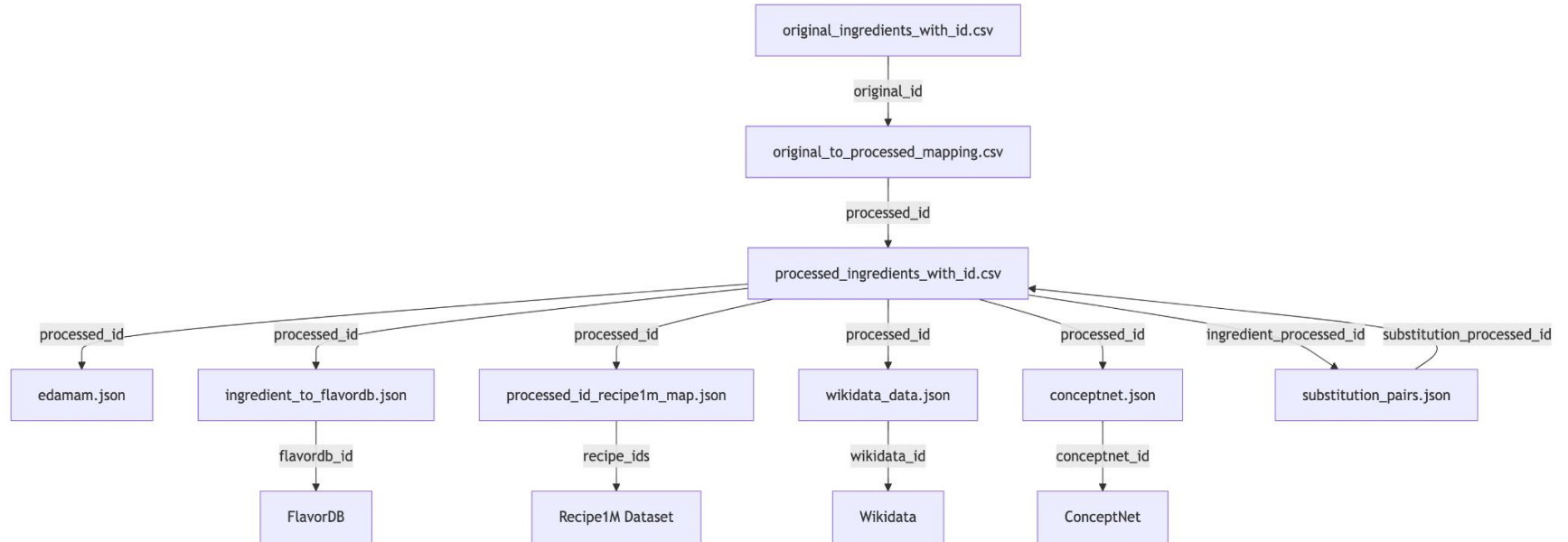


# AI-Powered Ingredient Substitute Recommender

Sami & Tim - group 174

# Complete Data Relationship Diagram



# Dataset 1

Used to find the substitution after model prediction

```
archive > ds_1.csv > data
1 processed_id,processed,processed_ENERC_KCAL,processed_PROCONT,processed_FAT,processed_CHOCDF,processed_FIBTG,processed_category,processed_label,processed_weight,
2 d5a8268e,ababai,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
3 a5bd8077,abalone,105.0,17.1,0.76,6.01,0.0,Generic foods,Serving,85.0,38,1.0
4 a8310f99,abalone sauce,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,38,0.7715134024620056
5 5d97368d,abalone steak,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,38,0.7832129001617432
6 7353da1f,abiu,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
7 9fb0f62f,abondance,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
8 8940b6ab,absinthe,231.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,50.0,39,0.9999999463558197
9 e5c3d2d4,absolut citron vodka,231.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,50.0,40,0.9291608929634094
10 8f379912,absolut kurant vodka,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,40,0.7470061182975769
11 a66607d3,absolut mandarin vodka,231.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,50.0,40,0.7385093569755554
12 318ff5e9,absolut pear flavored vodka,231.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,50.0,4708,0.7543702125549316
13 2d939705,absolut vodka,231.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,50.0,40,0.8422659635543823
14 b8a56b06,abura age,378.0,8.890000343327254,0.0,84.44000244140625,1.100000023841858,Packaged foods,Serving,45.0,-1,-1.0
15 ecf8fc66,aburage,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
16 31c3105e,acai pulp,80.0,2.0,6.0,4.0,3.0,Generic foods,Whole,100.0,-1,-1.0
17 351ec84f,accent seasoning,307.0,9.59,7.53,65.6,11.3,Generic foods,Serving,4.0,-1,-1.0
18 ea4dacd0,aceite,321.1594863380178,4.870625801645378,11.125212824701885,48.82743612905897,1.6259301938672803,Generic meals,Whole,71.59605640300624,7115,0.71968597
19 dacee146,aceite de aguacate,83.0,0.0,5.0,8.3299999923706055,3.299999952316284,Packaged foods,Serving,30.0,-1,-1.0
20 7c9a1e59,aceite de canola,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
21 7d17db74,aceite de oliva,928.3375247959824,0.0,99.97481036264428,0.0,0.0,Packaged foods,Serving,14.003527437778589,-1,-1.0
22 ab98643b,aceite de oliva antiadherente en aerosol,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
23 33bba7a7,aceite en aerosol,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
24 0d22a232,acelga,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
25 42d3775c,aceroia,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,7114,0.7332984209060669
26 7f4a5e7a,acessulfame k,378.0,18.7,1.7,73.0,1.1,Generic foods,Serving,31.0,-1,-1.0
27 d07f0be4,aceto balsamico,107.0,0.0,0.0,20.0,-1.0,Packaged foods,Serving,14.786674781,-1,-1.0
28 77374c91,achar,11.0,0.33,0.2,2.26,1.2,Generic foods,Serving,65.0,44,0.7015097141265869
29 bf4f4fc0,achiote,0.0,0.0,0.0,0.0,-1.0,Packaged foods,Serving,1.7999999523162842,44,0.9999997019767761
30 8e835d75,achiote oil,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,45,0.8849557638168335
31 263503c7,achiote paste,185.0,5.0,2.5,35.0,-1.0,Packaged foods,Serving,20.0,46,0.8851694464683533
32 84893660,achiote powder,0.0,0.0,0.0,0.0,-1.0,Packaged foods,Serving,1.7999999523162842,47,0.8880389332771301
33 3aeaa886,achiote seed,0.0,0.0,0.0,0.0,-1.0,Packaged foods,Serving,1.7999999523162842,44,0.7537855505943298
34 b85eba2f,achute water,0.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,237.0,-1,-1.0
35 359a1fee,acid blend,800.0,0.0,90.0,0.0,0.0,Packaged foods,Serving,9.857843188,-1,-1.0
36 d17fc028,acidophilis milk,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
37 be32739c,acidulated water,0.0,0.0,0.0,0.0,0.0,0.0,Generic foods,Serving,237.0,-1,-1.0
38 c7ada900,acini di pepe,371.0,13.0,1.51,74.7,3.2,Generic foods,Serving,300.0,49,0.7673333287239075
39 769481fe,acini di pepe pasta,371.0,13.0,1.51,74.7,3.2,Generic foods,Serving,300.0,49,0.8909724950790405
40 b3d800c7,ackee,66.0,0.83,0.44,16.5,1.3,Generic foods,Serving,190.0,-1,-1.0
41 86d71b2f,acorn,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,50,0.9999999403953552
42 38b1ae53,acorn flour,387.0,6.15,23.9,40.8,-1.0,Generic foods,Serving,28.35,50,0.7003854513168335
43 cc86eb7f,acorn meal,387.0,6.15,23.9,40.8,-1.0,Generic foods,Serving,28.35,50,0.816230833530426
44 019064ef,acorn squash,40.0,0.8,0.1,10.4,1.5,Generic foods,Serving,140.0,51,0.8988012671470642
45 dccc02fd,active compressed yeast,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,1487,0.7048083543777466
46 47684931,active dry,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,-1,-1.0
47 0ac2f42b,active dry yeast,325.0,40.4,7.61,41.2,26.9,Generic foods,Serving,1.0,2095,0.9207765460014343
48 7c31f4d3,active starter,259.0,9.49,2.73,48.1,2.1,Generic foods,Serving,57.0,53,0.8965895175933838
49 c4d9a80e,adobo de los chipotles,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,1256,0.7516887187957764
50 d29fab29,adobo sauce,125.88723216107005,4.133187961268788,4.669603004551377,19.933852190125688,6.7328380122382026,Generic meals,Whole,295.9735271865784,54,0.9228
51 48a733e58,adobo seasoning,-1.0,-1.0,-1.0,-1.0,-1.0,-1,-1,-1.0,55,0.9249363541603088
```

# Dataset 2

Used for model  
training

```
archive > ds_2.csv > data
1 processed_ENERC_KCAL,processed_PROCNT,processed_FAT,processed_CHOCHF,processed_FIBTG,processed_flavordb_id,processed_cosine_similarity,substitution_ENERC_KCAL,su
2 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,43.0,0.47,0.26,10.8,1.7,4600,1.0000001192092896
3 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,60.0,0.82,0.38,15.0,1.6,4052,1.0000003576278689
4 105.0,17.1,0.76,6.01,0.0,38,1.0,-1.0,-1.0,-1.0,-1.0,-1.0
5 105.0,17.1,0.76,6.01,0.0,38,1.0,86.0,14.7,0.96,3.57,0.0,-1.0
6 105.0,17.1,0.76,6.01,0.0,38,1.0,-1.0,-1.0,-1.0,-1.0,1489,0.7724593877792358
7 105.0,17.1,0.76,6.01,0.0,38,1.0,157.0,21.9,7.02,0.0,0.0,6511,0.8826214671134949
8 105.0,17.1,0.76,6.01,0.0,38,1.0,27.0,6.7,0.0,0.0,0.0,641,0.8387778997421265
9 -1.0,-1.0,-1.0,-1.0,-1.0,38,0.7715134024620056,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
10 -1.0,-1.0,-1.0,-1.0,-1.0,38,0.7832129001617432,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
11 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
12 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,3127,1.0000001192092896
13 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,2404,1.0000003576278689
14 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
15 231.0,0.0,0.0,0.0,0.0,39,0.999999463558197,-1.0,-1.0,-1.0,-1.0,-1.0,3251,0.9999998211860656
16 231.0,0.0,0.0,0.0,0.0,39,0.999999463558197,231.0,0.0,0.0,0.0,0.0,4761,1.0
17 231.0,0.0,0.0,0.0,0.0,39,0.999999463558197,-1.0,-1.0,-1.0,-1.0,-1.0,154,0.999999701976776
18 231.0,0.0,0.0,0.0,0.0,40,0.9291608929634094,14.0,0.0,0.04,3.59,0.0,4873,0.8864786028862
19 -1.0,-1.0,-1.0,-1.0,-1.0,40,0.7470061182975769,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
20 231.0,0.0,0.0,0.0,0.0,40,0.7385093569755554,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
21 231.0,0.0,0.0,0.0,0.0,4708,0.7543702125549316,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
22 231.0,0.0,0.0,0.0,0.0,40,0.8422659635543823,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
23 378.0,8.890000343322754,0.0,84.44000244140625,1.00000023841858,-1,-1.0,78.0,9.04,4.17,2.85,0.9,2612,0.8290473222732544
24 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
25 80.0,2.0,6.0,4.0,3.0,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
26 307.0,9.59,7.53,65.6,11.3,-1,-1.0,307.0,9.59,7.53,65.6,11.3,5700,0.8735607266426086
27 307.0,9.59,7.53,65.6,11.3,-1,-1.0,282.0,14.1,12.9,54.0,34.9,3099,0.8069978952407837
28 307.0,9.59,7.53,65.6,11.3,-1,-1.0,0.0,0.0,0.0,0.0,-1.0,5687,0.9040606021881104
29 321.1594863380178,4.870625801645378,11.125212824701885,48.82743612905897,1.6259301938672803,7115,0.719685971736908,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
30 83.0,0.0,5.0,8.329999923706055,3.299999952316284,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
31 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
32 928.3375247959824,0.0,99.97481036264428,0.0,0.0,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
33 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
34 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
35 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
36 -1.0,-1.0,-1.0,-1.0,-1.0,7114,0.7332984209060669,63.0,1.06,0.2,16.0,2.1,1111,0.9999998211860656
37 378.0,18.7,1.7,73.0,1.1,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
38 378.0,18.7,1.7,73.0,1.1,-1,-1.0,365.0,2.17,0.0,89.1,0.0,-1,-1.0
39 378.0,18.7,1.7,73.0,1.1,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
40 378.0,18.7,1.7,73.0,1.1,-1,-1.0,336.0,0.0,0.0,91.2,0.0,-1,-1.0
41 378.0,18.7,1.7,73.0,1.1,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
42 107.0,0.0,0.0,20.0,-1.0,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
43 11.0,0.33,0.2,2.26,1.2,44,0.7015097141265869,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
44 0.0,0.0,0.0,0.0,-1.0,44,0.999999701976776,130.0,21.6,4.81,0.12,0.0,4123,1.0
45 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,45,0.8849557638168335,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
46 185.0,5.0,2.5,35.0,-1.0,46,0.8851694464683533,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
47 0.0,0.0,0.0,0.0,-1.0,47,0.8880389332771301,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
48 0.0,0.0,0.0,0.0,-1.0,44,0.7537855505943298,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
49 0.0,0.0,0.0,0.0,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
50 800.0,0.0,90.0,0.0,0.0,-1,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0,-1.0
51 -1.0,-1.0,-1.0,-1.0,-1.0,-1.0,15.0,0.59,1.1,1.15,0.58,103,0.8999112844467163
```

# Features and targets

Features (7):

processed\_ENERC\_KCAL,processed\_PROCNT,processed\_FAT,processed\_CHOCDF,processed\_FIBTG,processed\_flavordb\_id,processed\_cosine\_similarity

Targets (7):

substitution\_ENERC\_KCAL,substitution\_PROCNT,substitution\_FAT,substitution\_CHOCDF,substitution\_FIBTG,substitution\_flavordb\_id,substitution\_cosine\_similarity

# Chosen models

- Linear Regression (baseline)
- Support Vector Machine (SVM)
- Random Forest Regression (RFR)
- Neural Network (NN)

# LR

## Training LR

```
train_LR.py > ...
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5 from sklearn.preprocessing import StandardScaler
6
7 # Load dataset
8 df = pd.read_csv("archive/ds_2.csv")
9
10 # Scale the features
11 scaler = StandardScaler()
12 scaled_arr = scaler.fit_transform(df)
13 df_scaled = pd.DataFrame(scaled_arr, columns=df.columns)
14
15 # Select the feature and target columns
16 X = df_scaled[['processed_ENERC_KCAL', 'processed_PROCNT', 'processed_FAT', 'processed_CHOCDF',
17 | | | 'processed_FIBTG', 'processed_flavordb_id', 'processed_cosine_similarity']]
18 y = df_scaled[['substitution_ENERC_KCAL', 'substitution_PROCNT', 'substitution_FAT', 'substitution_CHOCDF',
19 | | | 'substitution_FIBTG', 'substitution_flavordb_id', 'substitution_cosine_similarity']]
20
21
22 # Split the data into training and testing sets (80% train, 20% test)
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
24
25 # Initialize the linear regression model
26 model = LinearRegression()
27
28 # Train the model
29 model.fit(X_train, y_train)
30
31 # Make predictions on the test set
32 y_pred = model.predict(X_test)
33
34 # Evaluate the model using mean squared error
35 mse = mean_squared_error(y_test, y_pred)
36 print(f'Mean Squared Error: {mse}')
37
38
39 print("Coefficients:")
40 print(model.coef_)
41 print("Intercept:")
42 print(model.intercept_)
43
44
45
46 '''
47 Mean Squared Error: 0.8607558734932768
48 Coefficients:
49 [[ 0.37381254 -0.07561138  0.11793094 -0.05367961  0.00399168 -0.04434787
50    0.09655912]
51 [ 0.23662484  0.29234352 -0.25055133 -0.14491602  0.00345179 -0.04788253
```



# LR

## Training LR with automatic hyperparameter search

```
train_LR_auto.py > ...
22 # Split the data
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
24
25 # Define hyperparameter grid
26 param_grid = {'alpha': np.logspace(-4, 2, 10)} # Alpha values from 0.0001 to 100
27
28 # Grid search for ridge regression
29 ridge = Ridge()
30 ridge_search = GridSearchCV(ridge, param_grid, cv=3, scoring='neg_mean_squared_error', n_jobs=-1, verbose=1)
31 ridge_search.fit(X_train, y_train)
32 best_ridge = ridge_search.best_estimator_
33 print(f"Best ridge alpha: {ridge_search.best_params_}")
34
35 # Grid search for lasso regression
36 lasso = Lasso(max_iter=5000) # Ensure Lasso converges
37 lasso_search = GridSearchCV(lasso, param_grid, cv=3, scoring='neg_mean_squared_error', n_jobs=-1, verbose=1)
38 lasso_search.fit(X_train, y_train)
39 best_lasso = lasso_search.best_estimator_
40 print(f"Best lasso alpha: {lasso_search.best_params_}")
41
42 # Train final ridge and lasso models
43 best_ridge.fit(X_train, y_train)
44 best_lasso.fit(X_train, y_train)
45
46 # Make predictions
47 y_pred_ridge = best_ridge.predict(X_test)
48 y_pred_lasso = best_lasso.predict(X_test)
49
50 # Evaluate the models
51 mse_ridge = mean_squared_error(y_test, y_pred_ridge)
52 mse_lasso = mean_squared_error(y_test, y_pred_lasso)
53 print(f"Ridge Mean Squared Error: {mse_ridge}")
54 print(f"Lasso Mean Squared Error: {mse_lasso}")
55
56
57 '''
58 Fitting 3 folds for each of 10 candidates, totalling 30 fits
59 Best ridge alpha: {'alpha': 21.54434690031882}
60 Fitting 3 folds for each of 10 candidates, totalling 30 fits
61 Best lasso alpha: {'alpha': 0.0001}
62 Ridge Mean Squared Error: 0.8607540434816144
63 Lasso Mean Squared Error: 0.8607551729860069
64 '''
```



# SVM

## Training SVM

```
train_SVM.py > ...
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.svm import SVR
4 from sklearn.multioutput import MultiOutputRegressor
5 from sklearn.metrics import mean_squared_error
6 from sklearn.preprocessing import StandardScaler
7
8 # Load dataset
9 df = pd.read_csv("archive/ds_2.csv")
10
11 # Scale the features
12 scaler = StandardScaler()
13 scaled_arr = scaler.fit_transform(df)
14 df_scaled = pd.DataFrame(scaled_arr, columns=df.columns)
15
16 # Select the feature and target columns
17 X = df_scaled[['processed_ENERC_KCAL', 'processed_PROCNT', 'processed_FAT', 'processed_CHOCDF',
18 | 'processed_FIBTG', 'processed_flavordb_id', 'processed_cosine_similarity']]
19 y = df_scaled[['substitution_ENERC_KCAL', 'substitution_PROCNT', 'substitution_FAT', 'substitution_CHOCDF',
20 | 'substitution_FIBTG', 'substitution_flavordb_id', 'substitution_cosine_similarity']]
21
22
23 # Split the data into training and testing sets (80% train, 20% test)
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
25
26 # Initialize the model
27 svr = SVR(kernel='rbf', C=1.0, epsilon=0.1)
28
29 # Make SVR handle multiple outputs
30 multi_output_svr = MultiOutputRegressor(svr)
31
32 # Train the model
33 multi_output_svr.fit(X_train, y_train)
34
35 # Make predictions
36 y_pred = multi_output_svr.predict(X_test)
37
38 # Evaluate the model using mean squared error
39 mse = mean_squared_error(y_test, y_pred)
40 print(f"Mean Squared Error: {mse}")
41
42
43 '''
44 Mean Squared Error: 0.8866638501044166
45 '''
```

# SVM

## Training SVM with automatic hyperparameter search

```
train_SVM_auto.py > ...
23 # Split the data
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
25
26 # Define the hyperparameter search space
27 param_dist = {
28     'estimator__C': np.logspace(-3, 3, 10), # Regularization parameter
29     'estimator__epsilon': np.logspace(-3, 1, 10), # Epsilon in the loss function
30     'estimator__kernel': ['linear', 'poly', 'rbf', 'sigmoid'], # Kernel types
31     'estimator__gamma': ['scale', 'auto'] # Kernel coefficient
32 }
33
34 # Initialize the SVR model
35 base_svr = SVR()
36
37 # Make SVR handle multiple outputs
38 multi_output_svr = MultiOutputRegressor(base_svr)
39
40 # Search for best hyperparameters
41 random_search = RandomizedSearchCV(
42     multi_output_svr,
43     param_distributions=param_dist,
44     n_iter=20,
45     cv=3,
46     verbose=1,
47     n_jobs=-1,
48     scoring='neg_mean_squared_error',
49     random_state=42
50 )
51
52 # Fit the model
53 random_search.fit(X_train, y_train)
54
55 # Best parameters
56 print(f"Best hyperparameters: {random_search.best_params_}")
57
58 # Train the best model
59 best_svr = random_search.best_estimator_
60 best_svr.fit(X_train, y_train)
61
62 # Make predictions
63 y_pred = best_svr.predict(X_test)
64
65 # Evaluate the model
66 mse = mean_squared_error(y_test, y_pred)
67 print(f"Mean Squared Error: {mse}")
68
```

# RFR

## Training RFR

```
train_RFR.py > ...
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.metrics import mean_squared_error
5 from sklearn.preprocessing import StandardScaler
6
7 # Load dataset
8 df = pd.read_csv("archive/ds_2.csv")
9
10 # Scale the features
11 scaler = StandardScaler()
12 scaled_arr = scaler.fit_transform(df)
13 df_scaled = pd.DataFrame(scaled_arr, columns=df.columns)
14
15 # Select the feature and target columns
16 X = df_scaled[['processed_ENERC_KCAL', 'processed_PROCNT', 'processed_FAT', 'processed_CHOCDF',
17               'processed_FIBTG', 'processed_flavordb_id', 'processed_cosine_similarity']]
18 y = df_scaled[['substitution_ENERC_KCAL', 'substitution_PROCNT', 'substitution_FAT', 'substitution_CHOCDF',
19               'substitution_FIBTG', 'substitution_flavordb_id', 'substitution_cosine_similarity']]
20
21
22 # Split the data into training and testing sets (80% train, 20% test)
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
24
25 model = RandomForestRegressor(n_estimators=100, random_state=42)
26 model.fit(X_train, y_train)
27
28 # Make predictions
29 y_pred = model.predict(X_test)
30
31 # Evaluate the model using mean squared error
32 mse = mean_squared_error(y_test, y_pred)
33 print(f"Mean Squared Error: {mse}")
34
35
36 '''
37 Mean Squared Error: 0.673619717571143
38 '''
```

# RFR

## Training RFR with automatic hyperparameter search

```
train_RFR_auto.py > ...
22 # Split the data into training and testing sets (80% train, 20% test)
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
24
25 # Define hyperparameter grid
26 param_dist = {
27     'n_estimators': [50, 100, 200, 300], # Number of trees
28     'max_depth': [None, 10, 20, 30], # Depth of trees
29     'min_samples_split': [2, 5, 10], # Minimum samples to split a node
30     'min_samples_leaf': [1, 2, 4], # Minimum samples in leaf node
31     'max_features': ['auto', 'sqrt'], # Number of features per split
32     'bootstrap': [True, False] # Bootstrapping
33 }
34
35 # Search for the best parameters
36 random_search = RandomizedSearchCV(
37     RandomForestRegressor(random_state=42),
38     param_distributions=param_dist,
39     n_iter=20,
40     cv=3,
41     verbose=1,
42     n_jobs=-1,
43     scoring='neg_mean_squared_error'
44 )
45
46 # Fit the model
47 random_search.fit(X_train, y_train)
48
49 # Best hyperparameters
50 print(f"Best hyperparameters: {random_search.best_params_}")
51
52 # Train final model
53 best_rfr = RandomForestRegressor(**random_search.best_params_, random_state=42)
54 best_rfr.fit(X_train, y_train)
55
56 # Make predictions
57 y_pred = best_rfr.predict(X_test)
58
59 # Evaluate the model
60 mse = mean_squared_error(y_test, y_pred)
61 mae = mean_absolute_error(y_test, y_pred)
62 print(f"Mean Squared Error: {mse}")
63 print(f"Mean Absolute Error: {mae}")
64
65 '''
66 Best hyperparameters: {'n_estimators': 200, 'min_samples_split': 2, 'min_samples_leaf': 4, 'max_features': 'sqrt', 'max_depth': None, 'bootstrap': False}
67 Mean Squared Error: 0.6697726151041151
68 Mean Absolute Error: 0.5246457865163581
69 '''
70
```

# NN

## Training NN

```
train_NN.py > ...
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 import tensorflow as tf
4 from tensorflow import keras
5 from sklearn.metrics import mean_squared_error
6 from sklearn.preprocessing import StandardScaler
7
8 # Load dataset
9 df = pd.read_csv("archive/ds_2.csv")
10
11 # Scale the features
12 scaler = StandardScaler()
13 scaled_arr = scaler.fit_transform(df)
14 df_scaled = pd.DataFrame(scaled_arr, columns=df.columns)
15
16 # Select the feature and target columns
17 X = df_scaled[['processed_ENERC_KCAL', 'processed_PROCNT', 'processed_FAT', 'processed_CHOCDF',
18               'processed_FIBTG', 'processed_flavordb_id', 'processed_cosine_similarity']]
19 y = df_scaled[['substitution_ENERC_KCAL', 'substitution_PROCNT', 'substitution_FAT', 'substitution_CHOCDF',
20               'substitution_FIBTG', 'substitution_flavordb_id', 'substitution_cosine_similarity']]
21
22
23 # Split the data into training and testing sets (80% train, 20% test)
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
25
26 model = keras.Sequential([
27     keras.layers.Dense(64, activation='relu', input_shape=(X.shape[1],)),
28     keras.layers.Dense(32, activation='relu'),
29     keras.layers.Dense(y.shape[1])
30 ])
31
32 # Compile the model
33 model.compile(optimizer='adam', loss='mse', metrics=['mse'])
34
35 # Train the model
36 model.fit(X_train, y_train, epochs=100, batch_size=8, validation_data=(X_test, y_test), verbose=1)
37
38 # Make predictions on the test set
39 y_pred = model.predict(X_test)
40
41 # Evaluate the model using mean squared error (MSE)
42 mse = mean_squared_error(y_test, y_pred)
43 print(f"Mean Squared Error: {mse}")
44
45
46 '''
47 Mean Squared Error: 0.7408530116081238
48 '''
```

# NN

## Training NN with automatic hyperparameter search

```
train_NN_auto.py > ...
23 # Split data
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
25
26 # Model building function for keras tuner
27 def build_model(hp):
28     model = keras.Sequential()
29     model.add(keras.layers.Input(shape=(X.shape[1],)))
30
31     # Number of hidden layers (1-3) and neurons per layer (16-128)
32     for i in range(hp.Int('num_layers', 1, 3)):
33         model.add(keras.layers.Dense(
34             units=hp.Int(f'units_{i}', min_value=16, max_value=128, step=16),
35             activation=hp.Choice('activation', ['relu', 'tanh', 'selu'])
36         ))
37
38     model.add(keras.layers.Dense(y.shape[1]))
39
40     # Learning rate
41     lr = hp.Choice('learning_rate', [0.001, 0.0005, 0.0001])
42     optimizer = keras.optimizers.Adam(learning_rate=lr)
43
44     model.compile(optimizer=optimizer, loss='mse', metrics=['mse'])
45     return model
46
47 # Initialize keras tuner
48 tuner = kt.Hyperband(
49     build_model,
50     objective='val_mse',
51     max_epochs=50,
52     factor=3,
53     directory='tuner_results',
54     project_name='nn_regression_tuning'
55 )
56
57 # Search for the best hyperparameters
58 tuner.search(X_train, y_train, validation_data=(X_test, y_test), epochs=50, batch_size=8, verbose=1)
59
60 # Get the best model
61 best_hps = tuner.get_best_hyperparameters(num_trials=1)[0]
62 best_model = tuner.hypermodel.build(best_hps)
63
64 # Train the best model
65 best_model.fit(X_train, y_train, epochs=100, batch_size=8, validation_data=(X_test, y_test), verbose=1)
66
67 # Predict and evaluate
68 y_pred = best_model.predict(X_test)
69 mse = mean_squared_error(y_test, y_pred)
70 print(f"Best Model Mean Squared Error: {mse}")
71
72 '''
73 MSE: 0.7183371782302856
```

# Results

Best: RFR

Worst: SVM

Metric: MSE

## Model performance (Descending order)

1. RFR auto: 0.6697726151041151
2. RFR: 0.673619717571143
3. NN auto: 0.7183371782302856
4. NN: 0.7408530116081238
5. LR auto: 0.8607540434816144
6. LR: 0.8607558734932768
7. SVM: 0.8866638501044166
8. SVM auto: N/A

The dataset has 90283 rows, 14 cols.



# Usage

Using the model  
to find a substitution  
for any ingredient

```
predict_substitution.py > get_substitution_prediction
33 def get_substitution_prediction(ingredient_features):
34
35     arr = np.array(ingredient_features)
36     row = np.hstack([arr[0], np.zeros(7)])
37     df1 = pd.DataFrame([row], columns=columns)
38     df2 = scaler.transform(df1)
39     df3 = pd.DataFrame(df2, columns=columns)
40     scaled_features = df3.iloc[:, :7]
41
42     predicted_substitution = model.predict(scaled_features)
43
44     arr = np.array(predicted_substitution)
45     row = np.hstack([np.zeros(7), arr[0]])
46     df1 = pd.DataFrame([row], columns=columns)
47     df2 = scaler.inverse_transform(df1)
48     df3 = pd.DataFrame(df2, columns=columns)
49     inverse_scaled_features = df3.iloc[:, 7:]
50
51     substitution_features = inverse_scaled_features.iloc[0].to_numpy()
52
53     return substitution_features
54
55 def find_closest_ingredient(substitution_features):
56     ingredient_features = df_1[feature_cols].values
57     substitution_features = np.array(substitution_features).reshape(1, -1)
58     distances = cdist(ingredient_features, substitution_features, metric="euclidean")
59     closest_idx = np.argmin(distances)
60
61     closest_id = df_1.iloc[closest_idx]["processed_id"]
62     closest_name = df_1.iloc[closest_idx]["processed"]
63
64     return closest_id, closest_name
65
66 ingredient_features = get_ingredient_params(ingredient_name)
67
68 if ingredient_features is not None:
69     print('Ingredient features:')
70     print(ingredient_features)
71
72     substitution_features = get_substitution_prediction(ingredient_features)
73     print('Substitution features:')
74     print(substitution_features)
75
76     substitution_id, substitution_name = find_closest_ingredient([substitution_features])
77     print('Substitution pair:')
78     print(f'{ingredient_name} - {substitution_name}')
```