# Model Optimization and Tuning Phase Template

| Date | July 2024 |
|------|-----------|
| Team ID | 739859 |
| Project Title | Auto insurance fraud detection using michine 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<br><br>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>```python
# Logistic Regression
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
model = LogisticRegression()
model = model.fit(X_Train, Y_Train)
pred = model.predict(X_Test)

print('Accuracy:', accuracy_score(Y_Test, pred))
print('\n classification report:\n', classification_report(Y_Test, pred))
print('\n confusion matrix:\n', confusion_matrix(Y_Test, pred))
```<br><br>```
Accuracy: 0.75

classification report:
              precision    recall  f1-score   support

           0       0.00      0.00      0.00        48
           1       0.76      0.99      0.86       152

    accuracy                           0.75       200
   macro avg       0.38      0.49      0.43       200
weighted avg       0.58      0.75      0.65       200


confusion matrix:
[[  0  48]
 [  2 150]]
``` |

| 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> |
|---|---|
| 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> |

| | |
|---|---|
| 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> |
| 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><br>GridSearchCV<br>▸ estimator: RidgeClassifier<br>    ▸ RidgeClassifier |

| | |
|---|---|
| K- Nearest<br><br>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> |

# 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. |