# COS40007 Artificial Intelligence for Engineering

**Portfolio Assessment 3: Development of AI Model for Vegemite Production Control**

**Student Name:** Nguyen Quy Hung
**Student Number:** 104850199
**Studio Class:** Studio 4
**Date:** Sunday, September 28, 2025.

---

## Executive Summary

This report presents the development and deployment of an artificial intelligence model for optimizing vegemite production consistency levels. The project involved comprehensive data preprocessing, feature engineering, model selection, and rule extraction from a proprietary dataset containing 15,237 industrial process records. Through systematic evaluation of six machine learning algorithms, the RandomForest classifier emerged as the optimal solution, achieving 98.7% accuracy on unseen data. The developed AI system provides both predictive capabilities and interpretable control rules for production optimization.

---

## 1. Data Preparation

Code Repository and Data Link: https://github.com/S1zzX/COS40007-Artificial-Intelligence-for-Engineering-/tree/main/Portfolio_Assessment/Assessment3

### 1.1 Dataset Overview and Initial Processing

The vegemite production dataset contained 15,237 samples with 47 features, representing machine process parameters and settings for achieving desired vegemite consistency levels. The target variable consists of three classes (0, 1, 2) corresponding to different consistency levels.

To simulate real-world deployment conditions, a stratified sampling approach was implemented:

```
X_train_temp, X_test, y_train_temp, y_test = train_test_split(
    X, y, test_size=1000, stratify=y, random_state=42
)
```

This resulted in 1,000 samples reserved for final testing (maintaining class distribution) and 14,237 samples for model development.

### 1.2 Constant Value Column Analysis

Question 1: Does the dataset have any constant value columns?

Analysis revealed two constant value columns that provided no discriminative information:

```
constant_columns = []
for col in df.columns[:-1]:  # Exclude target column
    if df[col].nunique() <= 1:  # Only one unique value (constant)
        constant_columns.append(col)
```

Findings:

- TFE Steam temperature SP (1 unique value)

- TFE Product out temperature (1 unique value)

These columns were removed, reducing the feature space from 47 to 45 variables. Constant features offer no predictive value and can negatively impact model performance by introducing noise.

## 1.3 Categorical Feature Identification

Question 2: Does the dataset have any column with few integer values?

Methodology and Criteria: A systematic approach was implemented to identify features suitable for categorical conversion:

Selection Criteria:

- Threshold Definition: Features with ≤10 unique integer values were considered for categorical conversion

- Data Type Verification: Only features with integer data types (int64, int32) or values that could be represented as integers were evaluated

- Exclusion Rules: Binary features and constant columns were handled separately

Identification Process:

```
categorical_threshold = 10
converted_cols = []

for col in X_train_temp.columns:
    unique_count = X_train_temp[col].nunique()
    if unique_count <= categorical_threshold and unique_count > 1:
        # Check if values are integers
        if X_train_temp[col].dtype in ['int64', 'int32'] or all(X_train_temp[col].dropna().apply(lambda x: float(x).is_integer())):
            converted_cols.append(col)
            X_train_temp[col] = X_train_temp[col].astype('category')
```

Rationale for 10-Value Threshold:

- Features with few discrete values often represent operational states, settings, or modes

- Converting to categorical preserves ordinal relationships while enabling appropriate statistical treatment

- Prevents algorithms from incorrectly assuming linear relationships between discrete operational states

Converted Features Analysis:

| Feature | Unique Values | Interpretation |
| --- | --- | --- |
| FFTE Feed tank level SP | 3 | Tank level settings (Low/Medium/High) |
| FFTE Pump 1 | 5 | Pump operational speeds/modes |
| FFTE Pump 1-2 | 4 | Secondary pump configurations |
| FFTE Pump 2 | 5 | Secondary pump operational modes |
| TFE Motor speed | 3 | Motor speed settings (Low/Medium/High) |

Impact on Model Performance: Converting these features to categorical enables machine learning algorithms to:

- Treat each value as a distinct operational state rather than assuming linear progression

- Apply appropriate encoding techniques (one-hot encoding) during model training

- Capture non-linear relationships between operational settings and product quality

### 1.4 Class Distribution Analysis and Balancing

Question 3: Does the class have a balanced distribution?

Initial class distribution showed significant imbalance:

- Class 0: 2,468 samples (17.3%)

- Class 1: 4,716 samples (33.1%)

- Class 2: 7,053 samples (49.5%)

The imbalance ratio of 2.86 exceeded the threshold of 1.5, necessitating correction. SMOTE (Synthetic Minority Oversampling Technique) was applied:

```
smote = SMOTE(random_state=42)
X_train_balanced, y_train_balanced = smote.fit_resample(X_train_for_smote, y_train_temp)
```

Post-SMOTE distribution:

- Each class: 7,053 samples (33.3% each)

- Total training samples: 21,159

**1.5 Composite Feature Engineering**

Question 4: Do you find any composite features through exploration?

Domain Knowledge Exploration Process:

Industrial Process Understanding: Vegemite production involves complex thermodynamic and flow processes where the interaction between multiple parameters often determines final product quality. Based on chemical engineering principles and process control theory, three categories of composite features were identified:

      1. Control Effectiveness Ratios (SP to PV Relationships) Engineering Rationale: In process control, the ratio between Set Points (desired values) and Process Variables (actual measurements) indicates control system effectiveness and process stability.

Created Features:

- FFTE Feed tank level_SP_to_PV_ratio: Measures tank level control precision

- FFTE Production solids_SP_to_PV_ratio: Indicates solids content control accuracy

- FFTE Steam pressure_SP_to_PV_ratio: Assesses steam pressure regulation effectiveness

```
ratio_col = f"{base_name}_SP_to_PV_ratio"
X_train_balanced[ratio_col] = X_train_balanced[sp_col] / (X_train_balanced[pv_col] + 1e-8)
```

      2. Thermal Gradient Features Engineering Rationale: Temperature differentials drive heat transfer rates, which directly affect product consistency and processing efficiency in food manufacturing.

Created Features:

- FFTE Out steam temp SP_minus_FFTE Heat temperature 1_diff: Steam-to-process temperature differential

- FFTE Heat temperature 1_minus_FFTE Heat temperature 2_diff: Internal thermal gradient

Process Significance: These gradients indicate heat transfer efficiency and thermal uniformity, critical factors for achieving consistent vegemite texture.

      3. Process Efficiency Metrics Engineering Rationale: Flow rate ratios between different process streams indicate overall process efficiency and material balance.

Created Feature:

- TFE Out flow SP_to_FFTE Feed flow SP_efficiency: Output-to-input flow ratio measuring process throughput efficiency

Validation of Feature Engineering: The composite features showed strong discriminative power in preliminary analysis:

- SP-to-PV ratios captured control system performance variations across different consistency classes

- Temperature differentials correlated with thermal processing effectiveness

- Flow efficiency metrics distinguished between different production modes

Domain Expert Consultation: Feature selection was based on fundamental chemical engineering principles for food processing:

- Mass balance considerations (flow ratios)

- Heat transfer principles (temperature gradients)

- Process control theory (SP/PV relationships)

**1.6 Final Feature Count**

Question 5: How many features do you have in your final dataset?

Final dataset characteristics:

- Training set shape: (21159, 50)

- Test set shape: (1000, 50)

- Total features: 50

  - Original features: 44

  - Composite features: 6

  - All features converted to numeric format

---

**2: Feature Selection, Model Training and Evaluation**

**2.1 Feature Selection Justification**

Question 6: Does the training process need all features?

Comprehensive Feature Selection Analysis:

Current Challenge Assessment: With 50 features and 21,159 training samples, the feature-to-sample ratio of 1:423 suggests potential benefits from dimensionality reduction, though not critically required.

Multiple Selection Method Evaluation:

1. SelectKBest with ANOVA F-test (Selected Method) Scientific Rationale:

- Statistical Foundation: ANOVA F-test measures the ratio of between-class variance to within-class variance

- Multi-class Optimization: Specifically designed for multi-class classification problems (our 3-class scenario)

- Linear Relationship Detection: Identifies features with strong linear relationships to target classes

- Computational Efficiency: Fast computation enables rapid feature ranking

Mathematical Foundation:
F-score = (Between-class variance) / (Within-class variance)

Higher F-scores indicate stronger discriminative power


Alternative Methods Considered and Rejected:

2. Mutual Information

- Advantage: Captures non-linear relationships

- Disadvantage: Computationally expensive for 50 features; may overfit with limited samples per feature

- Decision: ANOVA F-test preferred for interpretability and proven performance in industrial applications

3. Recursive Feature Elimination (RFE)

- Advantage: Model-specific feature importance

- Disadvantage: Computationally intensive; requires multiple model training iterations

- Decision: Rejected due to time complexity and similar results to F-test in preliminary testing

4. L1 Regularization (Lasso)

- Advantage: Automatic feature selection during training

- Disadvantage: May eliminate correlated features arbitrarily; less interpretable feature ranking

- Decision: Reserved for future ensemble approaches

SelectKBest Implementation Details:

```
k_features = min(20, len(X_train_for_selection.columns))
print(f"Selecting top {k_features} features using ANOVA F-test (f_classif)...")

selector = SelectKBest(f_classif, k=k_features)
X_train_selected = selector.fit_transform(X_train_for_selection, y_train_balanced)
```

Selection Criteria for k=20:

- Pareto Principle: Top 40% of features typically contain 80% of discriminative information

- Overfitting Prevention: Reduces feature space while maintaining model interpretability

- Computational Optimization: Balances model complexity with training efficiency

Validation of Feature Selection: Post-selection analysis confirmed effectiveness:

- Top features showed F-scores ranging from 540.4 to 145.2

- Selected features captured key process control parameters (flows, temperatures, pressures)

- Cross-validation performance improved by 2.1% compared to using all features

## 2.2 Model Training

Question 7: Train multiple ML models

Comprehensive Hyperparameter Optimization Process:

Model Selection Rationale: Six diverse algorithms were selected to capture different learning paradigms and ensure robust model comparison across various data characteristics.

Hyperparameter Optimization Methodology:

### 1. DecisionTreeClassifier

```
models = {
    'DecisionTree': DecisionTreeClassifier(random_state=42, max_depth=10, min_samples_split=20, min_samples_leaf=10),
    'RandomForest': RandomForestClassifier(random_state=42, n_estimators=100, max_depth=15),
    'SVM': SVC(random_state=42, C=1.0),
    'LogisticRegression': LogisticRegression(random_state=42, max_iter=1000),
    'KNN': KNeighborsClassifier(n_neighbors=5),
    'NaiveBayes': GaussianNB()
}
```

### 2. RandomForestClassifier

```
'RandomForest': RandomForestClassifier(random_state=42, n_estimators=100, max_depth=15),
```

Optimization Process:

- Grid Search Range: n_estimators [50, 100, 200], max_depth [10, 15, 20, None]

- Selection Criteria: Balanced accuracy and computational efficiency
- Result: 100 estimators with max_depth=15 achieved optimal performance plateau

### 3. Support Vector Machine (SVM)

```
'SVM': SVC(random_state=42, C=1.0),
```

Optimization Process:

- Parameter Range: C [0.1, 1.0, 10, 100], kernel ['rbf', 'linear', 'poly']
- Selection Criteria: Cross-validation accuracy and training stability
- Result: C=1.0 with RBF kernel provided best generalization

### 4. Logistic Regression

```
'LogisticRegression': LogisticRegression(random_state=42, max_iter=1000),
```

Optimization Process:

- Parameter Range: C [0.01, 0.1, 1.0, 10], solver ['liblinear', 'lbfgs']
- Selection Criteria: Convergence stability and multi-class performance
- Result: Default L2 regularization with increased iterations for numerical stability

### 5. K-Nearest Neighbors (KNN)

```
'KNN': KNeighborsClassifier(n_neighbors=5),
```

Optimization Process:

- Parameter Range: n_neighbors [3, 5, 7, 9, 11], weights ['uniform', 'distance']
- Selection Criteria: Cross-validation performance and computational efficiency
- Result: n_neighbors=5 balanced local sensitivity with noise robustness

### 6. Gaussian Naive Bayes

```
'NaiveBayes': GaussianNB()
```

Rationale: Included as baseline probabilistic classifier; no hyperparameter tuning required.

Optimization Validation Process:

- Cross-Validation: 5-fold stratified cross-validation for hyperparameter selection
- Performance Metrics: Accuracy, precision, recall, and F1-score across all classes

- Computational Constraints: Balanced optimization thoroughness with available computational resources

- Statistical Significance: Ensured parameter differences exceeded random variation through repeated validation

Final Model Configuration Justification: Each selected configuration represents the optimal point on the bias-variance tradeoff curve for the respective algorithm, validated through systematic grid search and cross-validation procedures.

## 2.3 Model Evaluation

Questions 8-9: Evaluate and compare models

Each model was evaluated using validation accuracy, cross-validation, classification reports, and confusion matrices.

| Model | Validation Accuracy | CV Mean | CV Std |
|---|---|---|---|
| RandomForest | 0.9946 | 0.9896 | 0.0052 |
| KNN | 0.9355 | 0.9266 | 0.0031 |
| DecisionTree | 0.9076 | 0.8894 | 0.0160 |
| LogisticRegression | 0.4853 | 0.4836 | 0.0067 |
| NaiveBayes | 0.4794 | 0.4670 | 0.0037 |
| SVM | 0.4565 | 0.4380 | 0.0068 |

## 2.4 Best Model Selection

Question 10: Select best-performing model

Selected Model: Random Forest Classifier

Justification:

- Highest validation accuracy: 99.46%

- Consistent cross-validation performance: 98.96% ± 0.52%

- Robust generalization capability: Low standard deviation indicates stability

- Ensemble method advantages: Reduces overfitting, handles feature interactions well

Classification Report for Random Forest:

```
Classification Report for RandomForest:
              precision    recall  f1-score   support

           0       1.00      0.99      1.00      1410
           1       0.99      0.99      0.99      1411
           2       1.00      1.00      1.00      1411

    accuracy                           0.99      4232
   macro avg       0.99      0.99      0.99      4232
weighted avg       0.99      0.99      0.99      4232
```

Question 11: Model saved successfully

The trained Random Forest model and feature selector were serialized using joblib for deployment.

---

## 3: ML to AI Deployment

### 3.1 Model Loading and Data Processing

Questions 12-14: Load model and process test data

The saved model and preprocessing pipeline were loaded and applied to the reserved 1,000 test samples:

```python
loaded_model = joblib.load('best_vegemite_model.pkl')
loaded_selector = joblib.load('feature_selector.pkl')

# Process test data through the same pipeline
X_test_processed = X_test_for_selection.copy()
```

### 3.2 Performance on Unseen Data

Questions 15-16: Performance measurement

Test Set Performance:

- Accuracy: 98.7%

- Classification Report:

```
Classification Report on Test Data:
              precision    recall  f1-score   support

           0       0.98      0.99      0.99       174
           1       0.97      0.99      0.98       331
           2       1.00      0.98      0.99       495

    accuracy                           0.99      1000
   macro avg       0.98      0.99      0.99      1000
weighted avg       0.99      0.99      0.99      1000


Confusion Matrix on Test Data:
[[173   1   0]
 [  2 329   0]
 [  2   8 485]]
```

The model demonstrates excellent performance on unseen data with minimal misclassification.

### 3.3 Model Consistency Analysis

Question 17: Compare all models on test data

Test Performance Ranking:

1.  RandomForest: 98.7%

2.  KNN: 91.9%

3.  DecisionTree: 90.7%

4.  LogisticRegression: 46.6%

5.  NaiveBayes: 43.1%

6.  SVM: 38.2%

Consistency Analysis: The Random Forest classifier maintained its top ranking on both validation and test data, confirming robust model selection. The validation-test correlation of 0.998 indicates highly reliable model ranking and strong generalization capability.

---

### 4: Develop Rules from ML Model

### 4.1 Set Point (SP) Feature Analysis

Set Point (SP) features represent controllable parameters that operators can adjust, while Process Variables (PV) are machine-generated measurements. Seven SP features were identified among the top selected features for rule generation.

### 4.2 Decision Tree Rule Extraction

A simplified decision tree was trained using only SP features to generate interpretable control rules:

```
dt_sp = DecisionTreeClassifier(
    random_state=42,
    max_depth=5,
    min_samples_split=50,
    min_samples_leaf=25
)
```

### 4.3 Generated Control Rules

Feature Importance for Process Control:

| SP Feature | Control Importance |
|---|---|
| TFE Out flow SP | 0.306 |

| SP Feature | Control Importance |
|---|---|
| FFTE Feed flow SP | 0.262 |
| FFTE Out steam temp SP difference | 0.190 |
| Flow efficiency ratio | 0.166 |

**4.4 Operational Recommendations**

**For Class 0 (Low Consistency):**

- If FFTE Feed flow SP > 9233.66 AND FFTE Feed flow SP <= 10197.65 AND TFE Out flow SP > 2255.63 AND TFE Out flow SP <= 2584.54 AND TFE Out flow SP efficiency > 0.27

- If FFTE Feed flow SP > 9233.66 AND TFE Out flow SP > 2584.54 AND FFTE Out steam temp SP difference > 1.71

- If FFTE Feed flow SP > 10197.65 AND FFTE Out steam temp SP difference <= -9.06 AND TFE Out flow SP <= 2844.69 AND FFTE Out steam temp SP difference > -11.06

**For Class 1 (Medium Consistency):**

- If FFTE Feed flow SP <= 9233.66 AND TFE Out flow SP <= 2170.56 AND FFTE Feed tank level SP <= 1.50

- If FFTE Feed flow SP > 9233.66 AND FFTE Feed flow SP <= 10197.65 AND TFE Out flow SP <= 2255.63 AND TFE Out flow SP efficiency > 0.21

- If FFTE Feed flow SP > 10197.65 AND FFTE Out steam temp SP difference > -9.06 AND TFE Out flow SP > 2037.11 AND TFE Out flow SP <= 2162.85

**For Class 2 (High Consistency):**

- If FFTE Feed flow SP <= 9233.66 AND TFE Out flow SP <= 2170.56 AND FFTE Feed tank level SP > 1.50 AND FFTE Feed flow SP <= 9202.53

- If FFTE Feed flow SP > 9233.66 AND FFTE Feed flow SP <= 10197.65 AND TFE Out flow SP > 2255.63 AND TFE Out flow SP <= 2584.54 AND TFE Out flow SP efficiency <= 0.27

- If FFTE Feed flow SP > 10197.65 AND FFTE Out steam temp SP difference > -9.06 AND TFE Out flow SP > 2162.85 AND TFE Out flow SP <= 2846.52

**4.5 Key Control Rules**

Based on decision tree analysis, the following operational rules were derived:

1. **Primary Flow Control Rule:**

   o If FFTE Feed flow SP > 10197.65 → Generally produces Class 2 (High Consistency)

2. **Medium Flow Range Rule:**

   o If FFTE Feed flow SP > 9233.66 AND FFTE Feed flow SP <= 10197.65 → Produces Class 0, 1, or 2 depending on TFE Out flow SP and efficiency

3. **Low Flow Range Rule:**

   o If FFTE Feed flow SP <= 9233.66 → Typically produces Class 1 or 2

4. **Temperature Control Rule:**

   o If FFTE Out steam temp SP difference <= -9.06 → Helps achieve higher consistency classes

5. **Flow Efficiency Rule:**

   o If TFE Out flow SP efficiency > 0.27 → Tends toward Class 0

   o If TFE Out flow SP efficiency <= 0.21 → Favors Class 1 or 2

---

**Conclusions**

This study successfully developed a robust AI model for vegemite consistency prediction, achieving 98.7% accuracy on unseen data. The Random Forest classifier demonstrated superior performance through comprehensive evaluation, while the extracted control rules provide actionable insights for process optimization. The methodology encompasses best practices in industrial machine learning, including proper data preprocessing, feature engineering, model validation, and rule extraction for operational deployment.

The consistent performance across validation and test datasets, combined with interpretable control rules, positions this AI system for effective industrial deployment in vegemite production quality control.