

## Model Optimization and Tuning Phase Template

Date	21 July 2024
Team ID	SWTID1721319573
Project Title	Blueberry Yield Prediction
Maximum Marks	10 Marks

### Model Optimization and Tuning Phase

The Model Optimization and Tuning Phase involves refining machine learning 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 (6 Marks):

Model	Tuned Hyperparameters	Optimal Values
Linear Regression	<pre># Define the model and parameters for tuning lin_reg = LinearRegression() param_grid = {'fit_intercept': [True, False]}  # Perform GridSearchCV grid_search_lr = GridSearchCV(estimator=lin_reg, param_grid=param_grid, cv=5) grid_search_lr.fit(X_train, y_train)  # Get the best model from GridSearchCV best_lin_reg = grid_search_lr.best_estimator_ pred_linear = best_lin_reg.predict(X_test)</pre>	<p>Linear Regression - Best Hyperparameters: Best Hyperparameters: {'fit_intercept': False}</p>
Random Forest Regressor	<pre># Define the model and parameters for tuning rf_reg = RandomForestRegressor() param_grid_rf = {     'n_estimators': [100, 200],     'max_depth': [None, 10, 20],     'min_samples_split': [2, 5],     'min_samples_leaf': [1, 2] }  # Perform GridSearchCV grid_search_rf = GridSearchCV(estimator=rf_reg, param_grid=param_grid_rf, cv=5, n_jobs=-1) grid_search_rf.fit(X_train, y_train)  # Get the best model from GridSearchCV best_rf_reg = grid_search_rf.best_estimator_ pred_rf = best_rf_reg.predict(X_test)</pre>	<p>Random Forest Regressor - Best Hyperparameters: Best Hyperparameters: {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 200}</p>

Decision Tree Regressor	<pre># Define the model and parameters for tuning dt_reg = DecisionTreeRegressor() param_grid_dt = {     'max_depth': [None, 10, 20],     'min_samples_split': [2, 5],     'min_samples_leaf': [1, 2] }  # Perform GridSearchCV grid_search_dt = GridSearchCV(estimator=dt_reg, param_grid=param_grid_dt, cv=5, n_jobs=-1) grid_search_dt.fit(X_train, y_train)  # Get the best model from GridSearchCV best_dt_reg = grid_search_dt.best_estimator_ pred_dt = best_dt_reg.predict(X_test)</pre>	<p>Decision Tree Regressor - Best Hyperparameters: Best Hyperparameters: {'max_depth': None, 'min_samples_leaf': 2, 'min_samples_split': 5}</p>
XGBoost Regressor	<pre># Define the model and parameters for tuning xgb_reg = XGBRegressor(objective='reg:squarederror') param_grid_xgb = {     'n_estimators': [100, 200],     'max_depth': [3, 5, 7],     'learning_rate': [0.01, 0.1, 0.3] }  # Perform GridSearchCV grid_search_xgb = GridSearchCV(estimator=xgb_reg, param_grid=param_grid_xgb, cv=5, n_jobs=-1) grid_search_xgb.fit(X_train, y_train)  # Get the best model from GridSearchCV best_xgb_reg = grid_search_xgb.best_estimator_ pred_xgb = best_xgb_reg.predict(X_test)</pre>	<p>XGBoost Regressor - Best Hyperparameters: Best Hyperparameters: {'learning_rate': 0.1, 'max_depth': 3, 'n_estimators': 200}</p>

## Performance Metrics Comparison Report (2 Marks):

Model	Baseline Metric	Optimized Metric
Linear Regression	<p>Accuracy: 99.18%</p> <hr/> <p><b>Linear Regression:</b>  <b>MAE: 97.318</b>  <b>MSE: 16219.955</b>  <b>RMSE: 127.358</b>  <b>R-Square: 0.992</b>  <b>Accuracy: 99.18%</b></p>	<p>Linear Regression - Best Hyperparameters:  Best Hyperparameters: {'fit_intercept': False}  Performance Metrics:  MAE: 97.318  MSE: 16219.955  RMSE: 127.358  R-Square: 0.992  Accuracy: 99.18%</p>
Random Forest Regressor	<p>Accuracy: 98.84%</p>	<p>Accuracy: 98.84%</p> <hr/> <p>Random Forest Regressor - Best Hyperparameters:  Best Hyperparameters: {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 200}  Performance Metrics:  MAE: 116.756  MSE: 22876.699  RMSE: 151.250  R-Square: 0.988  Accuracy: 98.84%</p>

	<p>Random Forest Regressor:</p> <p>MAE: 117.197</p> <p>MSE: 22845.764</p> <p>RMSE: 151.148</p> <p>R-Square: 0.988</p> <p>Accuracy: 98.84%</p>	
Decision Tree Regressor	<p>Accuracy: 98.05%</p> <hr/> <p>Decision Tree Regressor:</p> <p>MAE: 148.381</p> <p>MSE: 38588.977</p> <p>RMSE: 196.441</p> <p>R-Square: 0.980</p> <p>Accuracy: 98.05%</p>	<p>Accuracy: 98.19%</p> <hr/> <p>Decision Tree Regressor - Best Hyperparameters: Best Hyperparameters: {'max_depth': None, 'min_samples_leaf': 2, 'min_samples_split': 5} Performance Metrics: MAE: 144.743 MSE: 35830.154 RMSE: 189.289 R-Square: 0.982 Accuracy: 98.19%</p>
XGBoost Regressor	<p>Accuracy: 99.09%</p> <hr/> <p>XGBoost Regressor:</p> <p>MAE: 106.537</p> <p>MSE: 17901.843</p> <p>RMSE: 133.798</p> <p>R-Square: 0.991</p> <p>Accuracy: 99.09%</p>	<p>Accuracy: 99.11%</p> <hr/> <p>XGBoost Regressor - Best Hyperparameters: Best Hyperparameters: {'learning_rate': 0.1, 'max_depth': 3, 'n_estimators': 200} Performance Metrics: MAE: 101.627 MSE: 17608.845 RMSE: 132.698 R-Square: 0.991 Accuracy: 99.11%</p>

### Final Model Selection Justification (2 Marks):

Final Model	Reasoning
Linear Regression	The Linear Regression model achieved the highest accuracy of 99.18% compared to other models. It provided a robust performance with the

best R-Square value of 0.992. Despite its simplicity, Linear Regression's high accuracy and efficiency make it the most suitable model for the given task.

1. **Highest R-Square Value:** Linear Regression achieved the highest R-Square value (0.992), indicating that it explains 99.2% of the variance in the target variable. This suggests that the model fits the data better than the other models.
2. **Lowest MAE and RMSE:** The Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) for Linear Regression are lower than those of the other models. This indicates that Linear Regression's predictions are closer to the actual values, making it more accurate and reliable.
3. **Simplicity and Interpretability:** Linear Regression is a simpler and more interpretable model compared to more complex models like XGBoost or Random Forest. Despite its simplicity, it outperformed the other models in terms of accuracy, making it a preferable choice for this particular problem.
4. **Consistency Across Metrics:** Linear Regression consistently showed the best performance across multiple metrics (MAE, MSE, RMSE, R-Square), proving its robustness and reliability as the best model for this task.

**Conclusion:** Linear Regression is chosen as the best model because it provides the highest accuracy (99.18%) and the best performance across various metrics. Its simplicity and interpretability further support

	its selection, ensuring both strong predictive power and ease of understanding.
--	---