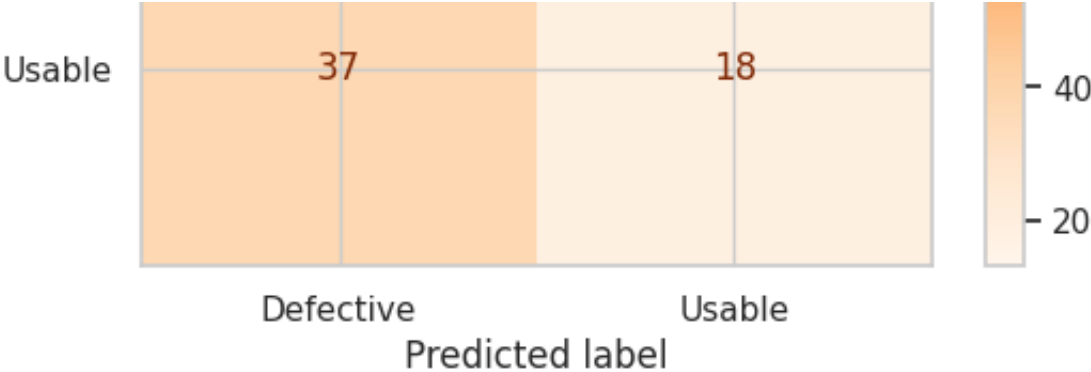


```
1 # Apply SMOTE for multi-class
2 smote = SMOTE(random_state=42)
3 X_train_clust_bal, y_train_clust_bal = smote.fit_resample(X_train_clust, y_train_clust)
4
```

```
1 # Logistic Regression
2 log_reg = LogisticRegression(max_iter=1000, random_state=42)
3 log_reg.fit(X_train_bin, y_train_bin)
4
5 # Predictions and Evaluation
6 log_reg_preds = log_reg.predict(X_test_bin)
7 print("Logistic Regression Binary Classification:")
8 print("Accuracy:", accuracy_score(y_test_bin, log_reg_preds))
9 print("F1 Score:", f1_score(y_test_bin, log_reg_preds))
10 ConfusionMatrixDisplay.from_estimator(log_reg, X_test_bin, y_test_bin, display_labels=[
11 plt.show()
12
```

➞ Logistic Regression Binary Classification:
Accuracy: 0.75
F1 Score: 0.4186046511627907



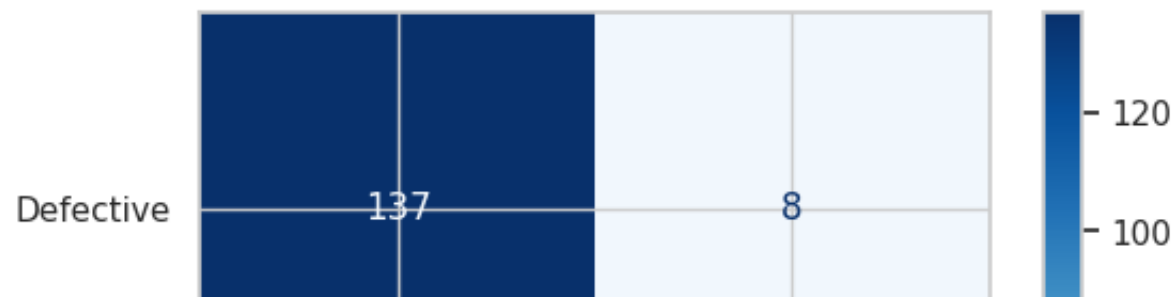


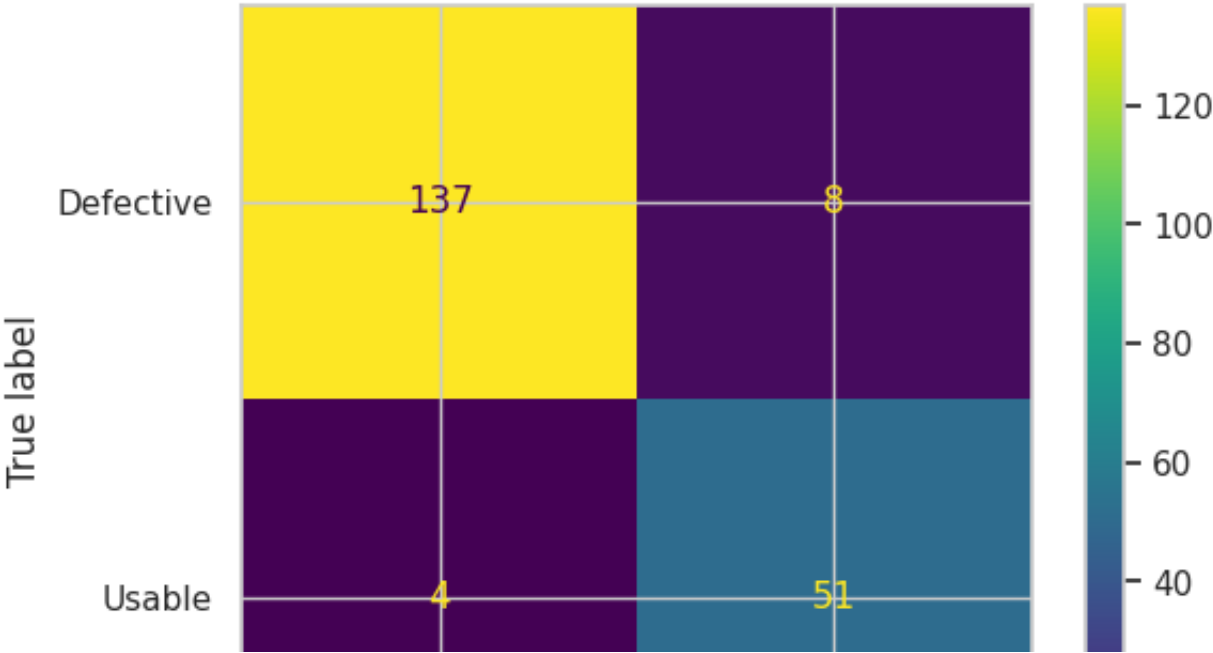
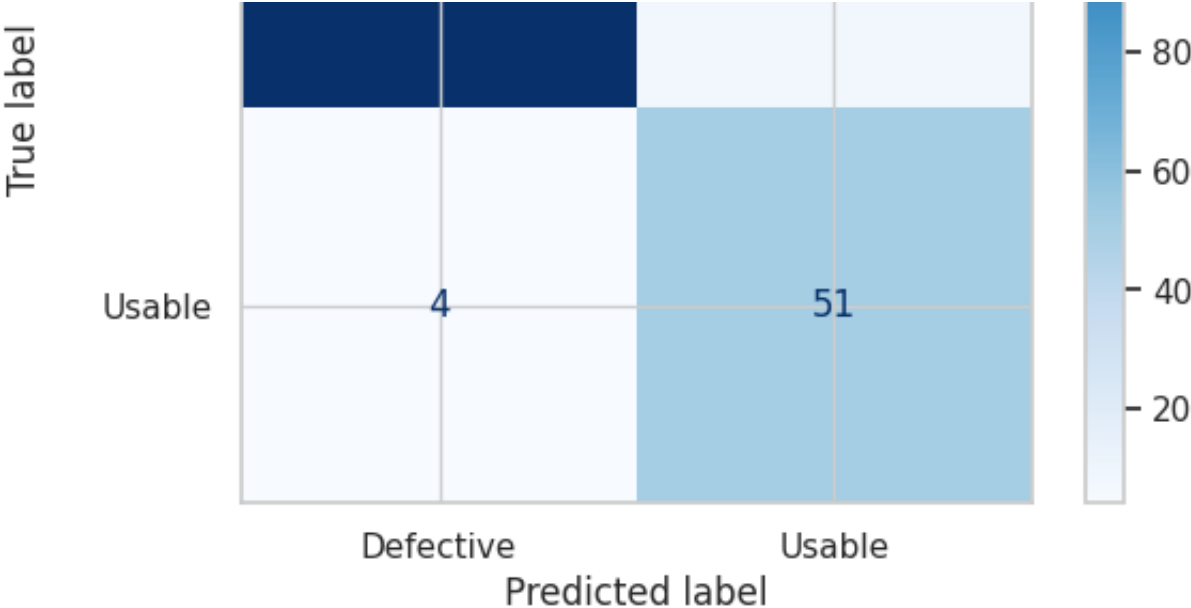
```

1 # Gradient Boosting
2 gb_clf = GradientBoostingClassifier(random_state=42)
3 gb_clf_params = {'n_estimators': [50, 100], 'learning_rate': [0.01, 0.1], 'max_depth':
4 gb_clf_grid = GridSearchCV(gb_clf, gb_clf_params, cv=3, scoring='accuracy')
5 gb_clf_grid.fit(X_train_bin, y_train_bin)
6
7 # Best model
8 best_gb = gb_clf_grid.best_estimator_
9 gb_preds = best_gb.predict(X_test_bin)
10 print("Gradient Boosting Binary Classification:")
11 print("Accuracy:", accuracy_score(y_test_bin, gb_preds))
12 print("F1 Score:", f1_score(y_test_bin, gb_preds))
13 ConfusionMatrixDisplay.from_estimator(best_gb, X_test_bin, y_test_bin, display_labels=[
14 plt.show()
15

```

➔ Gradient Boosting Binary Classification:
 Accuracy: 0.94
 F1 Score: 0.8947368421052632







```
1 # Logistic Regression
2 multi_log = LogisticRegression(random_state=42, multi_class='multinomial')
3 multi_log.fit(X_train_clust_bal, y_train_clust_bal)
4
5 # Predictions and Evaluation
6 multi_log_preds = multi_log.predict(X_test_clust)
7 print("Logistic Regression Multi-Class Classification:")
8 print(classification_report(y_test_clust, multi_log_preds))
9
```


Logistic Regression Multi-Class Classification:

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.93	0.97	0.95	58
1	0.91	0.90	0.90	77
2	0.94	0.92	0.93	65

accuracy			0.93	200
macro avg	0.93	0.93	0.93	200
weighted avg	0.92	0.93	0.92	200

```
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:469: ConvergenceWarning: lbfgs failed
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```

```

1 # Random Forest
2 rf_clf_params = {'n_estimators': [50, 100, 200], 'max_depth': [5, 10, 15]}
3 rf_clf_grid = GridSearchCV(RandomForestClassifier(random_state=42), rf_clf_params, cv=3)
4 rf_clf_grid.fit(X_train_clust_bal, y_train_clust_bal)
5 best_rf_clf = rf_clf_grid.best_estimator_
6 rf_preds = best_rf_clf.predict(X_test_clust)
7
8 print("Random Forest Multi-Class Classification:")
9 print(classification_report(y_test_clust, rf_preds))
10

```

➞ Random Forest Multi-Class Classification:

	precision	recall	f1-score	support
0	0.98	0.93	0.96	58
1	0.88	0.94	0.91	77
2	0.94	0.91	0.92	65
accuracy			0.93	200
macro avg	0.93	0.92	0.93	200
weighted avg	0.93	0.93	0.93	200

```

1 # Results Summary
2 results = {
3     'Model': ['Binary Logistic Regression', 'Binary Gradient Boosting', 'Multi-Class Lo
4     'Accuracy': [accuracy_score(y_test_bin, log_reg_preds), accuracy_score(y_test_bin,
5                     accuracy_score(y_test_clust, multi_log_preds), accuracy_score(y_test_c

```

```

6     'F1 Score': [f1_score(y_test_bin, log_reg_preds), f1_score(y_test_bin, gb_preds, av
7                   f1_score(y_test_clust, multi_log_preds, average='weighted'), f1_score(
8 }
9 results_df = pd.DataFrame(results)
10 print(results_df)
11 results_df.set_index('Model').plot(kind='bar', figsize=(12, 8))
12 plt.title('Model Performance Comparison')
13 plt.ylabel('Score')
14 plt.xticks(rotation=0)
15 plt.show()
16

```

↔

	Model	Accuracy	F1 Score
0	Binary Logistic Regression	0.750	0.418605
1	Binary Gradient Boosting	0.940	0.940633
2	Multi-Class Logistic Regression	0.925	0.924835
3	Multi-Class Random Forest	0.925	0.925457



