# Model Development Phase Template

| Date | 1 July 2025 |
|---|---|
| Team ID | NONE |
| Project Title | Employee Performance Prediction using Machine Learning |

**Initial Model Training Code, Model Validation and Evaluation Report**

Thee initial model training involved implementing and evaluating three regression algorithms — Linear Regression, Random Forest, and XGBoost — to predict employee productivity. The models were trained on historical work data and evaluated using the $R^2$ (coefficient of determination) score as the primary metric. A comparative analysis was performed to identify the best-performing model, with Random Forest achieving the highest $R^2$ score of approximately 0.46. Model evaluation results were visualized using a bar chart for clear comparison, and the trained Random Forest model was saved for integration into the Flask-based web application.

**Initial Model Training Code:**

```python
1   # Import required libraries
2   from sklearn.model_selection import train_test_split
3   from sklearn.linear_model import LinearRegression
4   from sklearn.ensemble import RandomForestRegressor
5   from xgboost import XGBRegressor
6   from sklearn.metrics import r2_score
7   import joblib
8
9   # Split the data
10  X = df.drop('actual_productivity', axis=1)
11  y = df['actual_productivity']
12  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_sta
13
14  # Initialize models
15  models = {
16      "Linear Regression": LinearRegression(),
17      "Random Forest": RandomForestRegressor(n_estimators=100, random_state=42),
18      "XGBoost": XGBRegressor(n_estimators=100, random_state=42)
19  }
20
```

```python
# Train and evaluate
results = {}
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    results[name] = r2_score(y_test, y_pred)
    print(f"{name} R² Score: {results[name]:.3f}")

# Select best model
best_model_name = max(results, key=results.get)
best_model = models[best_model_name]
print(f"\nBest Model: {best_model_name}")

# Save the best model
joblib.dump(best_model, 'model/best_model.pkl')
joblib.dump(X.columns.tolist(), 'model/model_features.pkl')
```

```python
1   # Plot model comparison
2   import matplotlib.pyplot as plt
3
4   plt.figure(figsize=(10, 5))
5   plt.bar(results.keys(), results.values(), color=['skyblue', 'lightgreen', 'orange'
6   plt.title('Model Comparison (R² Score)')
7   plt.ylabel('R² Score')
8   plt.ylim(0, 0.6)
9   for i, (key, value) in enumerate(results.items()):
10      plt.text(i, value + 0.01, f'{value:.3f}', ha='center', va='bottom')
11  plt.savefig('assets/model_accuracy.png', dpi=300, bbox_inches='tight')
12  plt.show()
```

**Validation and Evaluation Report:**

| Model | Classification Report | R² socre | Confusion Matrix |
|---|---|---|---|
| Random Forest | Best performing model with non-linear feature handling and high R² score. | 0.44 | (base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml><br>Linear Regression R² Score: 0.16816825663066545<br>Random Forest R² Score: 0.44671974539154335<br>XGBoost R² Score: 0.3538597397101051<br>Best model saved to model/best_model.pkl<br>(base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml> |
| Linear Regression | Baseline model; assumes linear relationships. Underperformed due to non-linear patterns in data. | 0.16 | (base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml><br>Linear Regression R² Score: 0.16816825663066545<br>Random Forest R² Score: 0.44671974539154335<br>XGBoost R² Score: 0.3538597397101051<br>Best model saved to model/best_model.pkl<br>(base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml> |
| XGBoost | Strong ensemble model, but slightly overfit and lower R² than Random Forest. | 0.35 | (base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml><br>Linear Regression R² Score: 0.16816825663066545<br>Random Forest R² Score: 0.44671974539154335<br>XGBoost R² Score: 0.3538597397101051<br>Best model saved to model/best_model.pkl<br>(base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml> |