



# **Smart Factory Automation: Data Analytics using semiconductor factory dataset**

**AI/ML Advanced end-to-end Solution Package**

# Contents

1. Introduction
2. Project Approach
3. Implementation
4. Comparison
5. Executive Summary and Conclusion
6. References

# 1. Introduction:

This is a data set from a semi-conductor manufacturing process. It had 1567 Number of Instances & 591 number of sensor instances. Semiconductor manufacturing is an extraordinarily complex process with hundreds of steps involved in it. Each of these steps has a feedback signal coming from the sensors. Since semiconductor manufacturing is expensive so even a small failure in one of the steps might result in catastrophic loss hence the large number of sensors to detect any anomaly. Process engineers are required to detect an anomaly during the manufacturing as soon as possible because these products are manufactured in an environment to achieve great precision on a scale of nanometer, but the existing univariate and multivariate control charts fail to do so with such high volume of data taken during the manufacturing process. This end-to-end solution package is a generic solution for a Smart Factory Operation using the semiconductor data set. Industries already have a Planned Maintenance schedule in the manufacturing sector which informs us when to carry out the maintenance of machineries involved in manufacturing depending on the running hours of that particular machinery or the time between overhauls but we need to have a predictive maintenance scheduling (especially in the case of Smart Factory Manufacturing) which after learning from the data available, predicts about the failure before its occurrence and thus reducing the time & money lost during the process of manufacturing. In other words, if we are provided with observations that has been generated by an unknown stochastic dependency then our goal is to infer a law that will be able to predict the future observations correctly based on the same dependency to maintain a high process yield in manufacturing. This is a classification problem of predicting a pass or fail for the manufactured semiconductor wafer based on the feedback coming from the sensors. There is total 590 sensors and one more feature of pass and fail state. Using these 590 sensors total 1567 observations are collected and classification model is made using the binary response of these sensors.

First, we performed **imputation** method to clean the data like- removing non-value column or rows, removing the non-available value cell with certain value by using different imputation method.

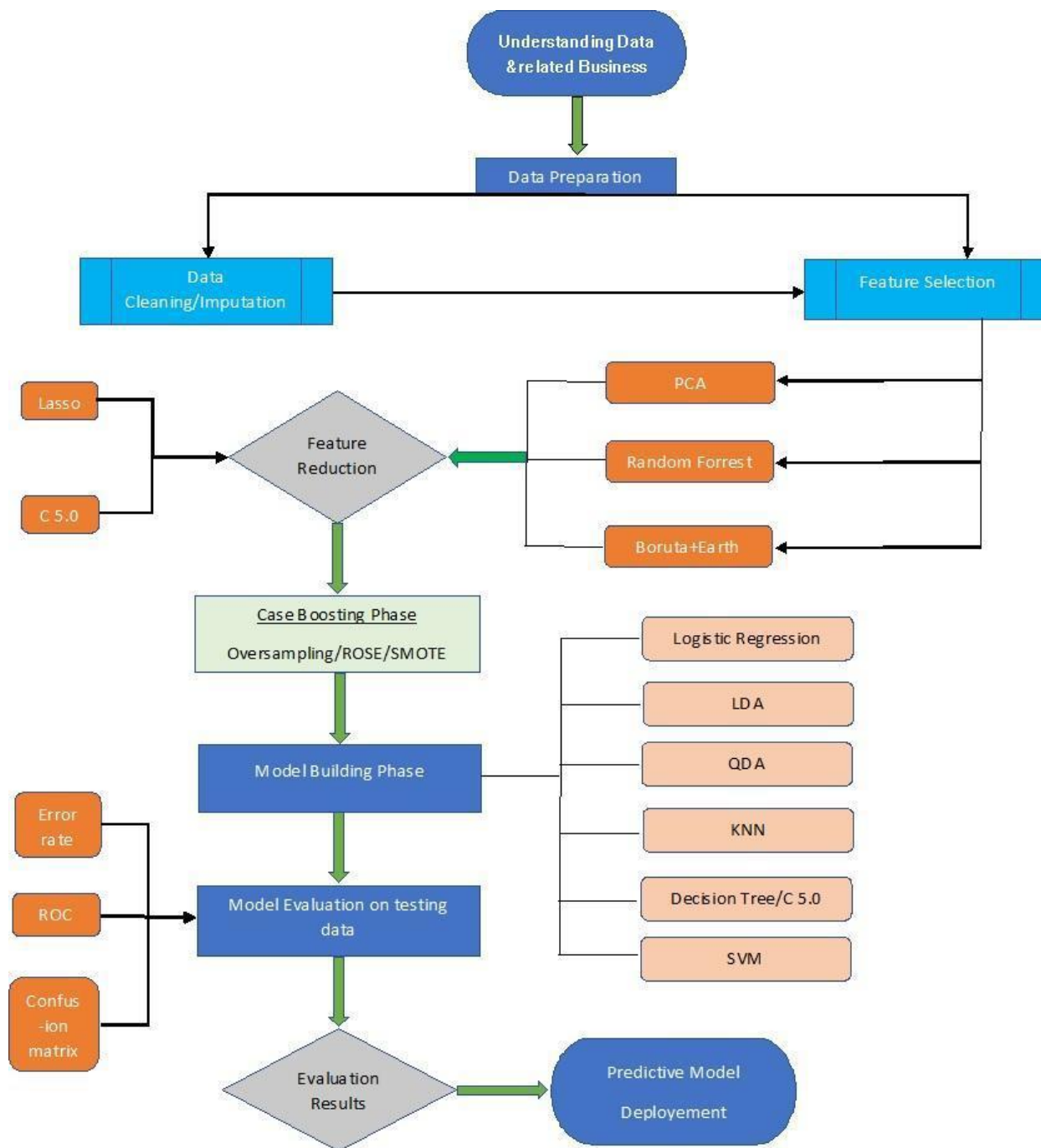
Second, the total observed data contains 1,463 pass cases with only 104 fail cases. so, to equalize these two different cases a **boosting** technique is devised to deal with highly imbalance between the pass and fail cases.

Third, **feature selection** was done using different models for selecting the important sensor for designing model.

Fourth, the **different methods** were made for making different models and then best model was implemented.

## 2. Project Approach:

### Flow Chart



We were provided with the data from a semiconductor manufacturing with a goal to classify the product in pass or fail class by creating a predictive model that follows the same dependency as the observations provided to us. With all the tools we have learned in and outside the class, we have approached in the following manner to tackle this problem: -

### 1. Data & Related Business Understanding:

The most important step in data analytics/machine learning is the first one, to understand the data and the underlying Business problem related to it. We know that cost, quality, and delivery time are key factors for firms to attain long term competition and specifically speaking, semiconductor industry is a capital-intensive market sector so the stakes in this case are high and there is a little margin for error (high precision required and high downtime penalty to be paid in case of failure). Therefore, creating an automatics & advanced process control method is required so that process engineer is aware of the failure that is yet to happen and take corrective action well before it occurs to reduce the downtime error and maintain a high process yield thus improving business technique.

### 2. Data Preparation:

The SECOM manufacturing data originally contained 1567 examples taken from a wafer fabrication production line with 590 manufacturing operation variables and 1 quality variable. There are only 104 fail cases among 1567 examples which is in a ratio of 1:14 so the classes are highly imbalanced.

Also, almost 4.5% of the values are missing so we cannot proceed further without cleaning our data. The Columns with more than 55% missing values are eliminated using “is.na” & “data.frame”. After removing these columns, we are left with a data set containing 3% of the missing values so we carried out data imputation techniques such as mean imputation, KNN imputation and MICE imputation to fill those missing values with randomly generated data by regressing on the data set multiple times.

### 3. Feature Selection:

With 590 features in the dataset, not all these signals coming from the sensors are equally valuable in a specific monitoring system. It contains relevant information as well as noise. To eliminate noise, it is important to carry out feature selection. Following steps were carried out for dimensionality reduction/feature selection: -

- All the features with more than 50% missing values were eliminated
- Principal Component Analysis was carried out on the imputed data and it was observed that only less than 100 features could explain most of the variances for the given data. Precisely speaking, 168 features explained 95% of the variance
- Random Forest was used on the data set and it was observed that only 27 features were important.
- Boruta is an all-relevant feature selection wrapper algorithm, capable of working with any classification method that output variable importance measure (VIM), by default

Boruta uses Random Forest. The method performs a top-down search for relevant features by comparing original attributes' importance with importance achievable at random, estimated using their permuted copies, and progressively eliminating irrelevant features to stabilize that test. This method provided us with 20 important features.

Reference: - <https://cran.r-project.org/web/packages/Boruta/Boruta.pdf>

### 1. Feature Reduction:

After knowing the important features, we need to select those features: -

Lasso: - This shrinkage method was used to equate the coefficients of those features which were not important equal to zero. For obtaining the optimum  $\lambda$  (lambda) value, cross validation was carried out and the result of cross validation was used as an input for Lasso's lambda value. The output of the Lasso shrinkage method will be used directly for model building as it would be cleaned, imputed as well as dimensionally reduced data set with much less noise.

C 5.0 This Package fit classification tree models or rule-based models using Quinlan's C5.0 algorithm and therefore helps us in identifying important predictors.

Reference: - <https://cran.r-project.org/web/packages/C50/C50.pdf>

### 1. Case Boosting Phase: -

This data presents us with an unusual form of an imbalanced case with a ratio of 1:14. In such cases models work well on the prediction of a majority class while fails on the prediction of a minority class. Therefore, we need to boost the minority class or in other words, we need to make the ratio as close to 1 as possible. For this case boosting, we use the following techniques: -

- Oversampling:**  
Oversampling duplicates data from the minority class by randomly sampling data with replacement from the same class
- SMOTE (Synthetic Minority Over-sampling Technique):**  
To shift the classifier bias towards the minority class, it artificially generates random set of minority class observations. For generating artificial data, bootstrapping and k-nearest neighbors are used.
- ROSE (Random Over Sampling Exercise):**  
ROSE also generates synthetic balanced samples which shifts the classifier bias towards minority class. The underlying concept of ROSE is also bootstrapping which aids classification in the presence of minority class.

After case boosting, we have a balanced data set on which we can start building our classification model.

1. Model Building: -

Since, we have a cleaned and balanced data set divided into training and testing and we know the important features, we can start building our classification model by using the following techniques we have learned in class:

i) Logistic Regression:

Since this is a binary classification problem so the intuition of using Logistic regression for classifying pass/fail was obvious. The data set was already divided into training and testing data. Logistic regression was applied on the training data and the results were tested on testing data. It shows an output with AUC 0.74

ii) LDA (Linear Discriminant Analysis):

Applying LDA to fit the model of pass/fail class on the signal coming from sensors (only the important features were used for generating this model) produced almost the same results giving us the idea that the variances of both the classes are almost equal and that the assumption that it follows normal distribution holds true.

iii) QDA (Quadratic Discriminant Analysis):

QDA was applied and the resulting models were compared based on confusion matrix/accuracy rate, ROC area and the classification error rate. It was observed that the LDA model outperforms QDA.

iv) Decision tree/C5.0:

It provided a result that can be interpreted well in a way that important features were quite obvious from the pruned tree, but the classification accuracy was not as high as compared to other classification technology

v) **SVM (Support Vector Machine):**

When SVM classification was applied for model building, it outperformed all the other classification techniques when the margin was varied and the optimum margin was selected for the model, i.e. a separating hyperplane with an optimum margin/cost/budget was created in the hidden feature space using polynomial programming (polynomial kernel) to find a unique solution.

All the models were compared, cross-validated and the best model was selected on the following criteria: -

- a) Confusion Matrix/Accuracy rate
- b) ROC area (AUC)
- c) Classification error rate

# Description of New Techniques:

## 1. IMPUTATION:

The dataset was obtained from a Semiconductor manufacturing process which has too many Missing Values and Redundant data. Hence to make the dataset usable we use a technique called Imputation. We used 3 types of imputation techniques on our dataset:

- ❖ KNN Imputation [9]
- ❖ MICE imputation
- ❖ MissForest Imputation

**KNN Imputation:** KNN imputation is used to substitute the missing data with its closest k neighbors provided the data is continuous, discrete, ordinal, or categorical. Certain properties must be kept in mind when using KNN imputation:

- The value of k: A lower value of k will increase the noise and the model will not be that generalizable. But a higher value of k will tend to diminish the properties of local neighbors which is exactly what we're looking for.
- The assigning method: kNN imputation uses Mean, Median and Mode of the k neighbors for numerical data and just the Mode for categorical data.

Usage:

```
library(DMwR)
data_Imputed = knnImputation(data_cleaned)
#this command creates a dataset data_Imputed which contains the new imputed dataset obtained using knnImputation.
```

**MICE Imputation:** MICE stands for Multivariate Imputation by Chained Equations. This method uses chained equations to perform imputation on missing data. It performs better while creating multiple imputations as compared to a single imputation. It assumes that all missing values are Missing at Random (MAR), i.e., the missing values can be predicted using the observed values. It imputes the dataset based on variable by variable and provides an imputation output for each of them. MICE can tackle continuous, discrete, binary variables.

The different methods used by this package are:

1. PMM (Predictive Mean Matching) – For numeric variables
2. Logreg (Logistic Regression) – For Binary Variables( with 2 levels)
3. Polyreg (Bayesian polytomous regression) – For Factor Variables ( $\geq 2$  levels)



#### 4. Proportional odds model (ordered, $\geq 2$ levels)

Usage:

```
library(mice)
data_IIpmm <- mice(data_c, m=1, maxit = 1, method = 'pmm', seed = 500)
data_II <- complete(data_IIpmm,1)
```

Here is an explanation of the parameters used:

1. m – Refers to 1 imputed data set
2. maxit – Refers to no. of iterations taken to impute missing values
3. method – Refers to method used in imputation. we used predictive mean matching(PMM).

**MissForest Imputation:** It uses random forest model which is trained on the observed value to predict the missing values in the data set. MissForest imputation is used when the data is of mixed type. random forest naturally constitutes a multiple imputation scheme by averaging over many unpruned regression or classification tree.

One of the advantage of MissForest imputation is that it yields an estimated out-of-bag(OOB) imputation error eliminating the need of a test set or an elaborate cross-validation technique. It also saves time by running in parallel. Additionally, MissForest exhibits attractive computational efficiency and can cope with high-dimensional data. we can impute continuous/categorical data having complex interactions or nonlinear relations.

Usage:

```
data_miss <- missForest(data_c)
data_III <- data_miss$ximp
```

## 2. CASE BOOSTING:

The total observed data contains 1,463 pass cases with only 104 fail cases. so, to equalize these two different cases a boosting technique is devised to deal with highly imbalance between the pass and fail cases. We have performed 3 method for case sampling-

- ❖ ROSE (Random Over Sampling Exercise)
- ❖ SMOTE (Synthetic Minority Over-sampling Technique)
- ❖ Oversampling/ Undersampling

**ROSE Boosting-** The function ROSE generates synthetic balanced samples and allows to strengthen the subsequent estimation of any binary classifier. ROSE (Random Over-Sampling Examples) is a bootstrap-based technique which aids the task of binary classification in the presence of rare classes. It handles both continuous and categorical data by generating synthetic examples from a conditional density estimate of the two classes.

Usage:

```
hacide.rose <- ROSE(cls ~ ., data=hacide.train, seed=123)$data #generating new balance data by ROSE boosting  
table(hacide.rose$cls) #check imbalance of new data
```

**SMOTE-** SMOTE is a well-known algorithm to fight the problem of unbalance cases in datasets. To counter such problem SMOTE algorithm creates artificial data based on feature space (rather than data space) similarities from minority samples. It generates a random set of minority class observations to shift the classifier learning bias towards minority class.

To generate artificial data, it uses bootstrapping and k-nearest neighbors. It basically takes the difference between feature vectors and its nearest neighbors and multiply the difference by random number in between 0 and 1 and add it to feature vector under consideration. This results in the selection of random point along the line segment between two specific features.

Usage:

```
train$Class = as.factor(train$Class)  
train = SMOTE(Class ~ ., train, perc.over = 300, perc.under=100) #generating new balance data by ROSE boosting  
table(train$Class) #check imbalance of new data
```

**UNDERSAMPLING-** This method reduces the number of observations from those class which are in majority in the data set so to balance the cases. This method mostly works better when the data is huge in size and reducing the number of training samples helps to reduce the modelling run time and storage troubles.

There are two types of Undersampling methods:

- Random
- Informative

Usage:

```
under_sampled_data = ovun.sample(Class ~ ., data = train, method = "under")$data # to reduce the number of  
observations from majority class data set.  
table(under_sampled_data$Class) # check the size of table
```

**OVERSAMPLING-** In a dataset, the class having the lowest number of observation is taken into account and all existing observations are taken and copied and extra observations are added by randomly sampling with replacement from this class. It is just opposite to undersampling.

Usage:

```
over_sampled_data = ovun.sample(Class~ ., data =train, method = "over")$data # to increase the number of  
observations of minority class data set.  
table(over_sampled_data$Class) # check the size of table
```

### 3. FEATURE SELECTION AND MODEL DEVELOPING

**C5.0-** C5.0 algorithm is a successor of C4.5 algorithm also developed by Quinlan (1994). To understand the working of C5.0, we must know how C4.5, a descendant of CLS and ID3, generates classifiers that can be used for the classification problems. It extends the ID3 algorithm by dealing with both continuous and discrete attributes, missing values and pruning trees after construction.

The algorithm for the C5.0 is presented as follows: -

**Input:** - A set of set S (continuous or discrete attributes) each belonging to one class.

**Output:** - A decision tree or a set of rules that assigns a class to a new case.

Algorithm: -

1. Check for the base case
2. Find the attribute with the highest information gain.
3. Partition S into S1, S2, S3.... according to the values of the attribute selected in Step 2
4. Repeat the steps for S1, S2, S3....

The base cases are the following

- All the examples from the training set belong to the same class (a tree leaf labelled with that class is returned).
- The training set is empty (returns a tree leaf called failure).
- The attribute list is empty (return a leaf labelled with the most frequent class or the disjunction of all the classes)

The attribute with the highest information gain is computed using the following formulas

Entropy:  $E(S) = \sum_{i=1}^n -Pr(C_i) * \log_2 Pr(C_i)$

Informational Gain:  $G(S,A) = E(S) - \sum_{i=1}^m Pr(A_i) E(S_i)$

$E(S) \Rightarrow$  Information entropy of S

$G(S,A) \Rightarrow$  gain of S after a split on attribute A

$N \Rightarrow$  Number of classes in S

$Pr(C_i) \Rightarrow$  frequency of class  $C_i$  in S

$m \Rightarrow$  Number of values of attribute A in S

$Pr(A_i) \Rightarrow$  frequency of classes that have  $A_i$  value in S

$E(S_i) \Rightarrow$  Subset of S with items that have  $A_i$  value

The C4.5 algorithm improves the ID3 algorithm by allowing numerical attributes, permitting missing values and performing Tree pruning.

C4.5 was superseded by See5/C5.0 (or C5.0 for short). The changes encompass new capabilities as well as much-improved efficiency, and include:

1. A variant of boosting, which constructs an ensemble of classifiers that are then voted to give a final classification. Boosting often leads to a dramatic improvement in predictive accuracy.
2. New data types (e.g., dates), “not applicable” values, variable misclassification costs, and mechanisms to pre-filter attributes.
3. Unordered rulesets—when a case is classified, all applicable rules are found and voted. This improves both the interpretability of rulesets and their predictive accuracy
4. Greatly improved scalability of both decision trees and (particularly) rulesets. Scalability is enhanced by multi-threading; C5.0 can take advantage of computers with multiple CPUs and/or cores.

Usage:- `modell=C5.0(data1,data2)`  
`modell`

Code: - **C5.0(X,Y)** where X is the predicting features while the Y is the Response.

The entries of the Y should be the factor to run the code, Hence **as.factor** is used to covert the entries into factor.

# Implementation Details:

"Data Set Input and initial analysis the dataset we started working on was obtained from the UCI repository. It initially had 1567 observations (rows) which were outputs of 590 sensor measurements (columns/variables) and a label of Yield Pass/Fail."

#Loading data to workspace

#Step 1: Removing Redundant data and Missing Values

#Step 2: Imputation using different techniques

#Step 3: Splitting into training & testing data

#Step 4: Feature Selection techniques

#Step 5: Case Boosting

#Step 6: MODEL BUILDING

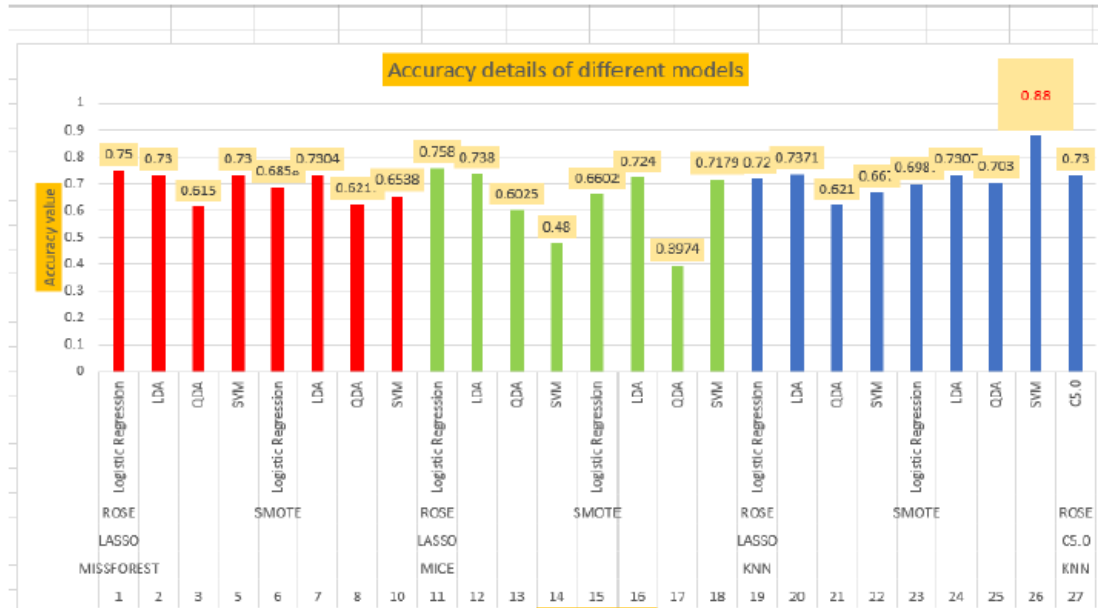
#SVM

```
library(e1071)
svm.fit=svm(Class~., data=train_rose[,c("Class",lasso_op)],
kernel="polynomial", cost=10,scale=FALSE)
svm.pred=predict(svm.fit,test)
summary(svm.fit)
## Call:
## svm(formula = Class ~ ., data = train_rose[, c("Class", lasso_op)],
##      kernel = "polynomial", cost = 10, scale = FALSE)
## Parameters:
##      SVM-Type:      C-classification
##      SVM-Kernel:    polynomial
##      cost:          10
##      degree:        3
##      gamma:         0.00625
##      coef.0:        0
## Number of Support Vectors: 604
## ( 308 296 )
## Number of Classes: 2
## Levels:
## pass fail
table(svm.pred,test$Class)
##      svm.pred      pass      fail
##      pass         98         5
##      fail         47         6
mean(svm.pred==test$Class)
## [1] 0.8846154 #Accuracy
```

#Hence, we finally selected SVM because of highest accuracy rate:

## Comparison between different Models:

Sl.No	Imputation	Feature Selection	Case Boosting	Model	Prediction Accuracy
1	MISSFOREST	LASSO	ROSE	Logistic Regression	0.75
2				LDA	0.73
3				QDA	0.615
5				SVM	0.73
6			SMOTE	Logistic Regression	0.6858
7				LDA	0.7304
8				QDA	0.6217
10				SVM	0.6538
11	MICE	LASSO	ROSE	Logistic Regression	0.758
12				LDA	0.738
13				QDA	0.6025
14				SVM	0.48
15			SMOTE	Logistic Regression	0.6602
16				LDA	0.724
17				QDA	0.3974
18				SVM	0.7179
19	KNN	LASSO	ROSE	Logistic Regression	0.72
20				LDA	0.7371
21				QDA	0.621
22				SVM	0.667
23			SMOTE	Logistic Regression	0.6987
24				LDA	0.7307
25				QDA	0.703
26				SVM	<b>0.88</b>
27	KNN	C5.0	ROSE	C5.0	0.73



We finally selected the SVM model using KNN imputation for our Predictions which uses LASSO for feature selection and ROSE for Case Boosting because out of different models' accuracy results, we found that this model has the highest accuracy of 88%. The data was first imputed using KNN, MICE and MISSFOREST and for each imputation we then did the feature selection using LASSO, RandomForest, PCA and then case boosting was done using ROSE and SMOTE. Using these all permutations and combinations we got different datasets for individual process path and then model was made using Logistic Regression, LDA, QDA, SVM and C5.0 to find out the best model with highest accuracy.

# Executive Summary:

This project deals with making predictive models for equipment fault detection in semiconductor manufacturing process. Semiconductor manufacturing is one of the most technologically complicated process and involves many variables that are monitored during manufacturing. In past many Machine Learning algorithms like multivariate analysis have been deployed for creating predictive models to detect faults. As per the semiconductor industry statistics the equipment usually suffers an 8% unscheduled downtime. If we can make a predictive model for preventive maintenance to reduce this 8% unscheduled downtime, it will improve the productivity significantly. Predictive Maintenance is the process of identifying when equipment needs maintenance to make sure that failure of equipment is avoided.

Our Objective is to make a model that classifies manufactured semiconductor wafer as yield pass/fail. We are training this model on a data set that consists of 590 sensor attributes and 1567 observations. These values are recorded by different sensors during the manufacturing process. So from this semiconductor data set analysis, we can appl this end-to-end solution to a Smart Factory Operation.

The first problem we faced was the missing values in the dataset. We had to first clean the data set which can be identified as the data cleaning phase. Columns with more than 40% missing values and columns consisting same data in all observations were removed. After initial cleaning of data, we learned new imputation techniques like Mice, KNN, missForest and applied them to fill the missing values. Another issue we faced while training the model was that the ratio of Yield Pass and Fails was 14:1 [1463 Pass and only 104 fails]. Training the model on this data would create a biased model. To encounter this issue, we used Case Boosting techniques such as ROSE and Smote that would equalize the pass/fail cases reducing the bias.

Monitoring all the variables/sensors is not a practical solution to decide if a wafer is pass/fail. After boosting we then do feature selecting using different models for selecting the important sensors for designing model. Out of all the models we got best result by doing KNN imputation, Lasso feature selection followed by ROSE Case boosting and finally modeling using SVM. This model has a prediction accuracy of 88%. We can use this model for preventive maintenance and predict with an accuracy of 88% to determine if a semiconductor wafer will pass or fail. This will reduce the 8% unscheduled downtime considerably.



# Conclusion:

Semiconductor manufacturing is an overly complex process, and the industry is one of the most capital-intensive industry. Process excursion during manufacturing can significantly impact product yield and manufacturing cost as well. So, for such processes, real time data is available, and our job is to build a predictive classification model since Detection of root cause gives process Engineers action plan to improve manufacturing robustness and prevent future process excursions. With such a large number of features in the manufacturing process, the most important step is the causal feature selection for effective monitoring & process control as this causal feature selection informs process engineers that which of features (sensor signals) are most important for the particular process so that preventive/predictive maintenance of the related machinery/equipment is carried out well before time in order to reduce any downtime that is supposed to occur because of that breakdown.

Therefore, we synthesized the classification model building techniques along with the feature selection methods to ensure that only the most important features (excluding the noises) are used for building classification model which will be used for prediction of pass/fail class to the semiconductor wafer manufactured. After data cleaning and imputation, we have tried various feature selection methods such as Principal Component analysis, Lasso, C5.0 and Random Forest. The best results obtained from Lasso was then used to build a classification model using various techniques such as Logistic Regression, LDA, QDA, KNN, Decision tree & SVM. After comparing based on confusion matrix/accuracy rate, AUC & Classification error rate and cross validating all the models, SVM was found to be the best among the above mentioned.

Hence, the SVM classification model was deployed and the prediction model was built. This would help the process Engineers in the semiconductor manufacturing industry to improve the yield by preventing breakdown thus reducing downtime and improving the business model. Our goal is to apply this end-to-end solution package to Smart Factory Operation (Industrial Revolution 4.0) with integrated sensor set to collect data then process the data set to achieve our goal and meet customer's requirements. Fell free to contact us if you have any requirement to implement this solution.

# References:

1. Author: Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani

Title: An Introduction to Statistical Learning

Edition: 7th edition

URL: <http://www-bcf.usc.edu/~gareth/ISL/ISLR%20First%20Printing.pdf>

2. Online site URL: <https://www.analyticsvidhya.com/blog/tag/statistics/>

3. Author: Sathyan Munirathinam and Balakrishnan Ramadoss

Journal: IACSIT International Journal of Engineering and Technology, Vol. 8, No.

4, August 2016

Title: Predictive Models for Equipment Fault Detection in the Semiconductor Manufacturing Process

URL: <http://www.ijetch.org/vol8/898-T10023.pdf>

4. Author: P.S. Frankwicz, S. E. Romano and T. Moutinho

Title: Process Excursion Detection using Statistical Analysis Methodologies in High Volume Semiconductor Production

URL: <http://www.lexjansen.com/nesug/nesug09/po/PO11.pdf>

5. Author: Kittisak Kerdprasop and Nittaya Kerdprasop

Journal: INTERNATIONAL JOURNAL OF MECHANICS

Title: A Data Mining Approach to Automate Fault Detection Model Development in the Semiconductor Manufacturing Process

URL: <http://www.naun.org/main/NAUN/mechanics/17-220.pdf>

6. Author: Michael McCann, Yuhua Li, Liam Maguire, Adrian Johnston Journal: Workshop and Conference Proceedings 6: 277-288 Title: Causality Challenge: Benchmarking relevant signal components for effective monitoring and process control

URL: <http://proceedings.mlr.press/v6/mccann10a/mccann10a.pdf>

7. Online site URL: <https://stackoverflow.com/questions/13114812/imputation-in-r>

8. KNN Imputation: <https://towardsdatascience.com/the-use-of-knn-for-missing-values-cf33d935c637>

9. R-programming

Online site URL: <https://www.datacamp.com/>

10. <http://archive.ics.uci.edu/ml/datasets/SECOM>

11. Reference:- <https://www.r-bloggers.com/imputing-missing-data-with-r-mice-package/>

12. Reference:- <https://cran.r-project.org/web/packages/Boruta/Boruta.pdf>