# Model Optimization and Tuning Phase Template

| Date | July 2024 |
|---|---|
| Team ID | 740295 |
| Project Title | Ecommerce Shipping Prediction using Machine Learning |
| Maximum Marks | 10 Marks |

## Model Optimization and Tuning Phase

The Model Optimization and Tuning Phase involves refining neural network models for peak performance. It includes optimized model code, fine-tuning hyperparameters, comparing performance metrics, and justifying the final model selection for enhanced predictive accuracy and efficiency.

**Hyperparameter Tuning Documentation (8 Marks):**

| Model | Tuned Hyperparameters |
|---|---|
| Logistic Regression | #importing the library for grid search<br>from sklearn.model_selection import GridSearchCV<br><br>The 'lr_param_grid' specifies different values for regularization strength (C), solvers (solver), and penalty types (penalty). GridSearchCV (lr_cv) is employed with 5-fold cross-validation (cv=5), evaluating model performance based on accuracy (scoring="accuracy"). The process uses all available CPU cores (n_jobs=-1) for parallel processing and provides verbose output (verbose=True) to track progress.<br><br>LOGISTIC REGRESSION HYPER PARAMETER TUNNING<br><br>`[54]` `#finding the grid search cv for logistic regression`<br>`lr=LogisticRegression(n_jobs=-1,random_state=0)`<br>`lr_param_grid={`<br>`    'C':[0.1,0.5,1,5,10],`<br>`    'solver':['liblinear','saga'],`<br>`    'penalty':['l1','l2']`<br><br>`}`<br>`lr_cv=GridSearchCV(lr,lr_param_grid,cv=5,scoring="accuracy",n_jobs=-1,verbose=T`<br>`lr_cv.fit(x_train,y_train)`<br><br>`Fitting 5 folds for each of 20 candidates, totalling 100 fits`<br>`/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:1211:`<br>`  warnings.warn(`<br>`        GridSearchCV`<br>`► estimator: LogisticRegression`<br>`     ► LogisticRegression` |

| | |
|---|---|
| Random Forest | The parameter grid (rfc_param_grid) for hyperparameter tuning. It specifies different values for the number of trees (n_estimators), splitting criterion (criterion), maximum depth of trees (max_depth), and maximum number of features considered for splitting (max_features). GridSearchCV (rfc_cv) is employed with 3-fold cross-validation (cv=3), evaluating model performance based on accuracy (scoring="accuracy").<br><br>RANDOM FOREST HYPER PARAMETER TUNNING<br><br>```python
[55] #finding the grid search cv for random forest classifier
     rfc=RandomForestClassifier()
     rfc_param_grid={
         'n_estimators':[100,200],
         'criterion':['entropy','gini'],
         'max_depth':[5,10],
         'max_features':['auto','sqrt']
     }
     rfc_cv=GridSearchCV(rfc,rfc_param_grid,cv=3,scoring="accuracy",n_jobs=-1,verbose=3)
     rfc_cv.fit(x_train,y_train)
```<br><br>Fitting 3 folds for each of 16 candidates, totalling 48 fits<br>/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_forest.py:424: FutureWarning:<br>  warn(<br><br>GridSearchCV<br>▸ estimator: RandomForestClassifier<br>  ▸ RandomForestClassifier |
| XGBoost | The (params) define a grid for hyperparameter tuning of the XGBoost Classifier (XGBClassifier), including min_child_weight, gamma, colsample_bytree, and max_depth. The XGBClassifier is configured with a learning rate of 0.5, 100 estimators, using a binary logistic regression objective, and utilizing 3 threads for processing. GridSearchCV (xg_cv) is used with 5-fold cross-validation (cv=5), refitting the best model (refit=True), evaluating based on accuracy (scoring="accuracy")<br><br>XGBOOST CLASSIFIER-HYPER PARAMETER TUNNIG<br><br>```python
#finding the grid search cv for xgboost
params={
    'min_child_weight':[10,20],
    'gamma':[1.5,2.0,2.5],
    'colsample_bytree':[0.6,0.8,0.9],
    'max_depth':[4,5,6]
}
xg=XGBClassifier(learning_rate=0.5,n_estimators=100,objective='binary:logistic',nthreads=3)
xg_cv=GridSearchCV(xg,param_grid=params,cv=5,refit=True,scoring="accuracy",n_jobs=-1,verbose=3)
xg_cv.fit(x_train,y_train)
```<br><br>Fitting 5 folds for each of 54 candidates, totalling 270 fits<br>/usr/local/lib/python3.10/dist-packages/xgboost/core.py:160: UserWarning: [14:07:26] WARNING: /works<br>Parameters: { "nthreads" } are not used.<br><br>  warnings.warn(smsg, UserWarning)<br><br>GridSearchCV<br>▸ estimator: XGBClassifier<br>  ▸ XGBClassifier |

| | |
|---|---|
| Decision Tree | The parameters (params) define a grid for hyperparameter tuning of the Decision Tree Classifier (DecisionTreeClassifier), including max_depth, min_samples_leaf, and criterion ('gini' or 'entropy'). GridSearchCV (dec_cv) is used with 5-fold cross-validation (cv=5), evaluating model performance based on accuracy (scoring="accuracy")<br><br>DECISION TREE CLASSIFIER-HYPER PARAMETER TUNNING<br><br>```python<br>[68] #finding grid search cv for decision tree classifier<br>    dec=DecisionTreeClassifier(random_state=42)<br>    params={ 'max_depth': [2, 3, 5, 10, 20],<br>        'min_samples_leaf': [5, 10, 20, 50, 100],<br>        'criterion': ['gini', 'entropy']<br>    }<br>    dec_cv=GridSearchCV(dec,param_grid=params,cv=5,n_jobs=-1,scoring="accuracy",verbose=3)<br>    dec_cv.fit(x_train,y_train)<br>```<br>Fitting 5 folds for each of 50 candidates, totalling 250 fits<br>    GridSearchCV<br>  estimator: DecisionTreeClassifier<br>    ▸ DecisionTreeClassifier |
| Ridge Classifier | The parameters (params) define a grid for hyperparameter tuning of the Decision Tree Classifier (DecisionTreeClassifier), including max_depth, min_samples_leaf, and criterion ('gini' or 'entropy'). GridSearchCV (dec_cv) is used with 5-fold cross-validation (cv=5), evaluating model performance based on accuracy (scoring="accuracy")<br><br>RIDGE-CLASSIFIER-HYPER PARAMETER TUNNING<br><br>```python<br>#finding the grid search cv for ridge classifier<br>rg=RidgeClassifier(random_state=42)<br>params={<br>    'alpha':(np.logspace(-8,8,100))<br>}<br>rg_cv=GridSearchCV(rg,param_grid=params,cv=5)<br>rg_cv.fit(x_train,y_train)<br>```<br>    GridSearchCV<br>  estimator: RidgeClassifier<br>    ▸ RidgeClassifier |

| | |
|---|---|
| K- Nearest Neighbors | The parameters (params) define a grid for hyperparameter tuning of the K-Nearest Neighbors Classifier (KNeighborsClassifier), including n_neighbors, weights ('uniform' or 'distance'), and metric ('minkowski', 'euclidean', or 'manhattan'). GridSearchCV (knn_cv) is used with 5-fold cross-validation (cv=5), evaluating model performance based on accuracy (scoring="accuracy")<br><br>K-NEAREST NEIGHBORS-HYPER PARAMETER TUNNING<br><br><pre>[69] #finding the grid search cv for k-nearest neighbors<br>     knn=KNeighborsClassifier()<br>     params={<br>         'n_neighbors':[3,5,7,9,11],<br>         'weights':['uniform','distance'],<br>         'metric':['minkowski','eculidean','manhattan']<br>     }<br>     knn_cv = GridSearchCV(knn, param_grid=params,cv=5, n_jobs=-1, verbose=3)<br>     knn_cv.fit(x_train, y_train)</pre><br>▸      **GridSearchCV**<br>▸ **estimator: KNeighborsClassifier**<br>   ▸ KNeighborsClassifier |

# Final Model Selection Justification (2 Marks):

| Final Model | Reasoning |
|---|---|
| **Random Forest** | Random Forest model is chosen for its robustness in handling complex datasets and its ability to mitigate overfitting while providing high predictive accuracy.<br><br>|  | Name | Accuracy | f1_score | Recall | Precision |<br>|---|---|---|---|---|---|<br>| 0 | Logistic Regression | 67.90 | 64.68 | 59.16 | 71.35 |<br>| 1 | Decision Tree Classifier | 73.88 | 66.60 | 52.41 | 91.32 |<br>| 2 | Random Forest | 74.68 | 66.70 | 51.03 | 96.24 |<br>| 3 | K-Nearest Nieghbors | 74.56 | 71.57 | 64.44 | 80.48 |<br>| 4 | Xgboost | 74.18 | 68.61 | 56.78 | 86.67 |<br>| 5 | Ridge Classifier | 68.39 | 63.91 | 56.32 | 73.87 |<br><br>Above all the models Random Forest model have the highest accuracy among all the models. |