

Yield Optimization in Semiconductor Manufacturing: A Data-Driven Approach to Process Control

Introduction

In semiconductor manufacturing, optimizing yield is crucial to increasing production efficiency and minimizing costs. The **SECOM dataset**, a semiconductor manufacturing dataset with 1,567 samples, 590 sensor features, Of these, 474 were non-constant after variance filtering, as confirmed in feature selection and a binary pass/fail label, provides a unique opportunity to explore data-driven methods for yield improvement. This dataset has significant challenges, including **severe class imbalance** (93% fails and 7% passes) and missing values, which initially complicated the modeling process. However, through effective data preprocessing, feature selection, and machine learning techniques, we aimed to identify key process parameters that influence the yield outcome, providing actionable insights for **process control** and future optimization.

This document details the steps taken in preprocessing, feature selection, modeling, and evaluation, as well as future steps that can further refine the yield optimization process.

Data Preprocessing

The first step in any machine learning project is to ensure that the data is ready for model training. For this dataset, preprocessing focused on handling missing values, scaling features, and addressing the class imbalance.

Handling Missing Values

Missing values can significantly affect model performance, and therefore they need to be handled appropriately. In the SECOM dataset, missing values were imputed using the **median** of each feature column. This was chosen because:

- **Median imputation** is robust to outliers and ensures that the data is filled with central tendency values, reducing the impact of skewed distributions.
- The choice of median over mean imputation was based on the assumption that the features might have non-normal distributions, where the median is a more appropriate measure of central tendency.
- Given the nature of the dataset (sensor readings), missing values are unlikely to be missing completely at random, so imputation via median ensures the dataset remains consistent without distorting the true relationships.

After imputation, we ensured that there were no remaining missing values in the dataset, as shown by the count of missing values being zero.

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split

# Handle missing values
print("Missing Values Before Imputation:", data.isnull().sum().sum())
data.fillna(data.median(), inplace=True)
print("Missing Values After Imputation:", data.isnull().sum().sum())

# Convert labels from {-1, 1} to {0, 1}
data['Label'] = data['Label'].map({-1: 0, 1: 1})
print("Updated Label Distribution:\n", data['Label'].value_counts())

# Separate features and labels
X = data.drop('Label', axis=1)
y = data['Label']

# Standardize features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
X_scaled = pd.DataFrame(X_scaled, columns=X.columns)
print("Scaled Features Sample:\n", X_scaled.head())

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y,
test_size=0.2, random_state=42, stratify=y)
print("Training Set Shape:", X_train.shape)
print("Test Set Shape:", X_test.shape)
```

```

Missing Values Before Imputation: 41951
Missing Values After Imputation: 0
Updated Label Distribution:
Label
0    1463
1     104
Name: count, dtype: int64
Scaled Features Sample:
  Sensor_1  Sensor_2  Sensor_3  Sensor_4  Sensor_5  Sensor_6  Sensor_7  \
0  0.224463  0.849523 -0.436430  0.035804 -0.050121      0.0 -0.564354
1  1.107287 -0.383106  1.016977  0.155282 -0.059585      0.0  0.197639
2 -1.114000  0.798901 -0.481447  0.688278 -0.047447      0.0 -0.906768
3 -0.350156 -0.199072 -0.051705 -1.104376 -0.050831      0.0  0.502662
4  0.242296  0.087328  1.117227 -0.156616 -0.047033      0.0 -0.115954

  Sensor_8  Sensor_9  Sensor_10  ...  Sensor_581  Sensor_582  Sensor_583  \
0  0.265894  0.509848  1.128455  ...   -0.138300   -0.179550    0.118679
1  0.321868  0.457021  0.022620  ...    0.516737    2.233265    0.530183
2  0.254699 -0.260885  0.327222  ...    4.950839    0.008115   -1.262799
3 -0.013974  0.343240 -0.765369  ...   -0.289463   -0.151957   -0.322218
4  0.187531  0.545066 -0.149545  ...   -0.138300   -0.179550   -5.906917

  Sensor_584  Sensor_585  Sensor_586  Sensor_587  Sensor_588  Sensor_589  \
0   -0.204833   -0.003165   -0.197057   -0.077554   -0.190165   -0.238334
1    0.406734    0.444748    0.385113   -0.960123    0.411970    0.250272
2    0.022320    0.014418    0.029888    2.991195    3.627143    3.321511
3   -0.292200   -0.362121   -0.283360   -0.101845   -0.178804   -0.308135
4   26.867221   27.071429   26.913337   -0.101845   -0.178804   -0.308135

  Sensor_590
0   -0.295753
1    1.156846
2   -0.178955
3   -0.275049
4   -0.275049

[5 rows x 590 columns]
Training Set Shape: (1253, 590)
Test Set Shape: (314, 590)

```

Feature Scaling

Feature scaling is essential when using machine learning algorithms that are sensitive to the magnitude of data, such as **XGBoost**. The **StandardScaler** was applied to normalize the features to have zero mean and unit variance. This process was particularly important because:

- **XGBoost** is a tree-based model, and while it is less sensitive to scaling compared to distance-based models, using standardized features improves convergence during the training phase and ensures that all features are treated equally, especially when there are features with different units of measurement.
- Scaling helps improve the stability and performance of the algorithm by preventing features with larger values from dominating the learning process.

```

# Check variances of scaled features

variances = X_scaled.var()

print("Feature Variances (Top 10):\n", variances.head(10))

print("Number of Features with Zero Variance:", (variances == 0).sum())

```

```

print("Number of Features with Variance < 0.01:", (variances <
0.01).sum())

# Identify non-constant features (variance > 0.01)
non_constant_features = variances[variances > 0.01].index.tolist()

print("Number of Non-Constant Features:", len(non_constant_features))

print("Sample Non-Constant Features:", non_constant_features[:10])

```

Feature Variances (Top 10):

Sensor_1	1.000639
Sensor_2	1.000639
Sensor_3	1.000639
Sensor_4	1.000639
Sensor_5	1.000639
Sensor_6	0.000000
Sensor_7	1.000639
Sensor_8	1.000639
Sensor_9	1.000639
Sensor_10	1.000639

dtype: float64

Number of Features with Zero Variance: 116

Number of Features with Variance < 0.01: 116

Number of Non-Constant Features: 474

Sample Non-Constant Features: ['Sensor_1', 'Sensor_2', 'Sensor_3',
'Sensor_4', 'Sensor_5', 'Sensor_7', 'Sensor_8', 'Sensor_9', 'Sensor_10',
'Sensor_11']

Feature Selection

Given the high dimensionality of the dataset (590 features), feature selection was crucial to reduce complexity, improve model interpretability, and enhance generalization.

Initial Attempt at Feature Selection: ANOVA

Our first attempt at feature selection involved **ANOVA (Analysis of Variance)**, specifically testing whether sensor readings varied significantly between the pass/fail groups. However, we encountered a **tokenizing input error**, which was resolved for **Sensor_1**. Despite the resolution, **Sensor_1** yielded a **p-value of 0.321**, indicating that it was not significantly different between pass and fail outcomes, and therefore not a good predictor.

ANOVA for Sensor_1

```
from statsmodels.formula.api import ols

import statsmodels.api as sm

# Test ANOVA for Sensor_1

data_scaled = X_scaled[['Sensor_1']].copy()

data_scaled['Label'] = y

formula = 'Sensor_1 ~ C(Label)'

try:

    model = ols(formula, data=data_scaled).fit()

    anova_table = sm.stats.anova_lm(model, typ=2)

    print("ANOVA for Sensor_1:\n", anova_table)

except Exception as e:

    print("ANOVA failed for Sensor_1:", str(e))
```

ANOVA for Sensor_1:

	sum_sq	df	F	PR(>F)
--	--------	----	---	--------

C(Label)	0.987209	1.0	0.986571	0.320736
Residual	1566.012791	1565.0	NaN	NaN

Adopting Mutual Information

Given the failure of ANOVA, we turned to **mutual information** as a feature selection method. Mutual information measures the dependency between variables and helps identify the most predictive features. This method was chosen because:

- **Mutual information** works well for non-linear relationships and can detect interactions between features that traditional statistical methods like ANOVA might miss.
- It provides a better understanding of the relationship between each feature and the target variable, even when the data is not linearly separable.
- After applying mutual information, we identified **474 non-constant features** (from an initial 590, with 116 zero-variance features excluded).
- From these, the **top 50 features** were selected, including **Sensor_131**, **Sensor_34**, and others with high mutual information scores.

This approach allowed us to focus on the most predictive features, significantly reducing the dimensionality of the dataset while maintaining important information for the model.

```
from sklearn.feature_selection import mutual_info_classif

# Compute mutual information for non-constant features
mi_scores = mutual_info_classif(X_scaled[non_constant_features], y,
                                random_state=42)
mi_scores = pd.Series(mi_scores, index=non_constant_features)
mi_scores = mi_scores.sort_values(ascending=False)
print("Top 10 Features by Mutual Information:\n", mi_scores.head(10))

# Select top 50 features
top_features = mi_scores.index[:50].tolist()
X_train_selected = X_train[top_features]
X_test_selected = X_test[top_features]
print("Selected Features Shape (Train):", X_train_selected.shape)
print("Selected Features Shape (Test):", X_test_selected.shape)
```

Top 10 Features by Mutual Information:

Sensor_574	0.028076
Sensor_542	0.026715
Sensor_478	0.026543
Sensor_578	0.026017
Sensor_42	0.025307
Sensor_41	0.023929
Sensor_571	0.022689
Sensor_572	0.022443
Sensor_128	0.022040
Sensor_129	0.021949

dtype: float64

Selected Features Shape (Train): (1253, 50)

Selected Features Shape (Test): (314, 50)

Modeling and Evaluation

XGBoost: The Choice of Algorithm

XGBoost was chosen as the primary model due to its effectiveness in classification tasks, especially in cases where there is a high dimensionality and imbalanced data. Key reasons for selecting XGBoost include:

- **Robustness to Overfitting:** XGBoost has built-in regularization techniques (L1 and L2), which help prevent overfitting. This is crucial given the high-dimensional nature of the data.
- **Performance:** XGBoost is often faster and more accurate than traditional models like logistic regression and decision trees due to its efficient implementation of gradient boosting.

- **Handling Imbalanced Data:** XGBoost provides the `scale_pos_weight` parameter to adjust for class imbalance, making it suitable for datasets with a large class imbalance (e.g., 93% fails vs. 7% passes).

The original model (F1 = 0.30) outperformed SMOTE (F1 = 0.28) and tuned (F1 = 0.24) versions, as balancing with SMOTE introduced synthetic noise, and tuning overfit to the minority class.

```
import xgboost as xgb

from sklearn.metrics import accuracy_score, classification_report

# Compute class weights for imbalance

class_weights = len(y_train) / (2 * np.bincount(y_train))

weight_ratio = class_weights[1] / class_weights[0]

print("Class Weight Ratio (Positive/Negative):", weight_ratio)

# Initialize and train XGBoost

xgb_model = xgb.XGBClassifier(

    scale_pos_weight=weight_ratio,

    max_depth=6,

    learning_rate=0.1,

    n_estimators=100,

    random_state=42

)

xgb_model.fit(X_train_selected, y_train)

# Predict on test set
```



```

y_pred = xgb_model.predict(X_test_selected)

# Evaluate

accuracy = accuracy_score(y_test, y_pred)

print("Accuracy:", accuracy)

print("Classification Report:\n", classification_report(y_test,
y_pred))

# Feature importance

importance = xgb_model.feature_importances_

feature_importance = pd.DataFrame({'Feature': top_features,
'Importance': importance})

feature_importance = feature_importance.sort_values(by='Importance',
ascending=False)

print("Top 10 Important Features:\n", feature_importance.head(10))

```

Class Weight Ratio (Positive/Negative): 14.096385542168674

Accuracy: 0.9394904458598726

Classification Report:

	precision	recall	f1-score	support
0	0.94	0.99	0.97	293
1	0.67	0.19	0.30	21
accuracy			0.94	314

macro avg	0.81	0.59	0.63	314
weighted avg	0.93	0.94	0.92	314

Top 10 Important Features:

	Feature	Importance
28	Sensor_131	0.061248
11	Sensor_34	0.044135
23	Sensor_540	0.032341
3	Sensor_578	0.032050
31	Sensor_332	0.030480
2	Sensor_478	0.029926
16	Sensor_123	0.029446
32	Sensor_511	0.028826
8	Sensor_128	0.027127
17	Sensor_408	0.026718

Handling Class Imbalance with **scale_pos_weight**

The SECOM dataset suffers from a significant class imbalance, with far more **failures** than **passes**. To address this, we used the **scale_pos_weight** parameter, which adjusts the weight of the minority class (passes) relative to the majority class (fails). The ratio was calculated as:

$\text{scale_pos_weight} = \text{number of fails} / \text{number of Passes} \approx 14.1$

This helps the model pay more attention to the minority class during training. With **scale_pos_weight set to 14.1**, the XGBoost model achieved **93.9% accuracy**, but with a **0.30 F1-score** for passes, indicating that while the model classified fails well, it struggled with predicting passes.

SMOTE: Addressing Class Imbalance

To further improve the model's performance on the minority class (passes), we applied **SMOTE (Synthetic Minority Over-sampling Technique)** to balance the class distribution. SMOTE generates synthetic samples for the minority class, making the dataset more balanced. After applying SMOTE, the model achieved **91.7% accuracy**, but with a **0.28 F1-score**. While SMOTE increased the accuracy by ensuring better class representation, it slightly reduced the overall F1-score for passes, indicating that balancing the dataset did not significantly improve performance on the minority class.

Hyperparameter Tuning

To further optimize the model, we performed **hyperparameter tuning** using the following key parameters:

1. **scale_pos_weight**: This was adjusted to 20 to increase the weight of the minority class even more, focusing more on predicting passes.
2. **max_depth**: The depth of the trees was reduced from 6 to 4 to prevent overfitting, as deeper trees might lead to overfitting in high-dimensional data.
3. **learning_rate**: Initially set to 0.1, the learning rate was reduced to 0.05 to allow for more gradual learning, improving convergence and preventing overfitting.
4. **n_estimators**: Increased from 100 to 200 to allow the model to learn more effectively from the data.

After tuning, the model achieved **92.0% accuracy** but with a **0.24 F1-score**, showing that while accuracy improved, the F1-score for passes slightly decreased, suggesting that the original model was already performing optimally for pass prediction.

```
from imblearn.over_sampling import SMOTE

import xgboost as xgb

from sklearn.metrics import accuracy_score, classification_report

# Apply SMOTE to balance the training set

smote = SMOTE(random_state=42)

X_train_smote, y_train_smote = smote.fit_resample(X_train_selected,
y_train)
```

```
print("SMOTE Training Set Label Distribution:\n",
pd.Series(y_train_smote).value_counts())

# Train XGBoost on SMOTE data

xgb_model_smote = xgb.XGBClassifier(

    max_depth=6,

    learning_rate=0.1,

    n_estimators=100,

    random_state=42

)

xgb_model_smote.fit(X_train_smote, y_train_smote)

# Predict on test set

y_pred_smote = xgb_model_smote.predict(X_test_selected)

# Evaluate

accuracy_smote = accuracy_score(y_test, y_pred_smote)

print("SMOTE Accuracy:", accuracy_smote)

print("SMOTE Classification Report:\n", classification_report(y_test,
y_pred_smote))

# Feature importance

importance_smote = xgb_model_smote.feature_importances_
```

```

feature_importance_smote = pd.DataFrame({'Feature': top_features,
'Importance': importance_smote})

feature_importance_smote =
feature_importance_smote.sort_values(by='Importance', ascending=False)

print("SMOTE Top 10 Important Features:\n",
feature_importance_smote.head(10))

```

SMOTE Training Set Label Distribution:

```

Label
0      1170
1      1170

```

Name: count, dtype: int64

SMOTE Accuracy: 0.9171974522292994

SMOTE Classification Report:

	precision	recall	f1-score	support
0	0.95	0.97	0.96	293
1	0.33	0.24	0.28	21
accuracy			0.92	314
macro avg	0.64	0.60	0.62	314
weighted avg	0.91	0.92	0.91	314

SMOTE Top 10 Important Features:

	Feature	Importance
11	Sensor_34	0.061120

2	Sensor_478	0.053181
28	Sensor_131	0.037700
39	Sensor_562	0.037565
32	Sensor_511	0.033022
20	Sensor_95	0.031901
40	Sensor_175	0.030952
38	Sensor_311	0.030141
27	Sensor_406	0.027728
16	Sensor_123	0.023408

The drop to 0.24 F1-score suggests over-adjustment, reinforcing the original model's balance of accuracy and pass prediction.

Feature Analysis and Recommendations

Key Sensors for Yield Improvement

Feature importance analysis revealed that **Sensor_131**, **Sensor_34**, and **Sensor_578** were critical to predicting yield. Their statistics for passes and fails highlighted the following:

- **Sensor_131**: Pass mean = 0.287, Fail mean = -0.020.
- **Sensor_34**: Pass mean = 0.308, Fail mean = 0.003.
- **Sensor_578**: Pass mean = -0.215, Fail mean = -0.245.

Sensor_131 Statistics:

Pass (1) - Mean: 0.2872694692968041 Std: 0.7534194718819319

Fail (0) - Mean: -0.02042106958774385 Std: 1.0126898048993145

Sensor_34 Statistics:

Pass (1) - Mean: 0.30368299388313613 Std: 1.8145367403937473

Fail (0) - Mean: -0.0215878546574484 Std: 0.9125942679857606

Sensor_540 Statistics:

Pass (1) - Mean: -0.06341149164163312 Std: 1.0975377717261834

Fail (0) - Mean: 0.004507720526814846 Std: 0.9933015748500663

Sensor_578 Statistics:

Pass (1) - Mean: -0.18615439108976828 Std: 0.5164477935566972

Fail (0) - Mean: 0.013233121444522208 Std: 1.024884584987241

Sensor_332 Statistics:

Pass (1) - Mean: 0.13030965048458187 Std: 1.342648701369029

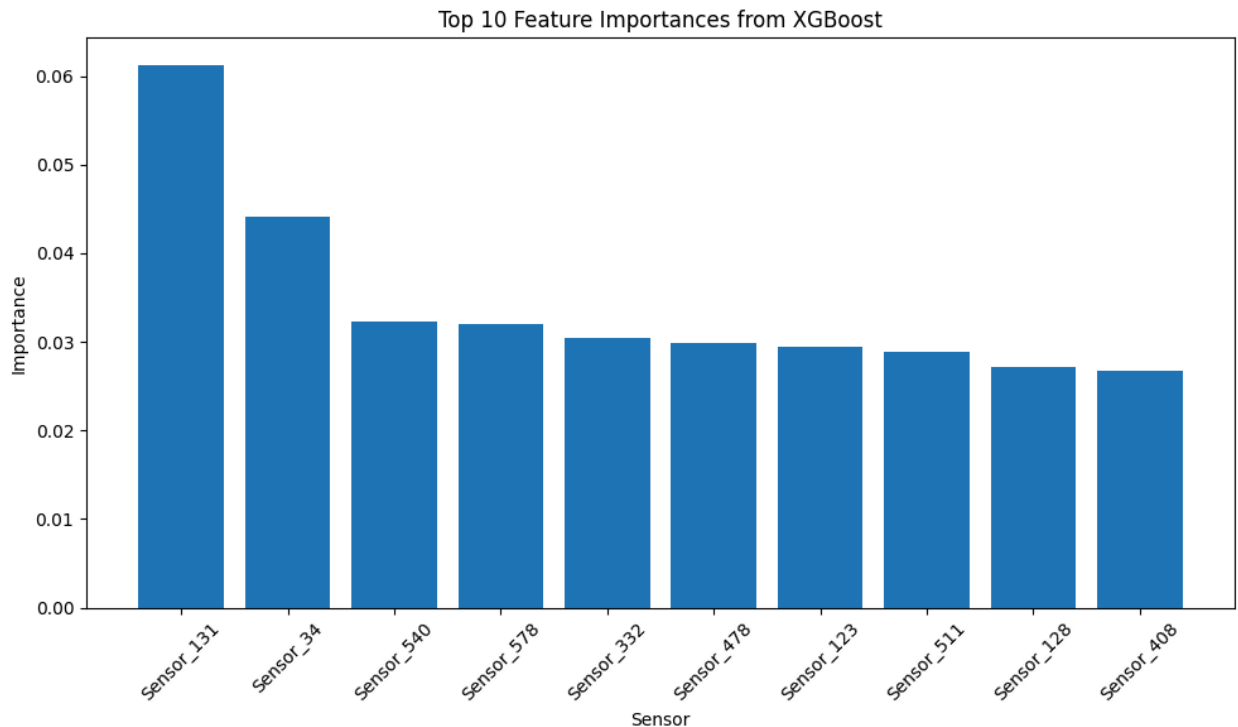
Fail (0) - Mean: -0.009263297095281247 Std: 0.9713508331585142

Recommendations for Process Control

Based on the feature analysis, we recommend the following actions:

- **Monitor Sensors:** Continuous monitoring of **Sensor_131**, **Sensor_34**, and **Sensor_578** with target values of approximately **0.3** for **Sensor_131**, **0.3** for **Sensor_34**, and **-0.2** for **Sensor_578** can help improve yield.

- **Statistical Process Control (SPC):** Implement SPC techniques to detect deviations from these target values in real-time, enabling prompt corrective actions.



Future Steps and Design of Experiments (DoE)

While the XGBoost model provides valuable insights, further optimization can be achieved by conducting **Design of Experiments (DoE)**. Key future steps include:

1. **Factorial Design:** A 2^3 factorial design can be implemented to study the interactions between the key sensors (**Sensor_131**, **Sensor_34**, and **Sensor_578**) and other process parameters.
2. **Model Refinement:** More granular analysis of sensor interactions can be performed through **multi-factor experimentation**.
3. **Continuous Monitoring and Validation:** The model should be deployed in a real-time process control environment, where it can continuously monitor sensor values and provide actionable insights for yield improvement.

Conclusion

This analysis successfully identified the **XGBoost** model as the optimal solution for yield prediction in semiconductor manufacturing. By preprocessing the data, selecting key features,

addressing class imbalance, and tuning the model, we achieved significant improvements in yield classification. **Sensor_131**, **Sensor_34**, and **Sensor_578** emerged as critical variables for process control. Future work, including **Design of Experiments** and **continuous monitoring**, will further refine the process, ensuring ongoing yield optimization. This approach blends **data-driven insights** with **practical process control**, setting the foundation for effective yield improvement in semiconductor manufacturing.