



Ain Shams University – Faculty of Engineering I-CHEP

# Machine Learning Model

Report

Introduction to Artificial Intelligence – CSE-281

December 2024

Submitted to:

Dr. Mariam Nabil

Eng. Mohammad Essam

---

Submitted by:

- 23P0270: Ahmed Sherif Nasri
  - 23P0314: Omar Ahmed Shabaan
  - 2300535: Mohamed Tarek Elsayed
  - 23P0182: Peter Maged Shokry
  - 23P0282: Mohamed Ahmed Mohamed
  - 23P0171: Omar Montaser Mahmoud
-

1 Contents

2 Data Analysis..... 2

2.1 Overview..... 2

2.2 Key Insights ..... 2

Feature X1 (Item Identifier): ..... 2

Feature X2 (Numerical): ..... 2

3 Feature Selection ..... 3

3.1 Initial Selection ..... 3

3.2 Methods Used ..... 3

3.3 Final Selection ..... 3

4 Feature Engineering..... 3

4.1 Transformations: ..... 3

4.2 Encoding: ..... 4

4.3 Interaction Terms: ..... 4

4.4 Missing Values: ..... 4

5 Models and Hyperparameter Selection..... 4

5.1 Models Evaluated: ..... 4

5.2 Hyperparameter Optimization: ..... 5

5.3 Model Performance:..... 5

6 Conclusion: ..... 6

7 Appendix ..... 6

7.1 Visualization ..... 6

7.2 Code ..... 14

7.2.1 CatBoost..... 14

7.3 Linear Model..... 24

7.4 SVR Model ..... 28

7.5 XGBoost Model ..... 33

## 2 Data Analysis

### 2.1 Overview

Dataset Size: The dataset contains 6000 entries with multiple numerical and categorical features.

Target Variable: Y (specifics not provided, assumed to be the dependent variable for prediction).

### 2.2 Key Insights

Feature X1 (Item Identifier):

- Contains 1553 unique identifiers.
- High cardinality necessitated feature splitting into `Category`, `Subcategory`, and `SKU`.
- Categories reduced to manageable values (e.g., 3 for `Category`).

Feature X2 (Numerical):

- Missing values (1006) handled by imputing the mean.
- Statistical Summary: Mean = 12.96, Std Dev = 4.65, Min = 4.55, Max = 21.35.

Feature X3 (Fat Content):

- Inconsistent labels cleaned and unified into `LF` (Low Fat) and `REG` (Regular).
- Distribution: 3889 (`LF`), 2111 (`REG`).

Feature X4 (Visibility):

- Positively skewed data with many zeros (360 occurrences).
- Log transformation considered for normalization.

Feature X7 (Outlet):

- Even distribution across 10 unique outlets.
- Encoding strategies tested: ordinal encoding and one-hot encoding.

Feature X8 (Establishment Year):

- Transformed to `PastYearsCont` (Years since 2009) and binned into categorical groups (`very old`, `old`, `recent`, `new`).

Other Key Features:

- X5: Product category (16 unique values).
- X9: Size (Small, Medium, High) transformed using ordinal encoding.
- X10: Location type (Tier 1, Tier 2, Tier 3), assumed to reflect luxury.

### 3 Feature Selection

#### 3.1 Initial Selection

Features were selected based on domain knowledge and their statistical relationship with the target variable.

#### 3.2 Methods Used

1. Correlation Thresholding:
  - Dropped features with low correlation to the target variable.
2. Recursive Feature Elimination (RFE):
  - Iteratively removed less impactful features.
3. Feature Importance (Tree-based Models):
  - Ranked features using Random Forest.
4. Mutual Information
5. Permutation Importance

#### 3.3 Final Selection

The final feature set included:

- High-correlation variables like Category, Fat Content, and Visibility.
- Features ranked as important across multiple methods.

### 4 Feature Engineering

#### 4.1 Transformations:

- Normalization/Standardization: Continuous variables were scaled using StandardScaler for models sensitive to feature magnitudes.
- Log Transformation: Applied to skewed variables to improve model interpretability and performance.
- Some custom transformers

## 4.2 Encoding:

- Categorical Variables: Encoded using one-hot encoding or label encoding, depending on the algorithm requirements.

## 4.3 Interaction Terms:

- New features representing interactions between key variables were created to capture non-linear relationships.

## 4.4 Missing Values:

- Imputed using mean, median, or mode strategies for numerical features and the most frequent category for categorical features.

# 5 Models and Hyperparameter Selection

## 5.1 Models Evaluated:

### 1. Linear Regression:

- Hyperparameters tuning: Fit intercept set to false.

### 2. Gradient Boosting (CATBoost):

- Hyperparameters tuning:

```
'iterations': 664,  
'learning_rate': 0.0362387234017892,  
'depth': 3,  
'subsample': 0.7999456005681143,  
'colsample_bylevel': 0.7763198164104148,  
'l2_leaf_reg': 5.466280028000805,  
'random_strength': 1.1183292594718606,  
'bagging_temperature': 3.5375530112539995,  
'border_count': 105,  
'loss_function': 'MAE',  
'verbose': 0,  
smoothing = 0.4966048630439446
```

## 3. Gradient Boosting (XGBoost):

## 4. Hyperparameters tuned:

```
'n_estimators': 670,  
'learning_rate': 0.06035467909986324,  
'max_depth': 6,  
'subsample': 0.923204157140139,  
'colsample_bytree': 0.4883434932394659,  
'reg_alpha': 5.410450954065299,  
'reg_lambda': 1.1186894717633518,  
'min_child_weight': 5,  
'gamma': 2.3596898720558754,  
'max_bin': 300,  
'objective': 'reg:absoluteerror',  
'eval_metric': 'mae',  
'random_state': 42,  
'enable_categorical': True
```

## 5. Support Vector Regressor (SVR):

## ○ Hyperparameters tuned

```
'C': 2.6311460970378824,  
'epsilon': 0.35225148545949075,  
'kernel': 'poly',  
'degree': 2,  
'gamma': 'scale'
```

## 5.2 Hyperparameter Optimization:

- Optuna Search was used for efficient tuning.
- Cross-validation ensured robust evaluation of hyperparameter configurations.

## 5.3 Model Performance:

- Models were evaluated using metrics such as MAE.

- XGBoost and CatBoost consistently outperformed other models.

## 6 Conclusion:

The analysis and modeling process provided a comprehensive approach to developing and evaluating machine learning models. Feature engineering and hyperparameter significantly improved model performance, with XGBoost and CatBoost emerging as the most effective models.

## 7 Appendix

### 7.1 Visualization

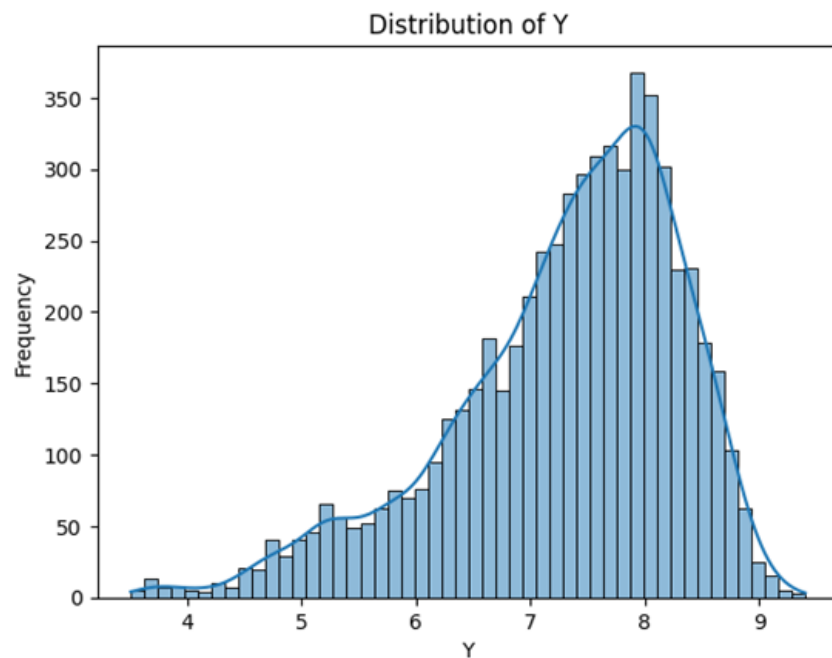


Figure 1 Correlation between Distribution and Frequency of Y

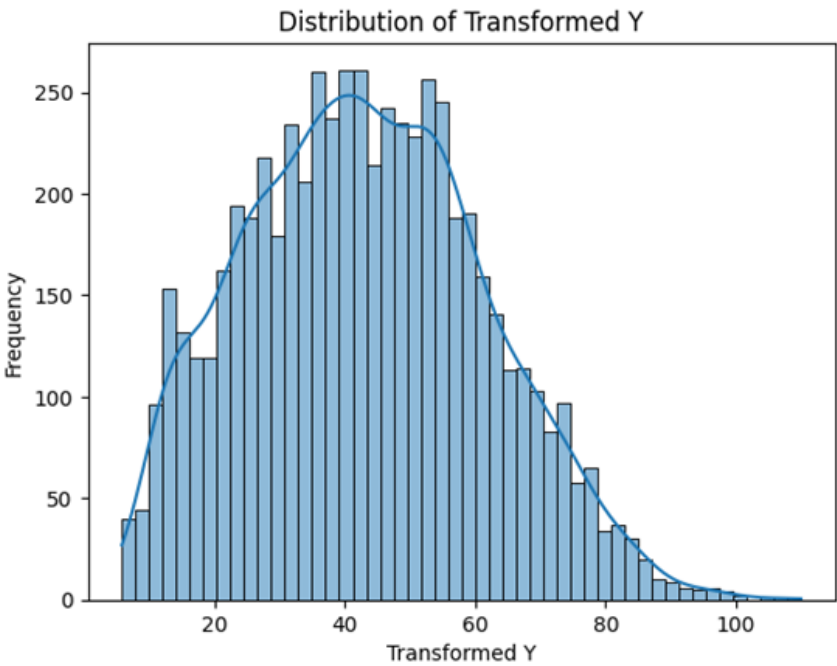


Figure 2 Correlation between Distribution and Frequency of Transformed Y

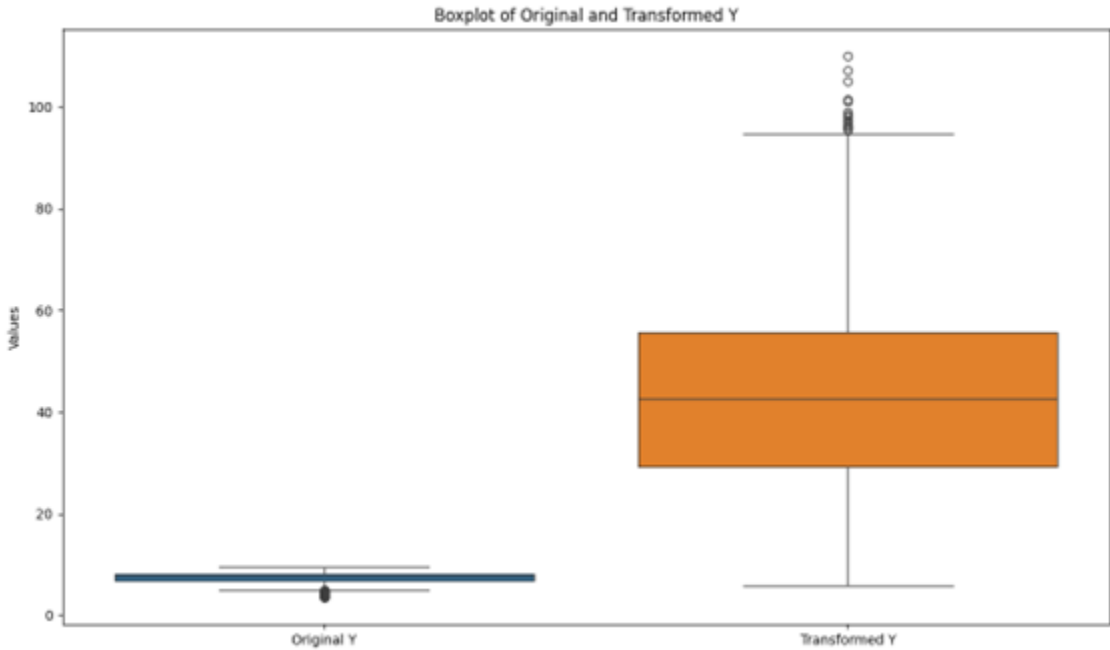


Figure 3 Boxplot of Original and Transformed Y



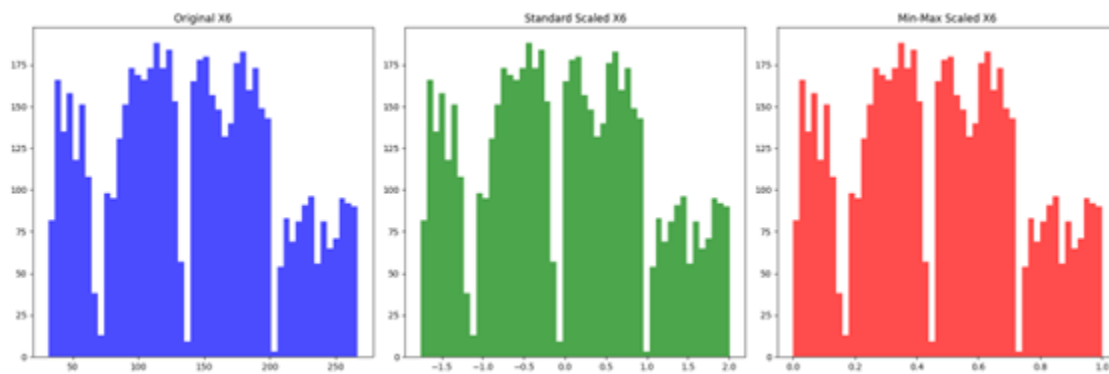


Figure 4 Apply Standard Scaler and Min-Max Scaler to the X6 column

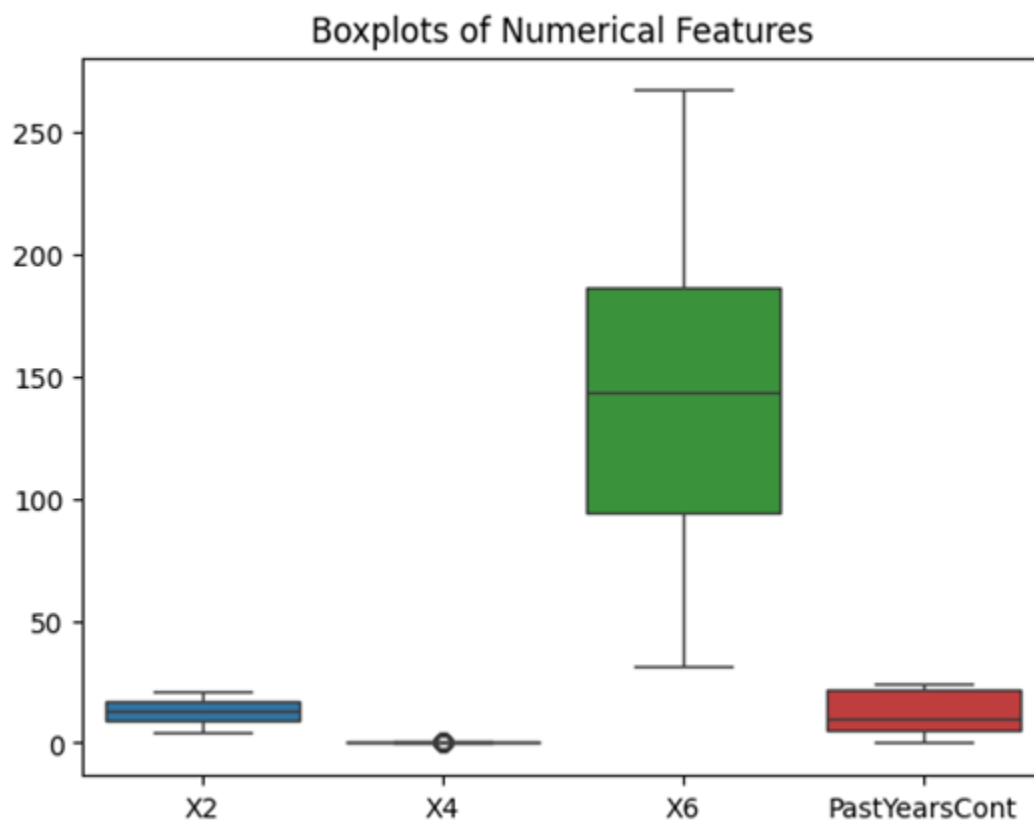
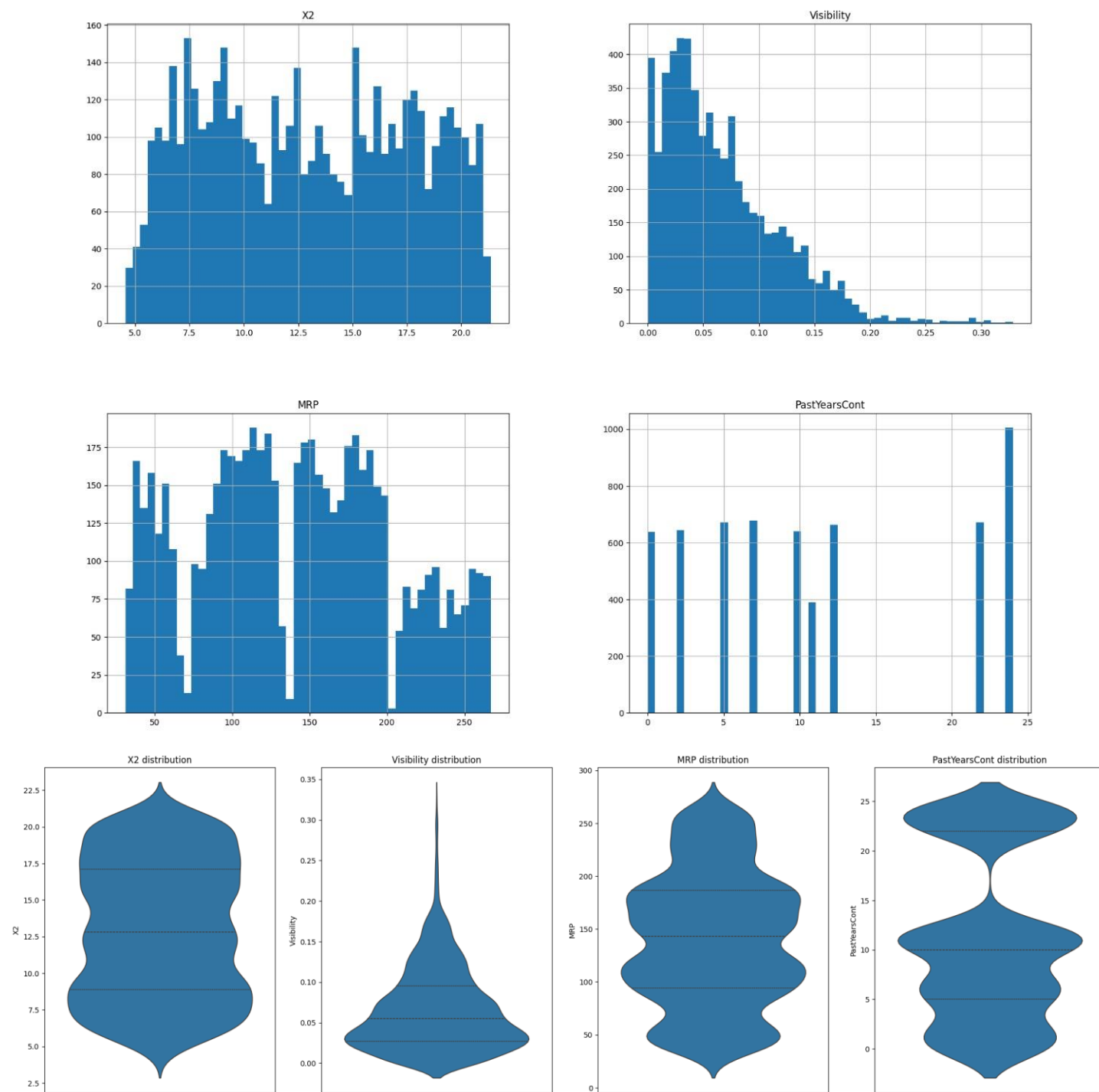
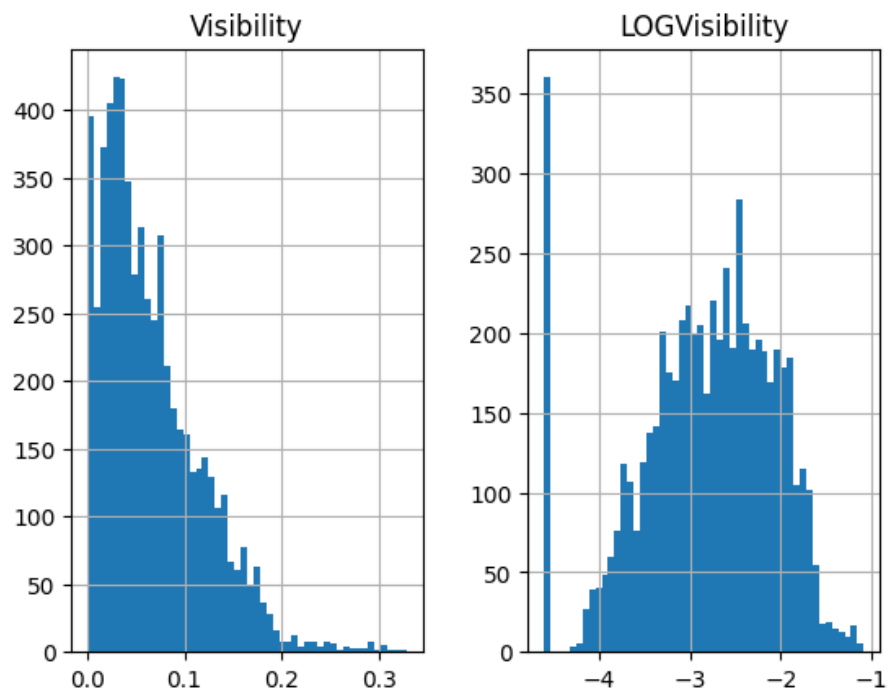


Figure 5 Overview of Numerical Features





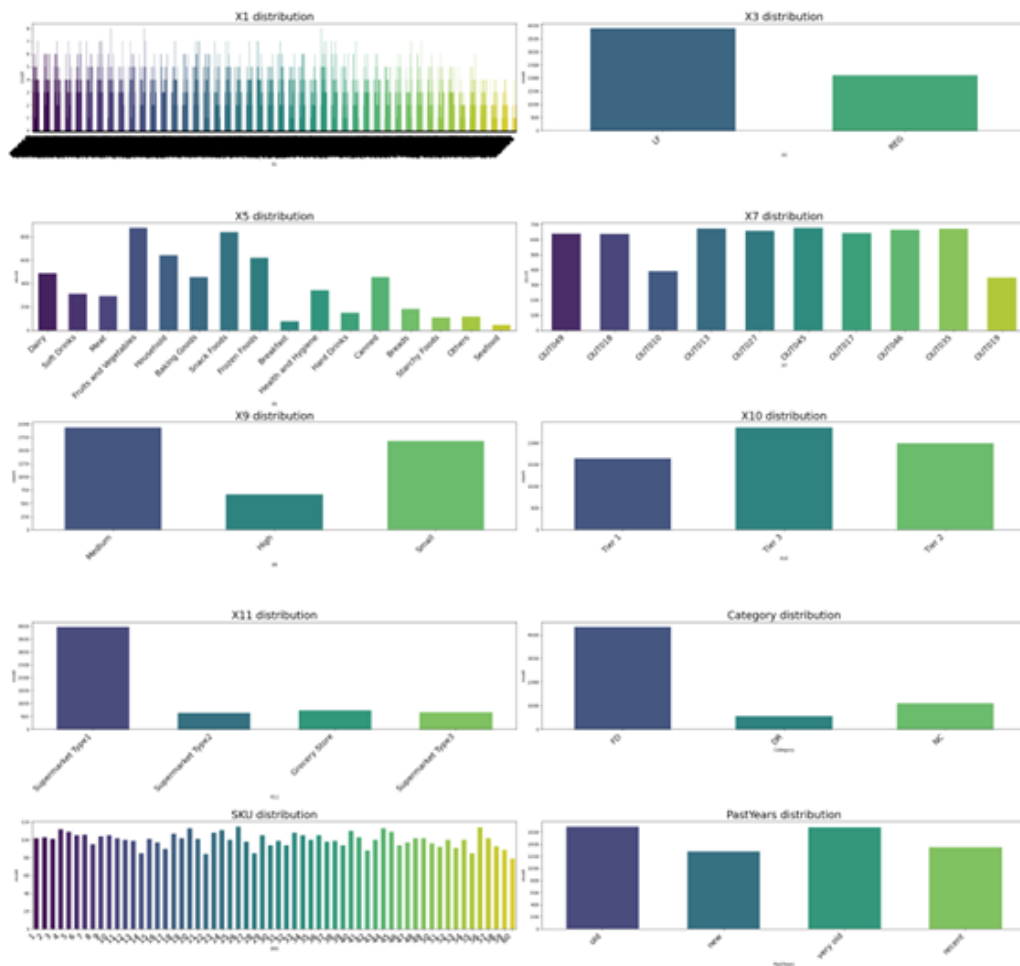


Figure 6 Overview of the all the features distributions

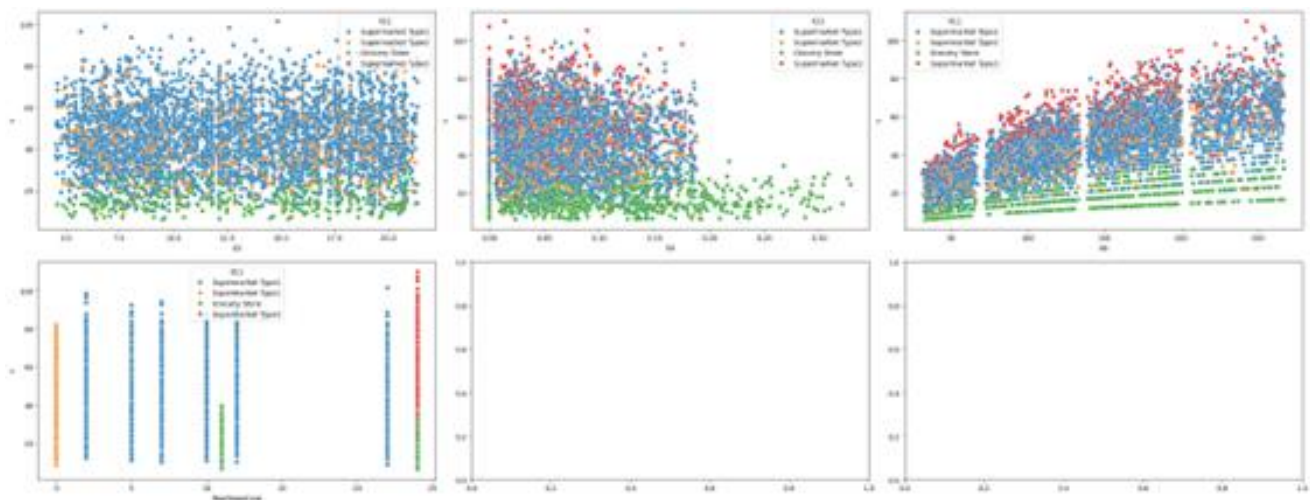
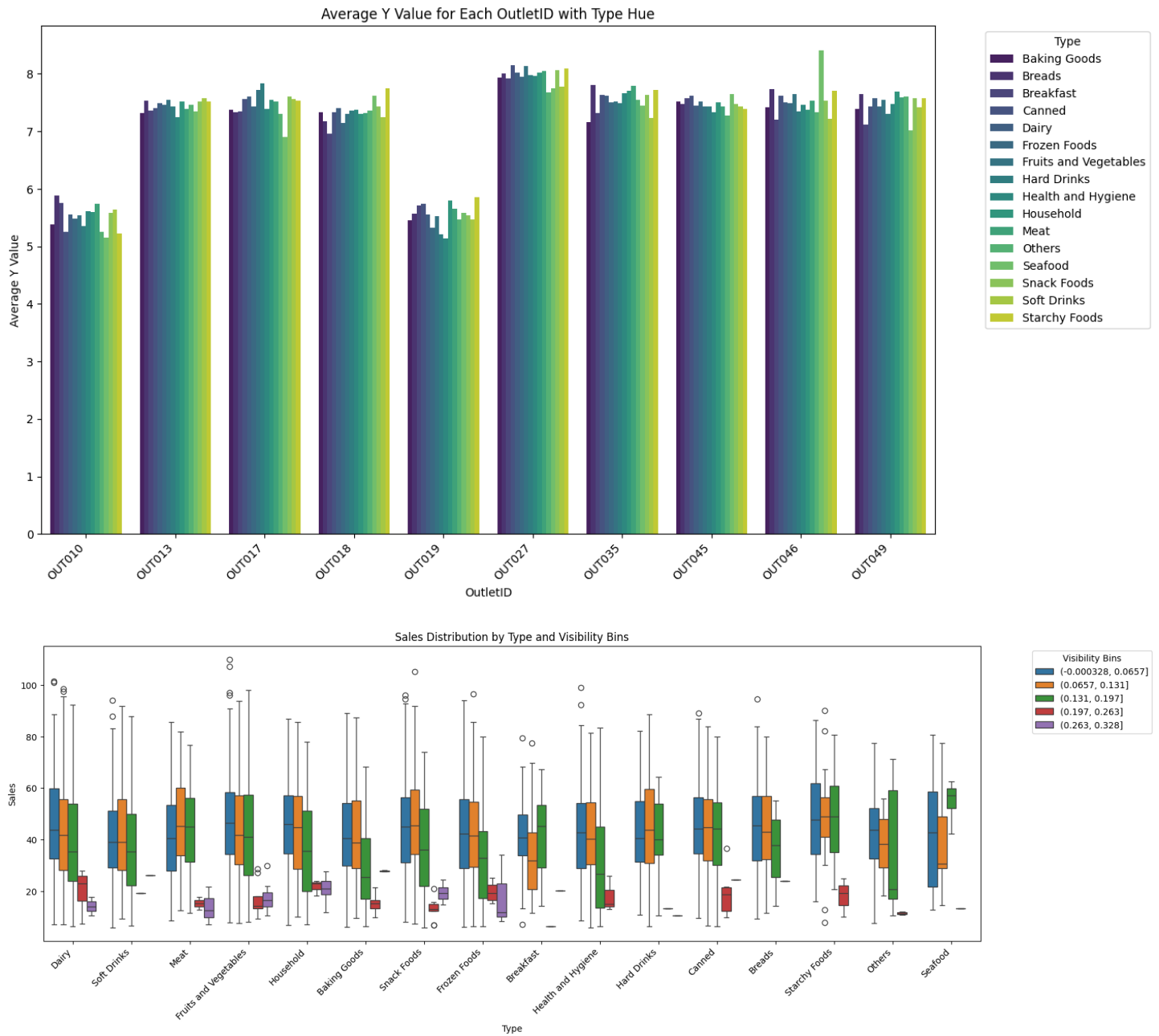


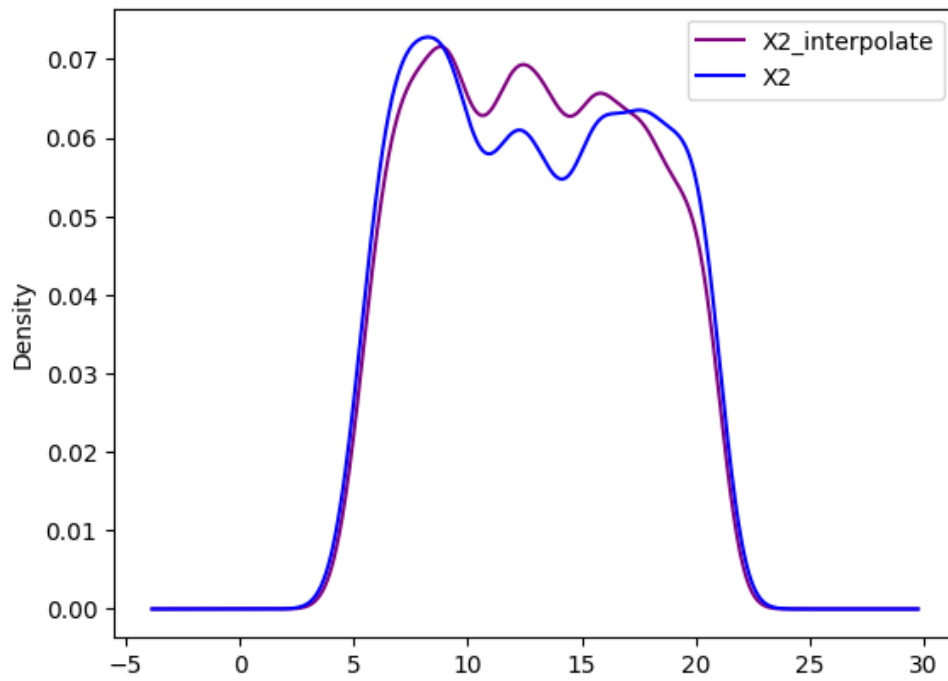
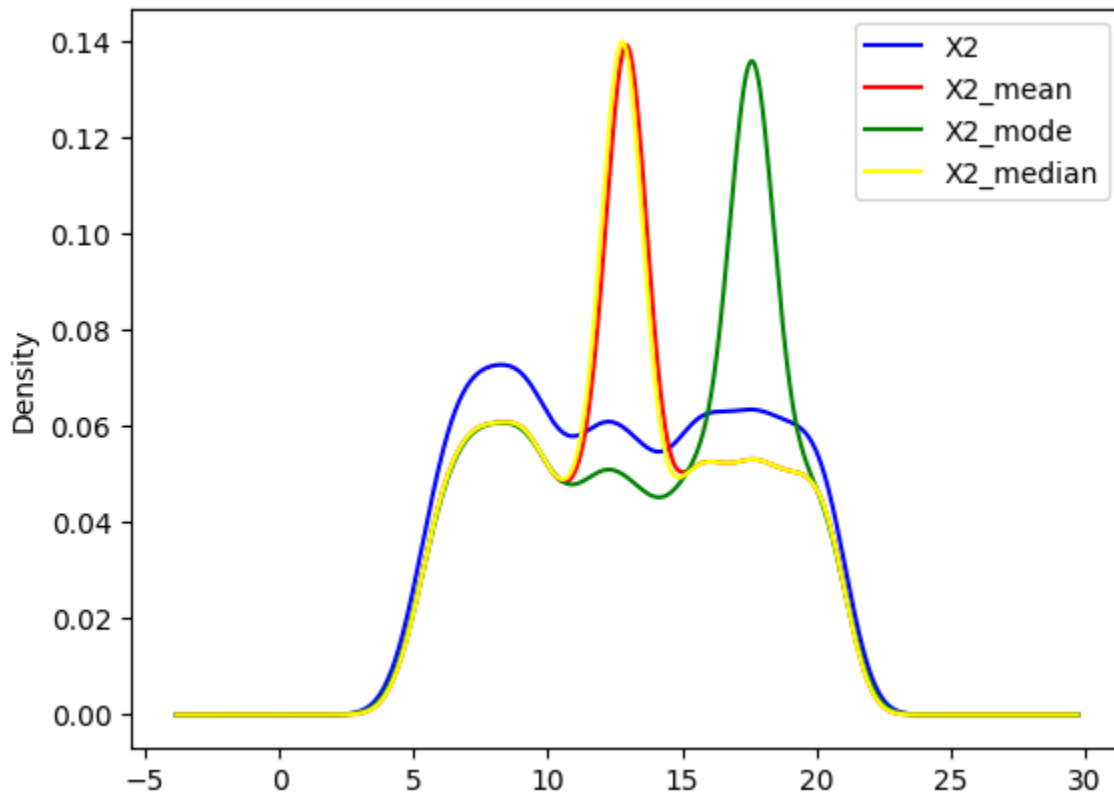
Figure 7 Correlations between features



## Handling Missing Values:

Univariate Analysis:

X2 imputation techniques:



## 7.2 Code

### 7.2.1 CatBoost

```
1. class TypeMeanPriceTransformer(BaseEstimator, TransformerMixin):
2.     def __init__(self, smoothing=0):
3.         self.smoothing = smoothing
4.         self.category_mean_price = None
5.         self.global_type_mean = None
6.
7.     def fit(self, X, y=None):
8.         if y is None:
9.             raise ValueError("Target variable `y` cannot be None in the fit method.")
10.
11.         data_ = X.copy()
12.         data_['Y'] = y
13.
14.         # Global mean for each Type across all OutletTypes
15.         self.global_type_mean = data_.groupby('Type')['Y'].mean()
16.
```

```

17.         # Local mean and count for each (OutletType, Type)
18.         local_mean = data_.groupby(['Outlet_Type', 'Type'])['Y'].mean()
19.         local_count = data_.groupby(['Outlet_Type', 'Type'])['Y'].count()
20.
21.         # Merge the global Type mean with the local statistics
22.         merged = local_mean.reset_index().merge(
23.             self.global_type_mean.reset_index(), on='Type', how='left', suffixes=('_local', '_global')
24.         )
25.
26.         # Apply smoothing
27.         merged['TypeMeanPrice'] = (
28.             (merged['Y_local'] * local_count.values) +
29.             (merged['Y_global'] * self.smoothing)
30.         ) / (local_count.values + self.smoothing)
31.
32.         self.category_mean_price = merged[['Outlet_Type', 'Type', 'TypeMeanPrice']]
33.         return self
34.
35.     def transform(self, X):
36.         X = X.copy()
37.         # Merge the smoothed means back into the original data
38.         X = pd.merge(X, self.category_mean_price, on=['Outlet_Type', 'Type'], how='left')
39.         #X['TypeMeanPrice'] = np.log1p(X['TypeMeanPrice'])
40.         return X[['TypeMeanPrice']]
41.

```

```

1. class VisibilityZerosImputerZerosImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.item_Visibility_mean = defaultdict(lambda: None) # Stores mean of Visibility per ItemID
4.         self.global_mean = None # Fallback global mean for Visibility
5.
6.     def fit(self, X, y=None):
7.         X = X.copy()
8.
9.         # Ensure the required columns exist
10.        if 'ItemID' not in X.columns or 'Visibility' not in X.columns:
11.            raise ValueError("Both 'ItemID' and 'Visibility' columns must be present in the dataset.")
12.
13.        # Ensure no None values in ItemID
14.        if X['ItemID'].isnull().any():
15.            raise ValueError("ItemID column contains None values.")
16.
17.        # Calculate mean Visibility for each ItemID
18.        item_Visibility_mean = X.groupby('ItemID')['Visibility'].mean()
19.        self.item_Visibility_mean.update(item_Visibility_mean.to_dict())
20.
21.        # Calculate global mean for the Visibility column
22.        if X['Visibility'].notnull().any():
23.            self.global_mean = X['Visibility'].mean()
24.        else:
25.            raise ValueError("Visibility column contains only NaN values.")
26.
27.        return self
28.
29.     def transform(self, X):
30.         X = X.copy()
31.
32.        # Ensure the required columns exist
33.        if 'ItemID' not in X.columns or 'Visibility' not in X.columns:

```



```

34.         raise ValueError("Both 'ItemID' and 'Visibility' columns must be present in the dataset.")
35.
36.     # Ensure no None values in ItemID
37.     if X['ItemID'].isnull().any():
38.         raise ValueError("ItemID column contains None values.")
39.
40.     # Safely impute zero Visibilitys based on ItemID or global mean
41.     X['Visibility'] = X.apply(
42.         lambda row: self.item_Visibility_mean.get(row['ItemID'], self.global_mean)
43.         if row['Visibility'] == 0
44.         else row['Visibility'],
45.         axis=1
46.     )
47.
48.     X['Visibility'] = X['Visibility'].replace(0, self.global_mean)
49.     X['Visibility'] = np.sqrt(X['Visibility'])
50.
51.     return X[['Visibility']]
52.

```

```

1. class X2NaNsImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.item_X2_mode = defaultdict(lambda: None) # Stores mode of X2 per ItemID
4.         self.global_mean = None # Fallback global mean for X2
5.
6.     def fit(self, X, y=None):
7.         X = X.copy()
8.         # Ensure the required columns exist
9.         if 'ItemID' not in X.columns or 'X2' not in X.columns:
10.            raise ValueError("Both 'ItemID' and 'X2' columns must be present in the dataset.")
11.
12.        # Ensure no None values in ItemID
13.        if X['ItemID'].isnull().any():
14.            raise ValueError("ItemID column contains None values.")
15.
16.        # Calculate mode X2 for each ItemID
17.        item_X2_mode = X.groupby('ItemID')['X2'].agg(lambda x: x.mode().iloc[0] if not x.mode().empty else
np.nan)
18.        self.item_X2_mode.update(item_X2_mode.to_dict())
19.
20.        # Calculate global mean for the X2 column
21.        if X['X2'].notnull().any():
22.            self.global_mean = X['X2'].mean()
23.        else:
24.            raise ValueError("X2 column contains only NaN values.")
25.        return self
26.
27.    def transform(self, X):
28.        X = X.copy()
29.
30.        # Ensure the required columns exist
31.        if 'ItemID' not in X.columns or 'X2' not in X.columns:
32.            raise ValueError("Both 'ItemID' and 'X2' columns must be present in the dataset.")
33.
34.        # Ensure no None values in ItemID
35.        if X['ItemID'].isnull().any():
36.            raise ValueError("ItemID column contains None values.")
37.
38.        # Safely impute NaN X2s based on ItemID or global mean

```

```

39.         X['X2'] = X.apply(
40.             lambda row: self.item_X2_mode.get(row['ItemID'], self.global_mean)
41.             if pd.isnull(row['X2'])
42.             else row['X2'],
43.             axis=1
44.         )
45.
46.     X['X2'] = X['X2'].fillna(self.global_mean)
47.     X['X2'] = X['X2'].astype(float)
48.
49.     return X[['X2']]
50.

```

```

1. class KMeansClusterTransformer(BaseEstimator, TransformerMixin):
2.     def __init__(self, num_clusters=5, random_state=42):
3.         """
4.         Parameters:
5.         - num_clusters: Number of clusters for KMeans.
6.         - random_state: Random state for reproducibility.
7.         """
8.         self.num_clusters = num_clusters
9.         self.random_state = random_state
10.        self.kmeans = None
11.
12.    def fit(self, X, y=None):
13.        # Ensure X is a DataFrame
14.        if not isinstance(X, pd.DataFrame):
15.            X = pd.DataFrame(X, columns=['Feature1', 'Feature2'])
16.
17.        # Fit the KMeans model
18.        self.kmeans = KMeans(n_clusters=self.num_clusters, random_state=self.random_state)
19.        self.kmeans.fit(X)
20.        return self
21.
22.    def transform(self, X):
23.        # Ensure X is a DataFrame
24.        if not isinstance(X, pd.DataFrame):
25.            X = pd.DataFrame(X, columns=['Feature1', 'Feature2'])
26.
27.        # Ensure the KMeans model is fitted
28.        if self.kmeans is None:
29.            raise ValueError("The KMeans model is not fitted yet. Please call 'fit' with appropriate arguments
before using this method.")
30.
31.        # Predict cluster labels
32.        cluster_labels = self.kmeans.predict(X)
33.
34.        # Return only the cluster labels as a DataFrame
35.        return pd.DataFrame(cluster_labels, columns=X.columns, index=X.index)
36.

```

```

1. class CustomFeatureTransformer(BaseEstimator, TransformerMixin):
2.     def __init__(self, column_names):
3.         self.column_names = column_names
4.         # self.scaler_price_per_unit_weight = StandardScaler()
5.         # self.scaler_X2 = StandardScaler()
6.         # self.scaler_pricing_strategy = StandardScaler()

```

```

7.
8.     def fit(self, X, y=None):
9.         # Convert to DataFrame with updated column names
10.        X = pd.DataFrame(X, columns=self.column_names)
11.        # Create new features
12.        # X['price_per_unit_weight'] = X['TypeMeanPrice_0'] / (X['X2'] + 0.0001)
13.        # X['pricing_strategy'] = X['MRP'] - (X['TypeMeanPrice_0'] * (X['Outlet_TypeOrdinal'] + 1) *
(X['Location_Type'] + 1))
14.
15.        # Fit the scalers on the training data
16.        # self.scaler_price_per_unit_weight.fit(X[['price_per_unit_weight']])
17.        # self.scaler_X2.fit(X[['X2']])
18.        # self.scaler_pricing_strategy.fit(X[['pricing_strategy']])
19.
20.        return self
21.
22.    def transform(self, X):
23.        # Convert to DataFrame with updated column names
24.        X = pd.DataFrame(X, columns=self.column_names)
25.        #X['OutletID'] = OutletID
26.        # X['OutletID'] = X['OutletID'].astype('category')
27.
28.        #X['Total'] = Total_X
29.        # Convert object type features to float
30.        for col in X.select_dtypes(include=['object']).columns:
31.            X[col] = X[col].astype(float)
32.        # Create new features
33.        #X['BigMac_index2'] = (X['VisibilityCont'] + 0.0001) * (X['OutletType'] + 1) * (X['TypeMeanPrice'] +
3.24995) * (X['LocationType'] + 1)
34.        #X['price_per_unit_weight'] = X['TypeMeanPrice'] / (X['dummy'] + 0.0001)
35.        #X['BigMac_index2'] = np.sqrt(X['BigMac_index2'])
36.        # X['price_per_unit_weight'] = np.log1p(X['price_per_unit_weight'])
37.        #X['Visibility'] = X['Visibility'].rank(ascending=False)
38.        #X['ItemCategory'] = ItemCategory
39.        #X['OutletID'] = OutletID.astype(str)
40.        # X['MRP_Age_Interaction'] = X['MRP'] / (2010 - X['EstablishmentYear'])
41.        # X['MRP_Age_Interaction'] = np.log1p(X['MRP_Age_Interaction'])
42.        X['EstablishmentYear'] = X['EstablishmentYear'].astype(str)
43.        X.MRP_cluster = X.MRP_cluster.astype(str)
44.        # X.Size = X.Size.astype(str)
45.        # X['ItemType'] = ItemType.astype(str)
46.        #X.ItemCategory = X.ItemCategory.astype('category')
47.        # X.FatContent = X.FatContent.astype(str)
48.        #X.X2 = X.X2.astype('category')
49.        #X.Visibility = X.Visibility.astype('category')
50.        #X.LocationType = X.LocationType.astype('category')
51.        # X['pricing_strategy'] = X['MRP'] - (X['TypeMeanPrice_0'] * (X['Outlet_TypeOrdinal'] + 1) *
(X['Location_Type'] + 1))
52.        # X['RegionalPotential'] = X['Location_Type'] * (X['PastYearsCont'] + 1)
53.        # # Standard scale the new features
54.        # X['price_per_unit_weight'] =
self.scaler_price_per_unit_weight.transform(X[['price_per_unit_weight']])
55.        # X['X2'] = self.scaler_X2.transform(X[['X2']])
56.        # X['pricing_strategy'] = self.scaler_pricing_strategy.transform(X[['pricing_strategy']])
57.        # Drop unwanted columns
58.        drop_columns =
['LocationType', 'Visibility', 'Size', 'FatContent', 'ItemType', 'MRP_Age_Interaction', 'EstablishmentYear']
59.        X.drop(columns=[col for col in drop_columns if col in X.columns], inplace=True)
60.
61.

```

```
62.         return X
63.
```

```
1. class FrequencyImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.freqs = {}
4.
5.     def fit(self, X, y=None):
6.         # Compute the frequency (count) of each category in each column
7.         if isinstance(X, pd.DataFrame):
8.             self.freqs = {col: X[col].value_counts() for col in X.columns}
9.         else:
10.            raise ValueError("Input must be a pandas DataFrame.")
11.        return self
12.
13.    def transform(self, X):
14.        # Replace missing values with the frequency of each category
15.        X = X.copy()
16.        for col in X.columns:
17.            freq = self.freqs.get(col, None)
18.            if freq is not None:
19.                # Replace NaN values with the frequency count
20.                X[col] = X[col].apply(lambda x: freq[x] if pd.notna(x) else freq.idxmax() if pd.isna(x) else x)
21.        return X
22.
```

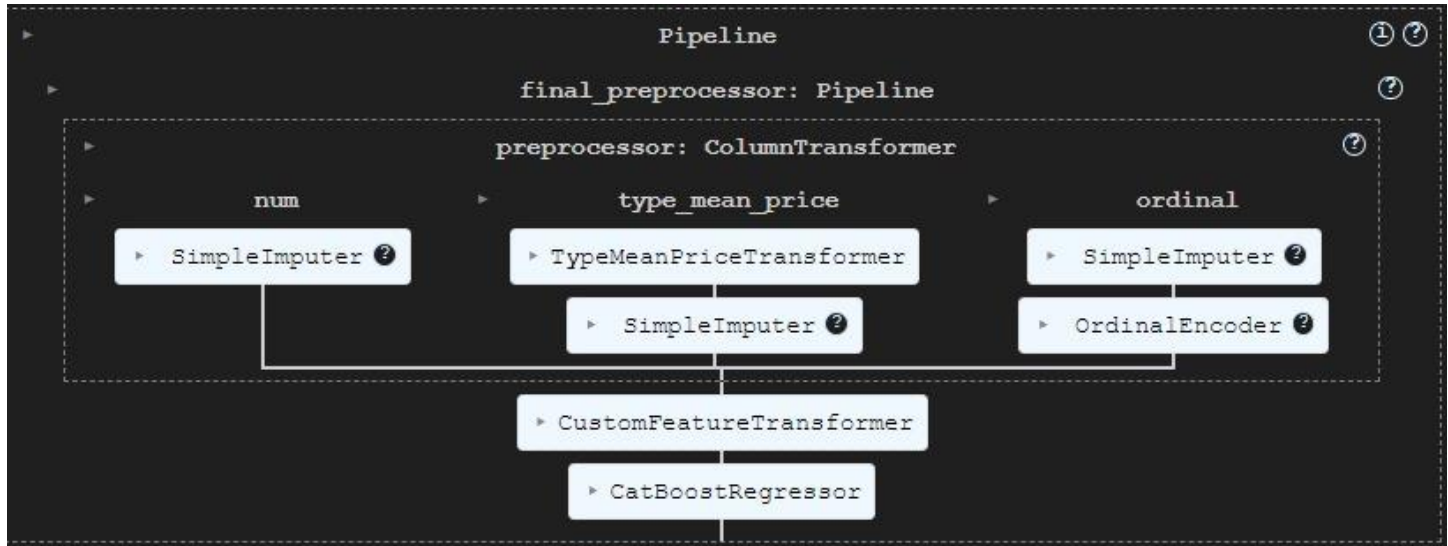
```
1. from sklearn.pipeline import FunctionTransformer
2. all_features = [ 'EstablishmentYear', 'MRP',
'Visibility', 'TypeMeanPrice', 'OutletType', 'MRP_cluster', 'LocationType', 'FatContent', 'Size' ]
3. target_encoder_cols = ['Type']
4. numerical_cols = ['EstablishmentYear', 'MRP', 'Visibility']
5. VisibilityZeros_cols = ['ItemID', 'Visibility']
6. X2_imputer_cols = ['X2', 'ItemID']
7. ordinal_cols = [ "Outlet_Type", 'MRP_cluster', 'Location_Type', 'FatContent', 'Size' ]
8. TypeMeanPriceTransformer_col = ['Outlet_Type', 'Type']
9. freq_imputer = ['Type'] # only one-hot encoding
10. ordinal_categories = [ # FatContent
11.     ['Grocery Store', 'Supermarket Type1', 'Supermarket Type2', 'Supermarket Type3'],
12.     ['very low', 'low', 'high', 'very high'],
13.     ['Tier 1', 'Tier 3', 'Tier 2'],
14.     ['nofat', 'LF', 'REG'],
15.     ['Small', 'Medium', 'High']
16.     # Size
17. ]
18. # Define frequency imputer pipeline
19. frequency_imputer_pipeline = Pipeline([
20.     ("freq_imputer", FrequencyImputer())
21. ])
22. # Define numerical pipeline
23. numerical_pipeline = Pipeline([
24.     ("imputer", SimpleImputer(strategy="most_frequent"))
25.     #, ("scaler", StandardScaler())
26. ])
27.
28. # size_imputer_pipeline = Pipeline([
29.     ("custom_size_imputer", SizeImputer()),
30.     ("imputer", SimpleImputer(strategy="most_frequent")),
```

```

31. #      ("ordinal", OrdinalEncoder(categories=[['Small', 'Medium', 'High']]))
32. # ])
33. # Define pipeline for X2 with X2 imputer and KMeans transformer
34. x2_pipeline = Pipeline([
35.     ("X2_imputer", X2NaNsImputer()),
36.     ("kmeans", KMeansClusterTransformer(num_clusters=5, random_state=42))
37. ])
38. type_mean_price_pipeline = Pipeline([
39.     ("type_mean_price_transformer", TypeMeanPriceTransformer(smoothing=0.5)),
40.     ("imputer", SimpleImputer(strategy="mean"))
41. ])
42.
43. # Define one-hot encoding pipeline
44. ordinal_pipeline = Pipeline([
45.     ("imputer", SimpleImputer(strategy="most_frequent")),
46.     ("ordinal", OrdinalEncoder(categories=ordinal_categories))
47. ])
48.
49. # Define one-hot encoding pipeline for onehot-transform columns
50. onehot_transform_pipeline = Pipeline([
51.     ("imputer", SimpleImputer(strategy="most_frequent")),
52.     ("onehot", OneHotEncoder(categories=[['Snack Foods', 'Frozen Foods', 'Fruits and Vegetables', 'Canned']],
handle_unknown="ignore"))
53. ])
54. # Define target encoding pipeline
55. target_encoding_pipeline = Pipeline([
56.     ("target_encoder", TargetEncoder(cols=target_encoder_cols , smoothing=0.5) )
57. ])
58.
59. # Combine all pipelines into a ColumnTransformer
60. preprocessor = ColumnTransformer([
61.     #("X2_imputer", X2NaNsImputer(), X2_imputer_cols),
62.     #("kmeans_X2", KMeansClusterTransformer(num_clusters=5, random_state=42), ['X2']),
63.     #("kmeans_visibility", KMeansClusterTransformer(num_clusters=5, random_state=42), ['Visibility']),
64.     ("num", numerical_pipeline, numerical_cols),
65.     ("type_mean_price", type_mean_price_pipeline, TypeMeanPriceTransformer_col),
66.     ("ordinal", ordinal_pipeline, ordinal_cols)
67.     #("freq_imputer", frequency_imputer_pipeline, ['TypeFreq'])
68. ])
69. ])
70. # Define the final pipeline
71. # Create the extended pipeline
72. final_pipeline = Pipeline([
73.     ("preprocessor", preprocessor),
74.     ("custom_features", CustomFeatureTransformer(column_names=all_features))
75.     #("kmeans", KMeansTransformer(num_clusters=5, random_state=42))
76. ])

```

77.



```

1. import optuna
2. from catboost import CatBoostRegressor
3. from sklearn.pipeline import Pipeline
4. from sklearn.model_selection import cross_val_score
5.
6. def objective_catboost(trial):
7.     param = {
8.         'iterations': trial.suggest_int('iterations', 100, 800),
9.         'learning_rate': trial.suggest_float('learning_rate', 0.01, 0.1),
10.        'depth': trial.suggest_int('depth', 3, 6),
11.        'subsample': trial.suggest_float('subsample', 0.5, 0.9),
12.        'colsample_bylevel': trial.suggest_float('colsample_bylevel', 0.5, 0.9),
13.        'l2_leaf_reg': trial.suggest_float('l2_leaf_reg', 1, 10),
14.        'random_strength': trial.suggest_float('random_strength', 1, 5),
15.        'bagging_temperature': trial.suggest_float('bagging_temperature', 1, 5),
16.        'border_count': trial.suggest_int('border_count', 32, 128),
17.        'loss_function': 'MAE'
18.    }
19.
20.    # Convert categorical columns to string
21.    # X_trans_catboost = X_trans.copy()
22.    # for col in X_trans_catboost.select_dtypes(include=['category']).columns:
23.    #     X_trans_catboost[col] = X_trans_catboost[col].astype(str)
24.    smoothing = trial.suggest_float('smoothing', 0.1, 5)
25.    preprocessor.set_params(
26.        type_mean_price__type_mean_price_transformer__smoothing=smoothing
27.    )
28.    # Create a pipeline with the final transformer and CatBoost regressor

```

```

29.     model = CatBoostRegressor(**param, verbose=0, cat_features=['MRP_cluster'], early_stopping_rounds=50)
30.     pipeline = Pipeline([
31.         ('final_transform', final_pipeline),
32.         ('catboost', model)
33.     ])
34.
35.     # Perform cross-validation
36.     scores = cross_val_score(pipeline, X, Y, cv=5, scoring='neg_mean_absolute_error')
37.     mae = -scores.mean()
38.     return mae
39.
40. # Create an Optuna study and optimize the objective function
41. study_catboost = optuna.create_study(direction='minimize')
42. study_catboost.optimize(objective_catboost, n_trials=10)
43.
44. # Print the best parameters and MAE
45. print("CatBoost Best Parameters:", study_catboost.best_params)
46. print("CatBoost Best MAE:", study_catboost.best_value)
47.

```

```

1. # Fit the pipeline
2. cat_pipeline.fit(X_train, Y_train)
3.
4. # Transform the training data
5. X_train_transformed = cat_pipeline.named_steps['final_preprocessor'].transform(X_train)
6.
7. # Calculate feature importances
8. mdi_importances = pd.Series(cat_pipeline.named_steps['catboost'].feature_importances_,
index=X_train_transformed.columns)
9. tree_importance_sorted_idx = np.argsort(cat_pipeline.named_steps['catboost'].feature_importances_)
10.
11. fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 8))
12.
13. # Plot Gini importance
14. mdi_importances.sort_values().plot.barh(ax=ax1)
15. ax1.set_xlabel("Gini importance")
16.
17. # Plot permutation importance
18. plot_permutation_importance(cat_pipeline.named_steps['catboost'], X_train_transformed, Y_train, ax2)
19. ax2.set_xlabel("Decrease in accuracy score")
20.
21. fig.suptitle("Impurity-based vs. permutation importances on multicollinear features (train set)")
22. fig.tight_layout()
23. plt.show()
24.

```

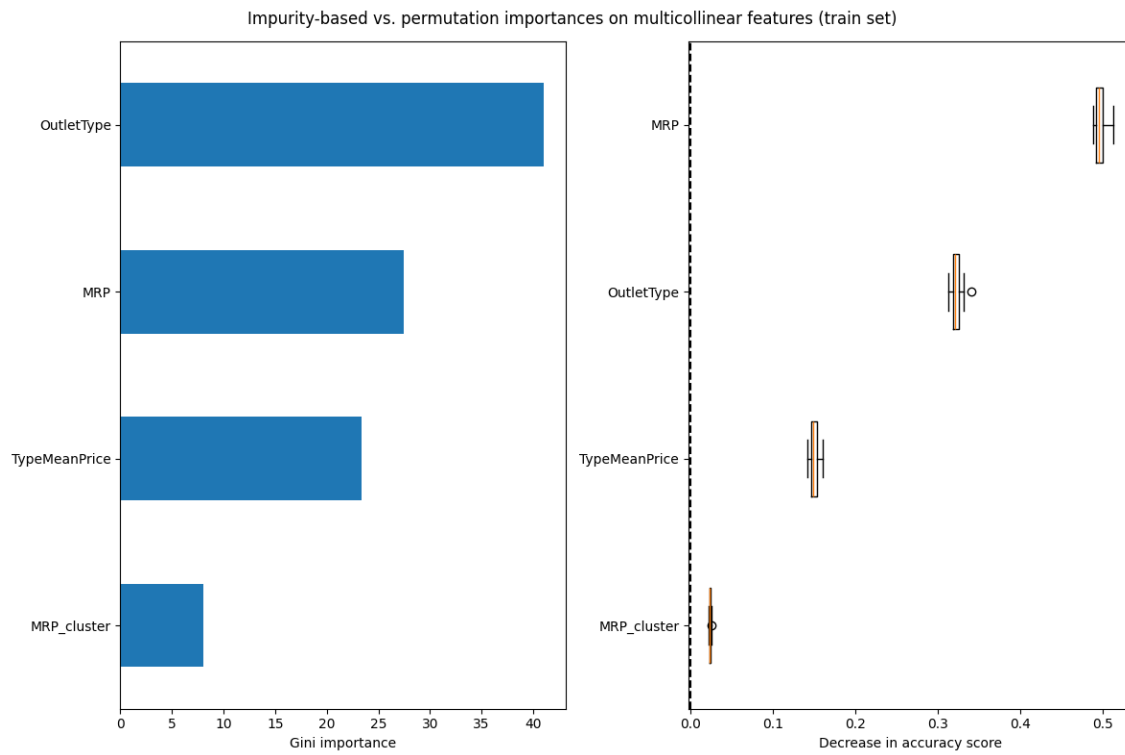


Figure 8 Impurity-based vs. permutation importances on multicollinear features



## 7.3 Linear Model

```

1. class ItemWeightImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.item_weight_mode = defaultdict(lambda: None) # Stores mode of Item_Weight per Item_Identifier
4.         self.global_mean = None # Fallback global mean for Item_Weight
5.
6.     def fit(self, X, y=None):
7.         X = X.copy()
8.
9.         # Ensure the required columns exist
10.        if 'Item_Identifier' not in X.columns or 'Item_Weight' not in X.columns:
11.            raise ValueError("Both 'Item_Identifier' and 'Item_Weight' columns must be present in the
dataset.")
12.
13.        # Ensure no None values in Item_Identifier
14.        if X['Item_Identifier'].isnull().any():
15.            raise ValueError("Item_Identifier column contains None values.")
16.
17.        # Calculate mode Item_Weight for each Item_Identifier
18.        item_weight_mode = X.groupby('Item_Identifier')['Item_Weight'].agg(lambda x: x.mode().iloc[0] if not
x.mode().empty else np.nan)
19.        self.item_weight_mode.update(item_weight_mode.to_dict())
20.
21.        # Calculate global mean for the Item_Weight column
22.        if X['Item_Weight'].notnull().any():
23.            self.global_mean = X['Item_Weight'].mean()
24.        else:
25.            raise ValueError("Item_Weight column contains only NaN values.")
26.
27.        return self
28.
29.    def transform(self, X):
30.        X = X.copy()
31.
32.        # Ensure the required columns exist
33.        if 'Item_Identifier' not in X.columns or 'Item_Weight' not in X.columns:
34.            raise ValueError("Both 'Item_Identifier' and 'Item_Weight' columns must be present in the
dataset.")
35.
36.        # Ensure no None values in Item_Identifier
37.        if X['Item_Identifier'].isnull().any():
38.            raise ValueError("Item_Identifier column contains None values.")
39.
40.        # Safely impute NaN Item_Weights based on Item_Identifier or global mean
41.        X['Item_Weight'] = X.apply(
42.            lambda row: self.item_weight_mode.get(row['Item_Identifier'], self.global_mean)
43.            if pd.isnull(row['Item_Weight'])
44.            else row['Item_Weight'],
45.            axis=1
46.        )
47.
48.        X['Item_Weight'] = X['Item_Weight'].fillna(self.global_mean)
49.
50.        return X[['Item_Weight']]

```

```

1. X['Category'] = X['Item_Identifier'].str[:2]
2. X["Item_Fat_Content"] = X["Item_Fat_Content"].replace({"low fat": "LF", "Low Fat": "LF", "Regular": "REG",
"reg": "REG"})
3. X.loc[X['Category'] == 'NC', 'Item_Fat_Content'] = 'nofat'
4. impute_sizes = {
5.     "OUT010": "Small",
6.     "OUT017": "Small",
7.     "OUT045": "Medium"
8. }
9. X['Outlet_Size'] = X.apply(
10.     lambda row: impute_sizes[row['Outlet_Identifier']] if pd.isnull(row['Outlet_Size']) else
row['Outlet_Size'], axis=1
11. )
12. X['MRP_cluster']=pd.cut(X["Item_MRP"],bins=[25,69,137,203,270],labels=['very low','low','high','very
high'],right=True)
13. weightimputer = ItemWeightImputer()
14. X["Item_Weight"] = weightimputer.fit_transform(X)
15. X["Weight_per_Unit_MRP"] = X["Item_Weight"]/X["Item_MRP"]
16.

```

```

1. #target_encoder_cols = ['Item_Type']
2. numerical_cols = ["Item_Weight","Weight_per_Unit_MRP","Item_Visibility"]
3. # ordinal_cols=["Outlet_Size"]
4. one_hot_columns = ["Outlet_Location_Type","Category","Item_Fat_Content","Outlet_Size","Outlet_Type"]
5. target_encoder_cols = ["Item_Identifier","Item_Type","Outlet_Identifier"]
6. ordinal_categories = [['very low', 'low', 'high', 'very high']]
7. # Numerical pipeline
8. # numerical_pipeline = Pipeline([
9. #     ("group_mean_imputer", FillNaWithGroupMode(group_col="Item_Identifier", target_col="Item_Weight")),
10. #     ("mean_imputer", SimpleImputer(strategy="mean")),
11. #     ("scaler", StandardScaler())
12. # ])
13. numerical_pipeline = Pipeline([
14.     # ("group_mode_imputer", FillNaWithGroupMode(group_col=group_col, target_col='Item_Weight')),
15.     ("mean_imputer", SimpleImputer(strategy='mean')),
16.     ("scaler", StandardScaler())
17. ])
18.
19. #Ordinal pipeline
20. # ordinal_pipeline = Pipeline([
21. #     ("imputer", SimpleImputer(strategy="most_frequent")),
22. #     ("ordinal", OrdinalEncoder(categories=ordinal_categories))
23. # ])
24. target_encoding_pipeline = Pipeline([
25.     ("target_encoder", TargetEncoder(cols=target_encoder_cols , smoothing=0.5) )
26. ])
27. # One-hot encoding pipeline
28. onehot_transform_pipeline = Pipeline([
29.     ("imputer", SimpleImputer(strategy="most_frequent")),
30.     ("onehot", OneHotEncoder(handle_unknown="ignore", sparse_output=False))
31. ])
32. ordinal_pipeline = Pipeline([
33.     ("imputer", SimpleImputer(strategy="most_frequent")),
34.     ("ordinal", OrdinalEncoder(categories=ordinal_categories))
35. ])
36.
37. # Target encoding pipeline
38. # target_encoding_pipeline = Pipeline([

```

```

39. #      ("target_encoder", TargetEncoder(cols=target_encoder_cols))
40. # ])
41.
42. # Combine all pipelines into a ColumnTransformer
43. preprocessor = ColumnTransformer([
44.     ("num", numerical_pipeline, numerical_cols),
45.     ("target_encoder", target_encoding_pipeline, target_encoder_cols),
46.     ("onehot", onehot_transform_pipeline, one_hot_columns),
47.     ("ordinal", ordinal_pipeline, ["MRP_cluster"])
48. ])
49. ])
50.
51. # Final pipeline
52. final_pipeline = Pipeline([
53.     ("preprocessor", preprocessor),
54. ])
55.

```

```

1. import optuna
2. from sklearn.linear_model import LinearRegression
3. from sklearn.pipeline import Pipeline
4. from sklearn.model_selection import cross_val_score
5.
6. def objective_linear_regression(trial):
7.     # LinearRegression does not have many hyperparameters to tune, but we can still create a pipeline
8.     param = {
9.         'fit_intercept': trial.suggest_categorical('fit_intercept', [True, False])
10.    }
11.
12.    # Create a pipeline with the final transformer and LinearRegression
13.    model = LinearRegression(**param)
14.    pipeline = Pipeline([
15.        ('final_transform', final_pipeline),
16.        ('linear_regression', model)
17.    ])
18.
19.    # Perform cross-validation with error handling
20.    try:
21.        scores = cross_val_score(pipeline, X, Y, cv=5, scoring='neg_mean_absolute_error', error_score='raise')
22.        mae = -scores.mean()
23.        print(f"Trial MAE: {mae}")
24.    except Exception as e:
25.        print(f"Error during cross-validation: {e}")
26.        mae = float('inf') # Assign a high error value if an exception occurs
27.
28.    return mae
29.
30. study_linear_regression = optuna.create_study(direction='minimize')
31. study_linear_regression.optimize(objective_linear_regression, n_trials=50)
32.
33. print("Linear Regression Best Parameters:", study_linear_regression.best_params)
34. print("Linear Regression Best MAE:", study_linear_regression.best_value)
35.
36. # Extract the best model
37. best_params = study_linear_regression.best_params
38. best_model = LinearRegression(**best_params)
39. best_pipeline = Pipeline([
40.     ('final_transform', final_pipeline),
41.     ('linear_regression', best_model)

```

```

42. ])
43.
44. # Fit the best pipeline on the entire dataset
45. best_pipeline.fit(X, Y)
46.
47. # Now you can use best_pipeline for predictions
48. predictions = best_pipeline.predict(X)
49. print(predictions)
50.

```

```

1. # Get the feature names from the preprocessor
2. feature_names = final_pipeline.named_steps['preprocessor'].get_feature_names_out()
3.
4. # Get the coefficients from the linear regression model
5. coefficients = best_model.coef_
6.
7. # Create a DataFrame to hold feature names and their corresponding coefficients
8. feature_importance = pd.DataFrame({
9.     'Feature': feature_names,
10.    'Coefficient': coefficients
11. })
12.
13. # Sort the DataFrame by the absolute value of the coefficients in descending order
14. feature_importance['Absolute_Coefficient'] = feature_importance['Coefficient'].abs()
15. feature_importance = feature_importance.sort_values(by='Absolute_Coefficient', ascending=False)
16.
17. # Print the top 10 most important features
18. print(feature_importance.head(10))
19.

```

```

1. # Get the feature names from the preprocessor
2. feature_names = final_pipeline.named_steps['preprocessor'].get_feature_names_out()
3.
4. # Get the coefficients from the linear regression model
5. coefficients = best_model.coef_
6.
7. # Create a DataFrame to hold feature names and their corresponding coefficients
8. feature_importance = pd.DataFrame({
9.     'Feature': feature_names,
10.    'Coefficient': coefficients
11. })
12.
13. # Sort the DataFrame by the absolute value of the coefficients in descending order
14. feature_importance['Absolute_Coefficient'] = feature_importance['Coefficient'].abs()
15. feature_importance = feature_importance.sort_values(by='Absolute_Coefficient', ascending=False)
16.
17. # Print the top 10 most important features
18. print(feature_importance.head(10))
19.

```

	Feature	Coefficient	Absolute_Coefficient	
2.	8	onehot__Outlet_Location_Type_Tier 3	-8.764891e+09	8.764891e+09
3.	7	onehot__Outlet_Location_Type_Tier 2	-8.764891e+09	8.764891e+09
4.	6	onehot__Outlet_Location_Type_Tier 1	-8.764891e+09	8.764891e+09
5.	12	onehot__Item_Fat_Content_LF	-8.764735e+09	8.764735e+09

6.	13	onehot__Item_Fat_Content_REG	-8.764735e+09	8.764735e+09
7.	11	onehot__Category_NC	-8.764199e+09	8.764199e+09
8.	14	onehot__Item_Fat_Content_nofat	-8.763889e+09	8.763889e+09
9.	10	onehot__Category_FD	-8.763352e+09	8.763352e+09
10.	9	onehot__Category_DR	-8.763352e+09	8.763352e+09
11.	17	onehot__Outlet_Size_Small	-8.762929e+09	8.762929e+09
12.				

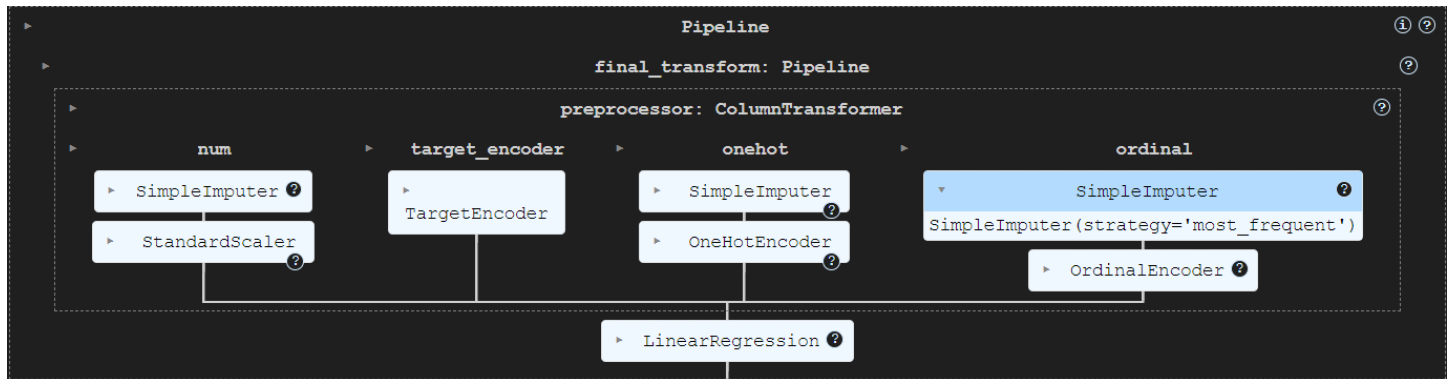


Figure 9 Linear Regression Pipeline

## 7.4 SVR Model

```

1. class ItemWeightImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.item_weight_mode = defaultdict(lambda: None) # Stores mode of Item_Weight per Item_Identifier
4.         self.global_mean = None # Fallback global mean for Item_Weight
5.
6.     def fit(self, X, y=None):
7.         X = X.copy()
8.
9.         # Ensure the required columns exist
10.        if 'Item_Identifier' not in X.columns or 'Item_Weight' not in X.columns:
11.            raise ValueError("Both 'Item_Identifier' and 'Item_Weight' columns must be present in the
dataset.")
12.
13.        # Ensure no None values in Item_Identifier
14.        if X['Item_Identifier'].isnull().any():
15.            raise ValueError("Item_Identifier column contains None values.")
16.
17.        # Calculate mode Item_Weight for each Item_Identifier
18.        item_weight_mode = X.groupby('Item_Identifier')['Item_Weight'].agg(lambda x: x.mode().iloc[0] if not
x.mode().empty else np.nan)
19.        self.item_weight_mode.update(item_weight_mode.to_dict())
20.
21.        # Calculate global mean for the Item_Weight column
22.        if X['Item_Weight'].notnull().any():
23.            self.global_mean = X['Item_Weight'].mean()
24.        else:
25.            raise ValueError("Item_Weight column contains only NaN values.")
26.
27.        return self
28.

```

```

29.     def transform(self, X):
30.         X = X.copy()
31.
32.         # Ensure the required columns exist
33.         if 'Item_Identifier' not in X.columns or 'Item_Weight' not in X.columns:
34.             raise ValueError("Both 'Item_Identifier' and 'Item_Weight' columns must be present in the
dataset.")
35.
36.         # Ensure no None values in Item_Identifier
37.         if X['Item_Identifier'].isnull().any():
38.             raise ValueError("Item_Identifier column contains None values.")
39.
40.         # Safely impute NaN Item_Weights based on Item_Identifier or global mean
41.         X['Item_Weight'] = X.apply(
42.             lambda row: self.item_weight_mode.get(row['Item_Identifier'], self.global_mean)
43.             if pd.isnull(row['Item_Weight'])
44.             else row['Item_Weight'],
45.             axis=1
46.         )
47.
48.         X['Item_Weight'] = X['Item_Weight'].fillna(self.global_mean)
49.
50.         return X[['Item_Weight']]
51.

```

```

1. class VisibilityZerosImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.item_visibility_mean = defaultdict(lambda: None) # Stores mean of Item_Visibility per
Item_Identifier
4.         self.global_mean = None # Fallback global mean for Item_Visibility
5.
6.     def fit(self, X, y=None):
7.         X = X.copy()
8.
9.         # Ensure the required columns exist
10.        if 'Item_Identifier' not in X.columns or 'Item_Visibility' not in X.columns:
11.            raise ValueError("Both 'Item_Identifier' and 'Item_Visibility' columns must be present in the
dataset.")
12.
13.        # Ensure no None values in Item_Identifier
14.        if X['Item_Identifier'].isnull().any():
15.            raise ValueError("Item_Identifier column contains None values.")
16.
17.        # Calculate mean Item_Visibility for each Item_Identifier
18.        item_visibility_mean = X.groupby('Item_Identifier')['Item_Visibility'].mean()
19.        self.item_visibility_mean.update(item_visibility_mean.to_dict())
20.
21.        # Calculate global mean for the Item_Visibility column
22.        if X['Item_Visibility'].notnull().any():
23.            self.global_mean = X['Item_Visibility'].mean()
24.        else:
25.            raise ValueError("Item_Visibility column contains only NaN values.")
26.
27.        return self
28.
29.     def transform(self, X):
30.         X = X.copy()
31.
32.         # Ensure the required columns exist

```

```

33.         if 'Item_Identifier' not in X.columns or 'Item_Visibility' not in X.columns:
34.             raise ValueError("Both 'Item_Identifier' and 'Item_Visibility' columns must be present in the
dataset.")
35.
36.         # Ensure no None values in Item_Identifier
37.         if X['Item_Identifier'].isnull().any():
38.             raise ValueError("Item_Identifier column contains None values.")
39.
40.         # Safely impute zero Item_Visibility based on Item_Identifier or global mean
41.         X['Item_Visibility'] = X.apply(
42.             lambda row: self.item_visibility_mean.get(row['Item_Identifier'], self.global_mean)
43.             if row['Item_Visibility'] == 0
44.             else row['Item_Visibility'],
45.             axis=1
46.         )
47.
48.         X['Item_Visibility'] = X['Item_Visibility'].replace(0, self.global_mean)
49.         X['Item_Visibility'] = np.sqrt(X['Item_Visibility'])
50.
51.         return X[['Item_Visibility']]
52.

```

```

1. X['Item_Type_Outlet_Type'] = X['Item_Type'] + "_" + X['Outlet_Type']
2. X['Item_Fat_Content_Outlet_Type'] = X['Item_Fat_Content'] + "_" + X['Outlet_Type']
3. X['Outlet_Size_Outlet_Location_Type'] = X['Outlet_Size'] + "_" + X['Outlet_Location_Type']
4. X['Item_Type_Item_Fat_Content'] = X['Item_Type'] + "_" + X['Item_Fat_Content']
5. X['Outlet_Type_Outlet_Location_Type'] = X['Outlet_Type'] + "_" + X['Outlet_Location_Type']
6.
7. # Numerical x Numerical Interactions
8. X['Item_MRP_Outlet_Establishment_Year'] = X['Item_MRP'] * X['Outlet_Establishment_Year']
9. X['Item_MRP_Item_Visibility'] = X['Item_MRP'] * X['Item_Visibility']
10. X['Years_Operating_Outlet_Type'] = X['age'] * X['Outlet_Type'].apply(lambda x: hash(x) % 1000)
11. X['Outlet_Identifier_Years_Operating'] = X['Outlet_Identifier'].apply(lambda x: hash(x) % 1000) * X['age']
12.
13. # Numerical x Categorical Interactions
14. X['Item_Visibility_Item_Type'] = X['Item_Visibility'] * X['Item_Type'].apply(lambda x: hash(x) % 1000)
15.

```

```

1. # Define the columns for each type of transformation
2. numerical_cols = ["Item_Weight", "Weight_per_Unit_MRP", "Item_Visibility",
"Item_MRP_Outlet_Establishment_Year", "Item_MRP_Item_Visibility", "Years_Operating_Outlet_Type",
"Outlet_Identifier_Years_Operating", "Item_Visibility_Item_Type"]
3. one_hot_columns = ["Outlet_Location_Type", "Category", "Item_Fat_Content", "Outlet_Size", "Outlet_Type",
"Item_Type_Outlet_Type", "Item_Fat_Content_Outlet_Type", "Outlet_Size_Outlet_Location_Type",
"Item_Type_Item_Fat_Content", "Outlet_Type_Outlet_Location_Type"]
4. target_encoder_cols = ['Item_Identifier', 'Item_Type', 'Outlet_Identifier']
5. ordinal_categories = [['very low', 'low', 'high', 'very high']]
6.
7. # Numerical pipeline
8. numerical_pipeline = Pipeline([
9.     ("mean_imputer", SimpleImputer(strategy='mean')),
10.    ("scaler", StandardScaler())
11. ])
12.
13. #Target encoding pipeline
14. target_encoding_pipeline = Pipeline([
15.    ("target_encoder", TargetEncoder(cols=target_encoder_cols, smoothing=0.5))

```

```

16. ])
17.
18. # One-hot encoding pipeline
19. onehot_transform_pipeline = Pipeline([
20.     ("imputer", SimpleImputer(strategy="most_frequent")),
21.     ("onehot", OneHotEncoder(handle_unknown="ignore", sparse_output=False))
22. ])
23.
24. # Ordinal encoding pipeline
25. ordinal_pipeline = Pipeline([
26.     ("imputer", SimpleImputer(strategy="most_frequent")),
27.     ("ordinal", OrdinalEncoder(categories=ordinal_categories))
28. ])
29.
30. # Combine all pipelines into a ColumnTransformer
31. preprocessor = ColumnTransformer([
32.     ("num", numerical_pipeline, numerical_cols),
33.     ("target_encoder", target_encoding_pipeline, target_encoder_cols),
34.     ("onehot", onehot_transform_pipeline, one_hot_columns),
35.     ("ordinal", ordinal_pipeline, ["MRP_cluster"])
36. ])
37.
38. # Final pipeline
39. final_pipeline = Pipeline([
40.     ("preprocessor", preprocessor)

```

```

1. import optuna
2. from sklearn.svm import SVR
3. from sklearn.pipeline import Pipeline
4. from sklearn.model_selection import cross_val_score
5.
6. def objective_svr(trial):
7.     param = {
8.         'C': trial.suggest_float('C', 0.1, 100.0),
9.         'epsilon': trial.suggest_float('epsilon', 0.01, 1.0),
10.        'kernel': trial.suggest_categorical('kernel', ['linear', 'poly', 'rbf', 'sigmoid']),
11.        'degree': trial.suggest_int('degree', 2, 5) if trial.suggest_categorical('kernel', ['linear', 'poly',
12.        'rbf', 'sigmoid']) == 'poly' else 3,
13.        'gamma': trial.suggest_categorical('gamma', ['scale', 'auto'])
14.    }
15.    # Create a pipeline with the final transformer and SVR
16.    model = SVR(**param)
17.    pipeline = Pipeline([
18.        ('final_transform', final_pipeline),
19.        ('svr', model)
20.    ])
21.
22.    # Perform cross-validation with error handling
23.    try:
24.        scores = cross_val_score(pipeline, X, Y, cv=5, scoring='neg_mean_absolute_error', error_score='raise')
25.        mae = -scores.mean()
26.        print(f"Trial MAE: {mae}")
27.    except Exception as e:
28.        print(f"Error during cross-validation: {e}")
29.        mae = float('inf') # Assign a high error value if an exception occurs
30.
31.    return mae
32.

```



```
33. study_svr = optuna.create_study(direction='minimize')
34. study_svr.optimize(objective_svr, n_trials=50)
35.
36. print("SVR Best Parameters:", study_svr.best_params)
37. print("SVR Best MAE:", study_svr.best_value)
```

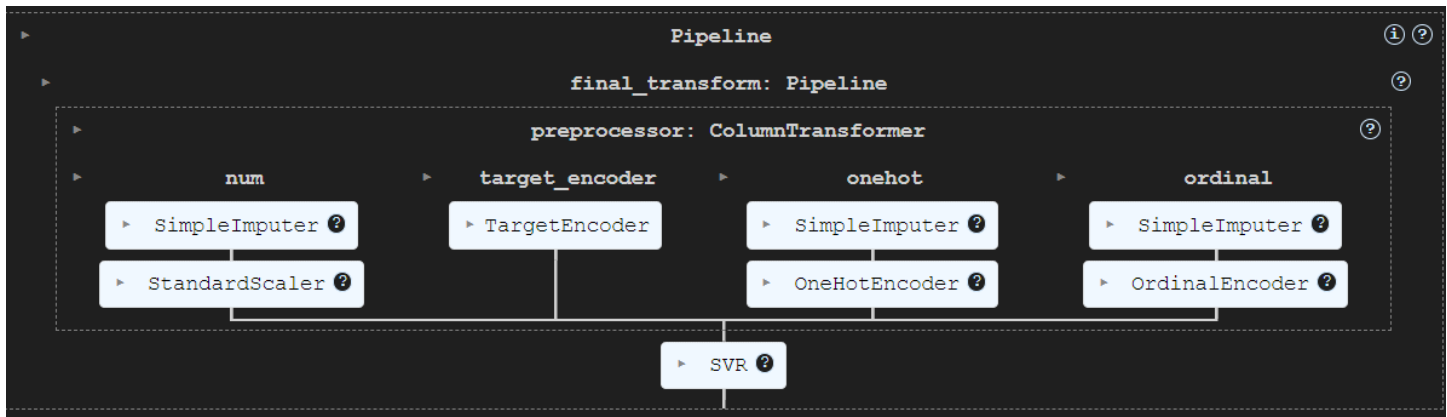


Figure 10 SVR Pipeline

## 7.5 XGBoost Stacking Model

```

1. class TypeMeanPriceTransformer(BaseEstimator, TransformerMixin):
2.     def __init__(self, smoothing=1):
3.         self.smoothing = smoothing
4.         self.category_mean_price = None
5.
6.     def fit(self, X, y=None):
7.         if y is None:
8.             raise ValueError("Target variable `y` cannot be None in the fit method.")
9.         data_ = X.copy()
10.        data_['Y'] = y
11.        # Calculate the mean price for each OutletID and Type
12.        mean_price = data_.groupby(['OutletID', 'Type'])['Y'].mean()
13.        count = data_.groupby(['OutletID', 'Type'])['Y'].count()
14.        global_mean = data_['Y'].mean()
15.
16.        # Apply smoothing
17.        self.category_mean_price = ((mean_price * count) + (global_mean * self.smoothing)) / (count + self.smoothing)
18.        self.category_mean_price = self.category_mean_price.reset_index()
19.        self.category_mean_price.columns = ['OutletID', 'Type', 'TypeMeanPrice']
20.        return self
21.
22.    def transform(self, X):
23.        X = X.copy()
24.        # Merge the mean price with the original dataframe
25.        X = pd.merge(X, self.category_mean_price, on=['OutletID', 'Type'], how='left')
26.        X['TypeMeanPrice'] = np.log1p(X['TypeMeanPrice'])
27.        return X[['TypeMeanPrice']]
28.

```

```

1. class VisibilityZerosImputerZerosImputer(BaseEstimator, TransformerMixin):
2.     def __init__(self):
3.         self.item_visibility_mean = defaultdict(lambda: None) # Stores mean of Visibility per ItemID
4.         self.global_mean = None # Fallback global mean for Visibility
5.
6.     def fit(self, X, y=None):
7.         X = X.copy()
8.
9.         # Ensure the required columns exist
10.        if 'ItemID' not in X.columns or 'Visibility' not in X.columns:
11.            raise ValueError("Both 'ItemID' and 'Visibility' columns must be present in the dataset.")
12.
13.        # Ensure no None values in ItemID
14.        if X['ItemID'].isnull().any():
15.            raise ValueError("ItemID column contains None values.")
16.
17.        # Calculate mean Visibility for each ItemID
18.        item_visibility_mean = X.groupby('ItemID')['Visibility'].mean()
19.        self.item_visibility_mean.update(item_visibility_mean.to_dict())
20.
21.        # Calculate global mean for the Visibility column
22.        if X['Visibility'].notnull().any():
23.            self.global_mean = X['Visibility'].mean()
24.        else:
25.            raise ValueError("Visibility column contains only NaN values.")
26.
27.        return self
28.

```

```

29.     def transform(self, X):
30.         X = X.copy()
31.
32.         # Ensure the required columns exist
33.         if 'ItemID' not in X.columns or 'Visibility' not in X.columns:
34.             raise ValueError("Both 'ItemID' and 'Visibility' columns must be present in the dataset.")
35.
36.         # Ensure no None values in ItemID
37.         if X['ItemID'].isnull().any():
38.             raise ValueError("ItemID column contains None values.")
39.
40.         # Safely impute zero Visibilitys based on ItemID or global mean
41.         X['Visibility'] = X.apply(
42.             lambda row: self.item_Visibility_mean.get(row['ItemID'], self.global_mean)
43.             if row['Visibility'] == 0
44.             else row['Visibility'],
45.             axis=1
46.         )
47.
48.         X['Visibility'] = X['Visibility'].replace(0, self.global_mean)
49.         X['Visibility'] = np.sqrt(X['Visibility'])
50.
51.         return X[['Visibility']]
52.

```

```

1. class CustomFeatureTransformer(BaseEstimator, TransformerMixin):
2.     def __init__(self, column_names):
3.         self.column_names = column_names
4.         self.scaler_price_per_unit_weight = StandardScaler()
5.         self.scaler_X2 = StandardScaler()
6.         self.scaler_pricing_strategy = StandardScaler()
7.
8.     def fit(self, X, y=None):
9.         # Convert to DataFrame with updated column names
10.        X = pd.DataFrame(X, columns=self.column_names)
11.        # Create new features
12.        X['price_per_unit_weight'] = X['TypeMeanPrice_0'] / (X['X2'] + 0.0001)
13.        X['pricing_strategy'] = X['MRP'] - (X['TypeMeanPrice_0'] * (X['Outlet_TypeOrdinal'] + 1) *
(X['Location_Type'] + 1))
14.
15.        # Fit the scalers on the training data
16.        self.scaler_price_per_unit_weight.fit(X[['price_per_unit_weight']])
17.        self.scaler_X2.fit(X[['X2']])
18.        self.scaler_pricing_strategy.fit(X[['pricing_strategy']])
19.
20.        return self
21.
22.     def transform(self, X):
23.         # Convert to DataFrame with updated column names
24.        X = pd.DataFrame(X, columns=self.column_names)
25.
26.        # Create new features
27.        X['BigMac_index2'] = (X['Visibility'] + 0.0001) * (X['Outlet_TypeOrdinal'] + 1) * (X['TypeMeanPrice_0']
+ 3.24995) * (X['Location_Type'] + 1)
28.        X['price_per_unit_weight'] = X['TypeMeanPrice_0'] / (X['X2'] + 0.0001)
29.        X['BigMac_index2'] = np.sqrt(X['BigMac_index2'])
30.        X['pricing_strategy'] = X['MRP'] - (X['TypeMeanPrice_0'] * (X['Outlet_TypeOrdinal'] + 1) *
(X['Location_Type'] + 1))
31.        X['RegionalPotential'] = X['Location_Type'] * (X['PastYearsCont'] + 1)

```

```

32.
33.     # Standard scale the new features
34.     X['price_per_unit_weight'] = self.scaler_price_per_unit_weight.transform(X[['price_per_unit_weight']])
35.     X['X2'] = self.scaler_X2.transform(X[['X2']])
36.     X['pricing_strategy'] = self.scaler_pricing_strategy.transform(X[['pricing_strategy']])
37.
38.     # Drop unwanted columns
39.     drop_columns = [
40.         "Category_DR", "Category_FD", "Outlet_Type_Supermarket Type2", 'FatContent',
41.         "Outlet_Type_Supermarket Type1", "BigMac_index3", 'Size',
42.         "Visibility", "BigMac_index", "Outlet_Type_Grocery Store", "Outlet_Type_Supermarket Type3"
43.     ]
44.     X.drop(columns=[col for col in drop_columns if col in X.columns], inplace=True)
45.
46.     return X
47.

```

```

1. class KMeansTransformer(BaseEstimator, TransformerMixin):
2.     def __init__(self, num_clusters=5, random_state=42):
3.         self.num_clusters = num_clusters
4.         self.random_state = random_state
5.
6.
7.     def fit(self, X, y=None):
8.         #X = X.drop("price_per_unit_weight", axis=1)
9.         # Fit KMeans using the original features
10.        self.feature_names_in_ = X.columns # Store original column names
11.        self.kmeans = KMeans(n_clusters=self.num_clusters, random_state=self.random_state)
12.        self.kmeans.fit(X)
13.        return self
14.
15.    def transform(self, X):
16.        #X = X.drop("price_per_unit_weight", axis=1)
17.        # Ensure the input columns match those seen during fit
18.        if list(X.columns) != list(self.feature_names_in_):
19.            raise ValueError(
20.                "The feature names should match those that were passed during fit.\n"
21.                f"Feature names unseen at fit time: {set(X.columns) - set(self.feature_names_in_)}\n"
22.                f"Feature names seen at fit time but not in transform: {set(self.feature_names_in_) -"
23.                f"set(X.columns)}"
24.            )
25.
26.        # Compute cluster labels
27.        cluster_labels = self.kmeans.predict(X)
28.
29.        # Compute distances to centroids
30.        centroid_distances = self.kmeans.transform(X)
31.
32.        # Create a new DataFrame with cluster labels and centroid distances
33.        new_features = pd.DataFrame(
34.            centroid_distances,
35.            columns=[f"Centroid_{i}" for i in range(centroid_distances.shape[1])],
36.            index=X.index
37.        )
38.        new_features['Cluster'] = cluster_labels
39.
40.        # Concatenate the original DataFrame with the new features
41.        X_transformed = pd.concat([X.reset_index(drop=True), new_features.reset_index(drop=True)], axis=1)

```

---

```
41.         return X_transformed
```

---

```

1. all_features = ['EstablishmentYear', 'MRP',"Outlet_Type",'MRP_cluster','Snack Foods', 'Frozen Foods', 'Fruits
and Vegetables', 'Canned','X2','Visibility']
2. target_encoder_cols = ['Type']
3. numerical_cols = ['EstablishmentYear', 'MRP']
4. VisibilityZeros_cols = ['ItemID','Visibility']
5. X2_imputer_cols = ['X2','ItemID']
6. ordinal_cols = [ "Outlet_Type",'MRP_cluster']
7. TypeMeanPriceTransformer_col = ['OutletID' , 'Type']
8. one_hot_columns = ['Type'] # only one-hot encoding
9. ordinal_categories = [ # FatContent
10.     ['Grocery Store', 'Supermarket Type1', 'Supermarket Type2', 'Supermarket Type3'],
11.     ['very low','low','high','very high']
12.
13.     # Size
14. ]
15.
16. # Define numerical pipeline
17. numerical_pipeline = Pipeline([
18.     ("imputer", SimpleImputer(strategy="most_frequent")),
19.     #,("scaler", StandardScaler())
20. ])
21.
22. # size_imputer_pipeline = Pipeline([
23. #     ("custom_size_imputer", SizeImputer()),
24. #     ("imputer", SimpleImputer(strategy="most_frequent")),
25. #     ("ordinal", OrdinalEncoder(categories=[['Small', 'Medium', 'High']]))
26. # ])
27.
28. type_mean_price_pipeline = Pipeline([
29.     ("type_mean_price_transformer", TypeMeanPriceTransformer(smoothing=0.5)),
30.     ("imputer", SimpleImputer(strategy="mean"))
31. ])
32.
33. # Define one-hot encoding pipeline
34. ordinal_pipeline = Pipeline([
35.     ("imputer", SimpleImputer(strategy="most_frequent")),
36.     ("ordinal", OrdinalEncoder(categories=ordinal_categories))
37. ])
38.
39. # Define one-hot encoding pipeline for onehot-transform columns
40. onehot_transform_pipeline = Pipeline([
41.     ("imputer", SimpleImputer(strategy="most_frequent")),
42.     ("onehot", OneHotEncoder(categories=[['Snack Foods', 'Frozen Foods', 'Fruits and Vegetables', 'Canned']],
handle_unknown="ignore"))
43. ])
44. # Define target encoding pipeline
45. target_encoding_pipeline = Pipeline([
46.     ("target_encoder", TargetEncoder(cols=target_encoder_cols , smoothing=0.5) )
47. ])
48.
49. # Combine all pipelines into a ColumnTransformer
50. preprocessor = ColumnTransformer([
51.     ("num", numerical_pipeline, numerical_cols),
52.     # ("type_mean_price", type_mean_price_pipeline, TypeMeanPriceTransformer_col),
53.     # ("target_encoder", target_encoding_pipeline, target_encoder_cols),
54.     ("ordinal", ordinal_pipeline, ordinal_cols),
55.     ("onehot_", onehot_transform_pipeline, one_hot_columns),

```

```

56.     ("X2_imputer", X2NaNsImputer(), X2_imputer_cols),
57.     ("VisibilityZeros", VisibilityZerosImputerZerosImputer(), VisibilityZeros_cols)
58. ])
59.
60. # Define the final pipeline
61. # Create the extended pipeline
62. final_pipeline = Pipeline([
63.     ("preprocessor", preprocessor)
64.     #("custom_features", CustomFeatureTransformer(column_names=all_features )),
65.     #("kmeans", KMeansTransformer(num_clusters=5, random_state=42))
66. ])
67.
68.
69.
70. # Fit and transform the data
71.
72. # Get the feature names after transformation
73.

```

```

1. import optuna
2. from xgboost import XGBRegressor
3. from lightgbm import LGBMRegressor
4. from sklearn.pipeline import Pipeline
5. from sklearn.model_selection import cross_val_score
6.
7. def objective_xgb(trial):
8.     param = {
9.         'n_estimators': trial.suggest_int('n_estimators', 50, 1000),
10.        'learning_rate': trial.suggest_float('learning_rate', 0.001, 0.3),
11.        'max_depth': trial.suggest_int('max_depth', 3, 15),
12.        'subsample': trial.suggest_float('subsample', 0.4, 1.0),
13.        'colsample_bytree': trial.suggest_float('colsample_bytree', 0.4, 1.0),
14.        'reg_alpha': trial.suggest_float('reg_alpha', 0, 10),
15.        'reg_lambda': trial.suggest_float('reg_lambda', 0, 10),
16.        'min_child_weight': trial.suggest_int('min_child_weight', 1, 10),
17.        'gamma': trial.suggest_float('gamma', 0, 5),
18.        'max_bin': trial.suggest_int('max_bin', 128, 512),
19.        'objective': 'reg:absoluteerror',
20.        'eval_metric': 'mae',
21.        'random_state': 42,
22.        'enable_categorical': True
23.    }
24.
25.    model = XGBRegressor(**param)
26.    scores = cross_val_score(model, X_trans, Y, cv=5, scoring='neg_mean_absolute_error')
27.    mae = -scores.mean()
28.    return mae
29.
30. study_xgb = optuna.create_study(direction='minimize')
31. study_xgb.optimize(objective_xgb, n_trials=70)
32.
33. print("XGBoost Best Parameters:", study_xgb.best_params)
34. print("XGBoost Best MAE:", study_xgb.best_value)
35.

```

```

1. best_xgb_model.fit(X_train, Y_train)
2.
3. # Calculate feature importances

```

```

4. mdi_importances = pd.Series(best_xgb_model.feature_importances_, index=X_train.columns)
5. tree_importance_sorted_idx = np.argsort(best_xgb_model.feature_importances_)
6.
7. fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 8))
8.
9. # Plot Gini importance
10. mdi_importances.sort_values().plot.barh(ax=ax1)
11. ax1.set_xlabel("Gini importance")
12.
13. # Plot permutation importance
14. plot_permutation_importance(best_xgb_model, X_train, Y_train, ax2)
15. ax2.set_xlabel("Decrease in accuracy score")
16.
17. fig.suptitle("Impurity-based vs. permutation importances on multicollinear features (train set)")
18. fig.tight_layout()
19. plt.show()

```

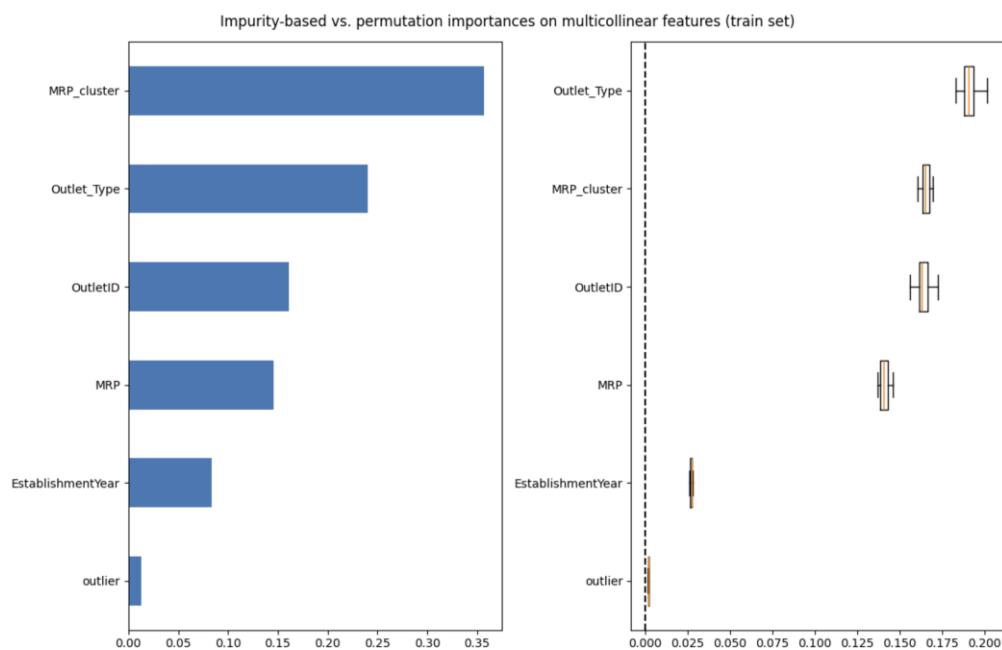
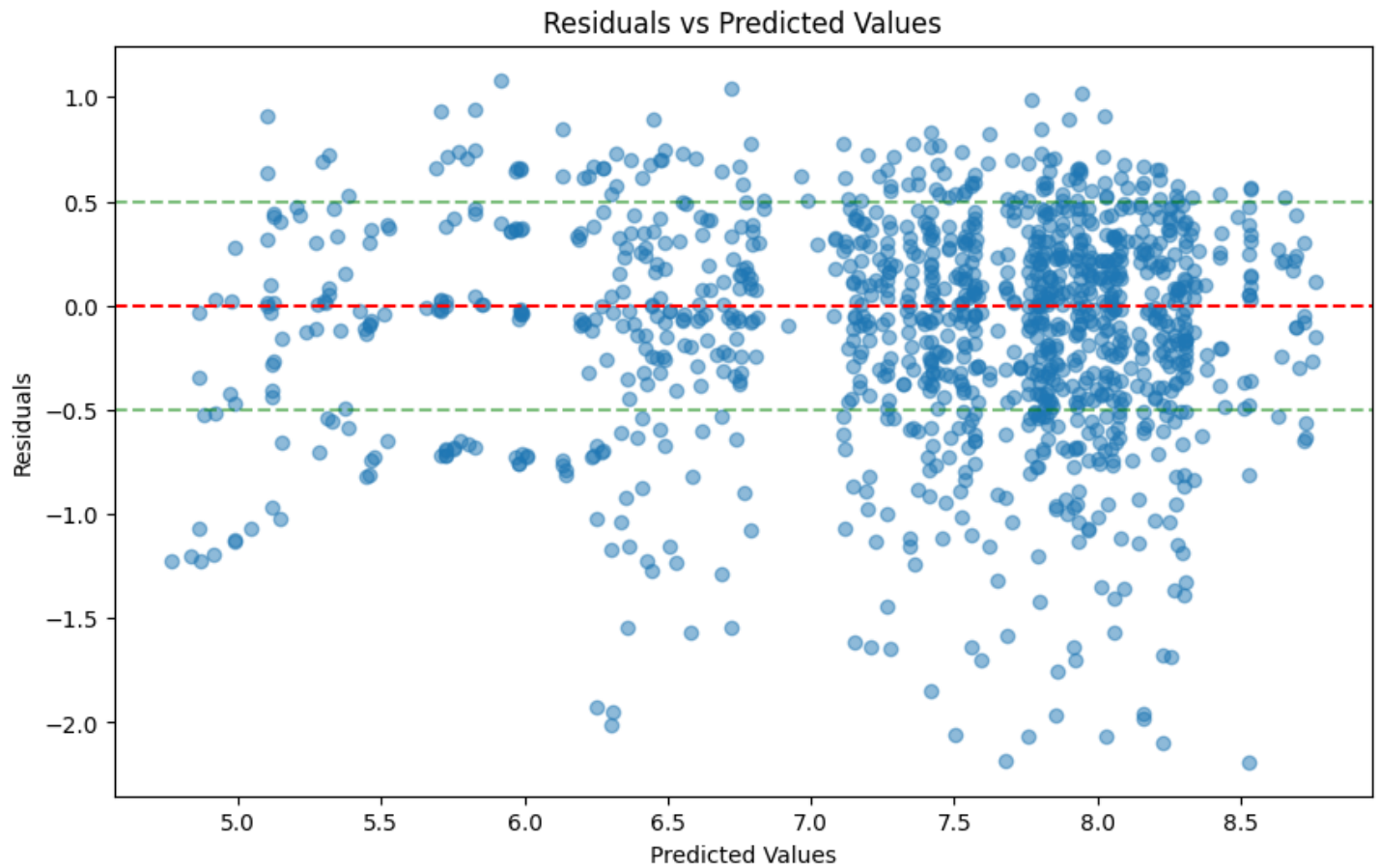


Figure 11 Impurity-based vs. permutation importances on multicollinear features

## Residual Model:



### Insights About Model Performance

#### Overall Bias:

- The residuals are centered around 0, meaning the model generally produces unbiased predictions.
- There is no clear indication of systematic over- or under-prediction across the full range of predicted sales prices.
- This is a positive indicator of the model's calibration.

#### Heteroscedasticity:

- The spread of residuals appears to increase for lower and higher predicted sales prices (e.g., below 6 and above 8). This suggests heteroscedasticity, where the model's prediction error varies with the magnitude of the predicted sales price.



- For low-price products, the model shows higher residual variability, indicating the predictions for these products are less accurate.
- Similarly, high-price products (above ~8) also exhibit larger prediction errors.

**Outliers:**

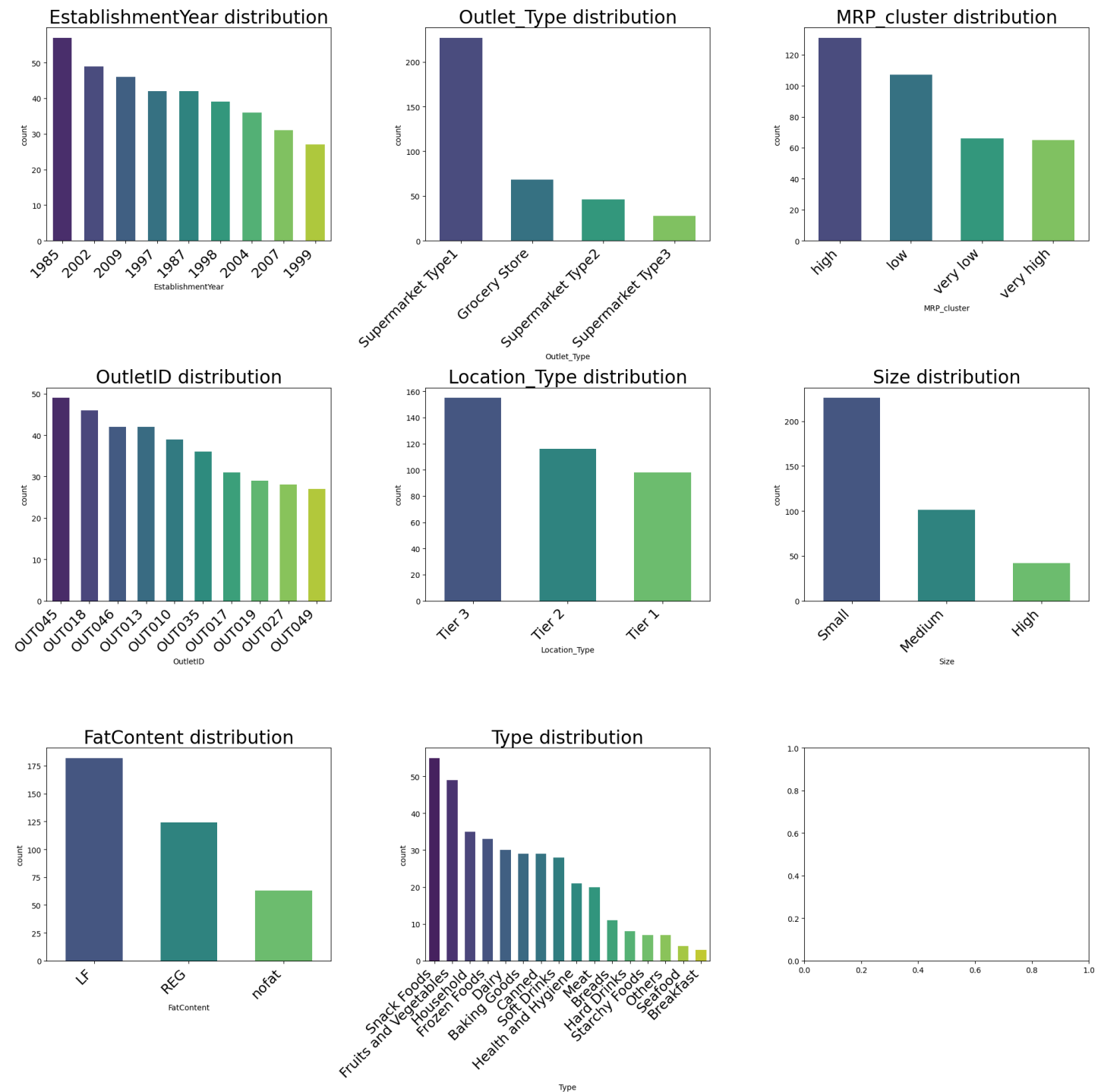
- There are a few outliers with large residuals (e.g.,  $<-1.5$  or  $>1.0$ ). These points correspond to extreme cases where the model predictions deviate significantly from the actual sales price. These outliers could result from:
  - Noise in the data (e.g., data entry errors or rare products).
  - Insufficient representation of these cases in the training data.

**Prediction Accuracy Across Ranges:**

- The residuals for mid-range predicted prices (6.5–7.5) are relatively well-behaved, with smaller spreads.
- This suggests that the model performs better for mid-priced products compared to low- or high-priced ones.

# Residuals > abs(0.5):

## High residual Dataframe:



```
X_trans['outlier'] = (X_trans['Snack Foods'] + X_trans['Frozen Foods'] + X_trans['Fruits and Vegetables'] + X_trans['Canned']) * (X_trans['X2'] + X_trans['Visibility'])
```

