

# Model

December 3, 2023

```
[1]: from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.metrics import accuracy_score
import pandas as pd
```

```
[2]: # Read data from Excel file into a Pandas DataFrame
file_path = 'dm_mimic_pathways.csv'
df = pd.read_csv(file_path)
```

```
[3]: column_name_mapping = {'person_id': 'Person',
                             'race_concept_id': 'Race',
                             'gender_concept_id': 'Gender',
                             'age_group': 'Age Group',
                             'pathways': 'Treatment Regimen'}

race_mapping = {8527: 'White/ Hispanic',
                 8516: 'Black',
                 8515: 'Asian',
                 0: 'Unknown',
                 38003592: 'Asian',
                 4077359: 'Other',
                 4218674: 'Unknown',
                 4188159: 'White/ Hispanic',
                 38003599: 'Black',
                 38003574: 'Asian',
                 4212311: 'Asian',
                 38003600: 'Black',
                 8557: 'Other',
                 38003584: 'Asian',
                 38003578: 'Asian',
                 4087921: 'Other',
                 38003615: 'Other',
                 38003581: 'Asian',
                 8657: 'Other',
                 38003579: 'Asian',
```

```

38003605: 'Black',
38003614: 'White',
4213463: 'White'}

gender_mapping = {8507: 'Male',
                  8532: 'Female'}

age_mapping = {'10 - 19': 'Teens',
               '20 - 29': 'Twenties',
               '30 - 39': 'Thirties',
               '40 - 49': 'Forties',
               '50 - 59': 'Fifties',
               '60 - 69': 'Sixties',
               '70 - 79': 'Seventies',
               '80 - 89': 'Eighties',
               '> 90': 'Nineties'}

```

```

[4]: df = df.rename(columns=column_name_mapping)
df['Race'] = df['Race'].replace(race_mapping)
df['Gender'] = df['Gender'].replace(gender_mapping)
df['Age Group'] = df['Age Group'].replace(age_mapping)
df['Age Group'].fillna('Unknown', inplace=True)

```

```

[5]: df = df[(df['Age Group'] != 'Unknown') & (df['Race'] != 'Unknown')]

```

```

[6]: print(len(df))
n = 3
values_to_preserve = df['Treatment Regimen'].value_counts().head(n)
print(values_to_preserve)

```

```

1746
Treatment Regimen
19071700          463
19071700,40166274  197
40164929           73
Name: count, dtype: int64

```

```

[7]: def preserve_or_change(value, value_set, replacement_value):
      return value if value in value_set else replacement_value

```

```

[8]: df['Treatment Regimen'] = df['Treatment Regimen'].apply(lambda x:
    ↪ preserve_or_change(x, values_to_preserve, 'Other'))
df.head(5)
len(df['Treatment Regimen'].unique())

```

```

[8]: 4

```

```
[9]: X = df[['Age Group', 'Race', 'Gender']]
     y = df['Treatment Regimen']
```

```
[10]: preprocessor = ColumnTransformer(
        transformers=[
            ('cat', OneHotEncoder(), ['Age Group', 'Race', 'Gender'])
        ],
        remainder='passthrough'
    )
    pipeline = Pipeline([
        ('preprocessor', preprocessor),
        ('classifier', LogisticRegression(multi_class='multinomial', class_weight = 'balanced'))
    ])
```

```
[11]: # Split the data into training and testing sets
     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
[12]: # Train the model
     pipeline.fit(X_train, y_train)
```

```
[12]: Pipeline(steps=[('preprocessor',
                        ColumnTransformer(remainder='passthrough',
                                           transformers=[('cat', OneHotEncoder(),
                                                         ['Age Group', 'Race',
                                                         'Gender'])])),
                      ('classifier',
                       LogisticRegression(class_weight='balanced',
                                           multi_class='multinomial'))])
```

```
[13]: # Make predictions on the test set
     y_pred = pipeline.predict(X_test)

     # Evaluate the accuracy
     accuracy = accuracy_score(y_test, y_pred)
     print(f'Accuracy: {accuracy:.2f}')

     # Create a DataFrame with actual and predicted values
     df_predictions = pd.DataFrame({
         'Actual': y_test,
         'Predicted': y_pred
     })

     print("Actual vs Predicted:")
     print(df_predictions)
```

Accuracy: 0.25

Actual vs Predicted:

	Actual	Predicted
408	19071700	Other
387	19071700	40164929
803	19071700	19071700
81	Other	40164929
942	19071700	40164929
...	...	...
596	19071700,40166274	19071700
1710	Other	Other
894	Other	40164929
1226	Other	Other
1466	Other	40164929

[350 rows x 2 columns]

```
[14]: # Access the one-hot encoder from the pipeline
encoder = pipeline.named_steps['preprocessor'].named_transformers_['cat']

# Get feature names after one-hot encoding
feature_names_after_encoding = list(encoder.get_feature_names_out(X.
    ↳select_dtypes(include=['object']).columns))

# Concatenate feature names with numeric features
all_feature_names = X.select_dtypes(include=['number']).columns.tolist() +
    ↳feature_names_after_encoding

# Access the model from the pipeline
model = pipeline.named_steps['classifier']

# Get coefficients
coefficients = model.coef_

# Display coefficients in a DataFrame
df_coefficients = pd.DataFrame(coefficients, columns=all_feature_names)
df_coefficients['Intercept'] = model.intercept_
df_coefficients['Class'] = model.classes_
df_coefficients.set_index('Class', inplace=True)

print("Coefficients:")
print(df_coefficients)
```

Coefficients:

	Age Group_< 90	Age Group_Eighties	Age Group_Fifties \
Class			
19071700	-0.818003	-0.081705	0.310013
19071700,40166274	-0.680751	-0.675324	-0.301469

40164929	0.576713	0.503740	-0.071875
Other	0.922041	0.253290	0.063331

	Age_Group_Forties	Age_Group_Seventies	Age_Group_Sixties \
Class			
19071700	0.128468	-0.214983	-0.093177
19071700,40166274	0.208711	-0.461090	-0.390980
40164929	-0.122421	0.249454	0.404228
Other	-0.214758	0.426618	0.079928

	Age_Group_Teens	Age_Group_Thirties	Age_Group_Twenties \
Class			
19071700	-0.310912	0.259881	0.820227
19071700,40166274	0.789999	0.605955	0.904688
40164929	-0.285045	-0.131003	-1.123763
Other	-0.194042	-0.734833	-0.601152

	Race_Asian	Race_Black	Race_Other	Race_White \
Class				
19071700	0.382025	0.004256	-0.276349	-0.067690
19071700,40166274	-0.791962	0.358433	0.358460	-0.134346
40164929	0.203091	-0.164447	0.099285	0.215719
Other	0.206845	-0.198242	-0.181396	-0.013682

	Race_White/ Hispanic	Gender_Female	Gender_Male	Intercept
Class				
19071700	-0.042432	0.010430	-0.010620	0.074985
19071700,40166274	0.209155	-0.017111	0.016850	0.104389
40164929	-0.353620	0.052054	-0.052026	0.088801
Other	0.186897	-0.045373	0.045797	-0.268174