GOOGLE MAPS

Submitted by:                                    Submitted to:
Team Amigos                                        Dr. Shivanjali Khare

Team Members:
Member 1: Jayanth Reddy Manda (00764288) jmand5@unh.newhaven.edu
Member 2: Snigdha Reddy Yeruva (00762316) syeru2@unh.newhaven.edu
Member 3: Jeevan Kumar Konduru (00806735) Jkond6@unh.newhaven.edu

GitHub link: https://github.com/JayanthReddy4/Phase-6-AMIGOS.git

Research question:
Recommendation to the customer in a path from source to destination (Ex: let us Assume google maps knows user is from India, assume user is travelling from "Newhaven railway station " to "Hartford Railway Station" via car. If the user needs coffee, we recommend best coffee available stores near user based on business data available in google maps such as rating, price levels, user rating, etc....)

Dataset:



As shown in the above image we deal with multiple attributes such as:

Latitude & Longitude: Coordinates of a business so that we able to predict if user is available in available range of coordinates.

Name: Name of different businesses at coordinates.

Vicinity: Human readable location of a business.

Type: Indicates different attributes of a store such as availability of food, coffee, restaurant….

Rating: Indicates the rating of a business by the customer, it's one of the key parameter for the business model, it's one of the key parameter whether a new customer is willing to visit or not.

User rating: Number of ratings users made on business. so that if the business has more ratings with God rating it indicates good business.

List of Data Mining Techniques used in optimization Phase:

1) Content Based Recommendation System.
2) Random Forest Classifier

Details of Model parameters and Hyperparameters used in Optimization datamining techniques:

| Data Mining Techniques | Model Parameters | Hyperparameters |
|---|---|---|
| Content Based Recommendation System. | Features, rating | Name, storeId, userId |
| Random Forest Classifier | Rating, user_ratings_total | Price_level |

**Content Recommendation System:**

The model parameters used in previous content recommendation system are Features and technique used is cosine similarity, in optimized technique the model parameters are Features and rating and technique used is cosine similarity with Tf-idf. we were able to recommend better stores by including rating with Features parameter.

In the optimized technique we will be using precision to evaluate the recommendation system with an value of 75% .

```
# 5. Test the recommendation system
movie_title = 'West Haven'  # Replace with the store title you want recommendations for
recommendations = recommend_movies(movie_title)
print("Recommended stores based on", movie_title, "are:")
print(recommendations)
```

```
Recommended stores based on West Haven are:
1                        Hilton Garden Inn Milford
9                                             ALDO
5                Super 8 by Wyndham Milford/New Haven
6                                           Zumiez
7                            Connecticut Post Mall
8                                        AT&T Store
2                    Hyatt Place Milford / New Haven
10                                    Hollister Co.
3        Courtyard by Marriott New Haven Orange/Milford
4                            Hampton Inn Milford
Name: name, dtype: object
```

```
West Haven: Precision = 0.40, Recall = 0.67
Average Precision = 0.40, Average Recall = 0.67
```

From the above precision level west have has got the 40% precision rate and 67% Recall.

## Heat Map for Cosine Similarity Matrix:



## Heat map to visualize ratings data:

**Feature Frequency:**



**Random Forest Classifier:**

```
In [28]:  #Create classifier
          rf_classifier = RandomForestClassifier(n_jobs=-1)

In [29]:  # set different parameter values to tune
          param_grid = {
              "n_estimators": [100, 200, 300, 400],
              "max_depth": [1, 3, 5, 7, 9],
              "criterion": ["gini", "entropy"],
          }

In [31]:  # train the model with gridserchCV
          model.fit(X_scaled,y)

          Fitting 5 folds for each of 40 candidates, totalling 200 fits
          [CV] END ......criterion=gini, max_depth=1, n_estimators=100; total time=
          3.4s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=100; total time=
          0.1s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=100; total time=
          0.2s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=100; total time=
          0.1s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=200; total time=
          0.3s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=200; total time=
          0.4s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=200; total time=
          0.3s
          [CV] END ......criterion=gini, max_depth=1, n_estimators=200; total time=
          0.3s

In [32]:  # print the best score and estimator
          print(model.best_score_)
          print(model.best_estimator_.get_params())

          0.9052631578947368
          {'bootstrap': True, 'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gi
          ni', 'max_depth': 9, 'max_features': 'auto', 'max_leaf_nodes': None, 'max_sa
          mples': None, 'min_impurity_decrease': 0.0, 'min_samples_leaf': 1, 'min_samp
          les_split': 2, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 200, 'n_job
          s': -1, 'oob_score': False, 'random_state': None, 'verbose': 0, 'warm_star
          t': False}
```
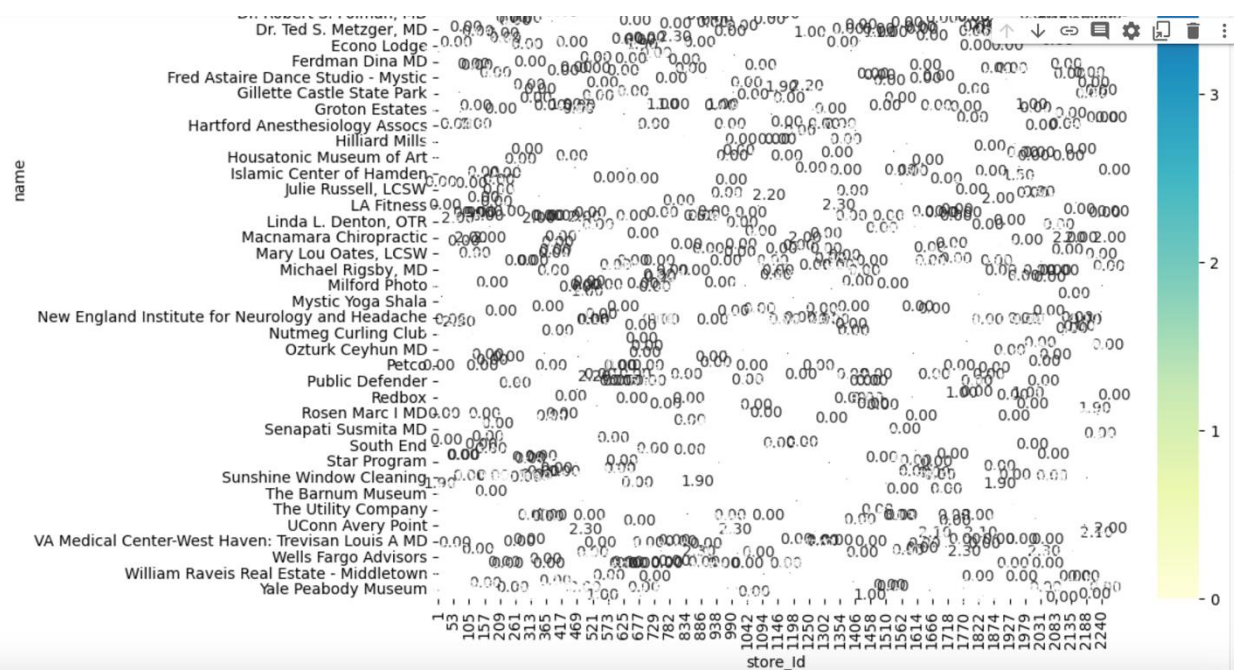
```
In [70]:  # perform optimization
          result = gp_minimize(
              func=evaluate_model,
              dimensions=search_space,
              n_calls=30,
              random_state=42,
              verbose=True,
              n_jobs=1,
          )

          Iteration No: 1 started. Evaluating function at random point.
          Iteration No: 1 ended. Evaluation done at random point.
          Time taken: 3.8747
          Function value obtained: -0.8882
          Current minimum: -0.8882
          Iteration No: 2 started. Evaluating function at random point.
          Iteration No: 2 ended. Evaluation done at random point.
          Time taken: 3.8083
          Function value obtained: -0.8368
          Current minimum: -0.8882
          Iteration No: 3 started. Evaluating function at random point.
          Iteration No: 3 ended. Evaluation done at random point.
          Time taken: 2.3749
          Function value obtained: -0.8368
          Current minimum: -0.8882
          Iteration No: 4 started. Evaluating function at random point.
          Iteration No: 4 ended. Evaluation done at random point.
          Time taken: 3.8165
          Function value obtained: -0.8873

In [71]:  # summarizing finding:
          print('Best Accuracy: %.3f' % (abs(result.fun)))
          print('Best Parameters: %s' % (result.x))

          Best Accuracy: 0.907
          Best Parameters: [400, 'gini', 9]
```

The above image depicts random forest with a best score when the evaluation pattern is "gini index" to build decision tree with 9 levels in building a decision tree.
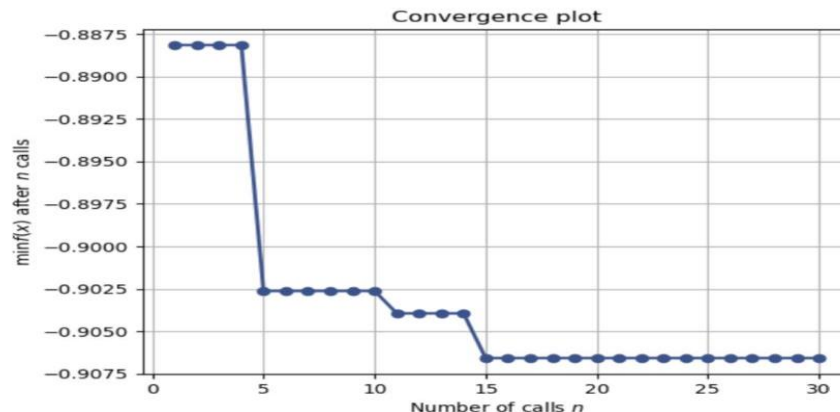
**Convergence Plot for the Random Forest Classifier:**

```
In [73]:  # plot convergence
          from skopt.plots import plot_convergence
          plot_convergence(result)

Out[73]:  <AxesSubplot:title={'center':'Convergence plot'}, xlabel='Number of calls $n
          $', ylabel='$\\min f(x)$ after $n$ calls'>
```



**Conclusion**:

The ability of Recommending a store to the user has improved by using content-based recommendation by tuning its model parameters with additional model parameters which gives a better efficient way compared to previous technique. The optimized content based recommendation system creates cosine similarity for an transformed data by using tf-idf which gives more accuracy, whereas in random forest classifies the usage of "gini" over "entropy " for better results.