

Configuration Manual

MSc Research Project
Data Analytics

Omkar Doke
Student ID: x18179525

School of Computing
National College of Ireland

Supervisor: Dr. Muhammad Iqbal

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Omkar Doke
Student ID:	x18179525
Programme:	Data Analytics
Year:	2020
Module:	MSc Research Project
Supervisor:	Dr. Muhammad Iqbal
Submission Due Date:	17/08/2020
Project Title:	Configuration Manual
Word Count:	1744
Page Count:	19

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	10th August 2020

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Omkar Doke
x18179525

1 Introduction

This configuration manual presents the software and hardware requirements along with the details of programming codes written for model implementation in research project: “Data Mining for Enhancing Silicon Wafer Fabrication”

2 System Configuration

2.1 Hardware Specifications

Table 1 represents hardware specification of the system on which the research was accomplished.

Table 1: Hardware Specification

RAM	8 GB
Processor	Intel i7 8550U
Speed	1.99 GHz
Operating System	Windows 10, 64 Bit
Storage	1 TB HDD
GPU	NVIDIA GeForce MX150

2.2 Software Specifications

- **Microsoft Excel 2019:**

Both the data-sets used in the research were downloaded and stored in csv (comma separated values) in excel. It was used for quick evaluation and exploratory plot.

- **Jupyter Notebook from Anaconda Distribution:**

Anaconda Navigator is an open source software downloaded from the anaconda distribution website ¹. It supports jupyter notebooks to implement machine learning models on research data. Latest version of jupyter notebook (version 5.7.4) was used in the research for data preprocessing, exploratory data analysis (EDA), manipulation of data, transformation and implementation of models.

¹<https://www.anaconda.com/products/individual>

3 Development of Project

Python programming was used to accomplish the research in various phases viz. data pre-processing of both the data-sets, EDA, merging of both the data-sets, addressing class imbalance and normalization of data to overcome the impact of outliers. It was followed by splitting data into train and test set for predictive modelling using classification-based machine learning algorithms and their cross validation using stratified K-fold validation technique. Sk-Learn (scikit-learn) and Keras were primary libraries used along with numpy, panda, matplotlib for executing the code.

3.1 Data Preparation

Both data-sets² downloaded from different websites³ have been uploaded onto jupyter notebook in csv format. Following sections provide a detail insight of data-processing, EDA, feature engineering, dimensionality reduction performed on both data-sets followed by merging of data-sets for implementation and evaluation of models after addressing class imbalance.

3.1.1 UCI SEMCOM Dataset

Pre-processing of UCI SEMCOM dataset involves handling missing values. UCI SEMCOM dataset consists of 591 attributes with 27 attributes having more than 50% of missing values which were dropped as it didn't lead to data loss. Apart from that, attributes with zero variance (i.e. no effect of dependent variable) were dropped as their presence or absence didn't have any impact on research. Attributes with less than 50% of missing values were imputed with median as the attributes had outliers and data has skew symmetric distribution. Thereafter, dataset was normalized using MinMaxScaler library for scaling because attributes consisted of outliers as well as the attribute values were in different range. The dependent variable of UCI SEMCOM dataset consists of pass category defined as '-1' and fail category as '+1'. Code for preprocessing of UCI SEMCOM dataset is highlighted in Figure 1.

3.1.2 WAFER Dataset

Pre-processing of WAFER dataset involves handling missing values. WAFER dataset consists of 154 attributes with no attributes having more than 50% of missing values thereby none of the attributes were dropped. Also, when checked for impact of attributes on dependent variable, it was found that none of the attributes had zero variance. Attributes with less than 50% of missing values were imputed with mean as the attributes didn't have outliers. Thereafter, dataset was normalized using MinMaxScaler library for scaling. The dependent variable of WAFER dataset consisted of pass category defined as '+1' and fail category as '-1'. To have standardized definition of pass and fail classes in dependent variable, we interchanged the designation for WAFER dataset thereby assigning '-1' to pass class and '+1' to fail class. Code for preprocessing of WAFER dataset is highlighted in Figure 2.

²<http://www.timeseriesclassification.com/description.php?Dataset=Wafer>

³<https://archive.ics.uci.edu/ml/datasets/SECOM>

Data Cleaning

```
1 ## Dropping 1st row
2 Wafer_fabrication_df = Wafer_fabrication_df.drop(Wafer_fabrication_df.index[0])
3
4 ##Removing Columns with more than 50% NaN Values
5 cols = Wafer_fabrication_df.columns[Wafer_fabrication_df.isnull().mean()>0.5]
6 Wafer_df = Wafer_fabrication_df.drop(cols, axis=1)
7
8 ## Dropping 1 column with date as time
9 Wafer_df = Wafer_df.drop(columns= [0], axis = 1)
10
11 ##Dropping columns with 0 varainace
12 Wafer_df = Wafer_df.loc[:,Wafer_df.apply(pd.Series.nunique) != 1]
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

Figure 1: Data Cleaning of UCI SEMCOM Dataset

```
1
2
3 Wafer_fabrication2 = pd.read_csv('D:/NCI - Research Project/Data/Wafer/csv_result-Wafer.csv', header = None)
4 Wafer_fabrication2_df = pd.DataFrame(Wafer_fabrication2)
5
6 ## Dropping 1st row
7 Wafer_fabrication2_df = Wafer_fabrication2_df.drop(Wafer_fabrication2_df.index[0])
8
9 ##Removing Columns with more than 50% NaN Values
10 cols = Wafer_fabrication2_df.columns[Wafer_fabrication2_df.isnull().mean()>0.5]
11 Wafer2_df = Wafer_fabrication2_df.drop(cols, axis=1)
12
13 ## Dropping 1st column with sequece number
14 Wafer2_df = Wafer2_df.drop(columns= [0], axis = 1)
15
16 ##Dropping columns with 0 varainace
17 Wafer2_df = Wafer2_df.loc[:,Wafer2_df.apply(pd.Series.nunique) != 1]
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

Figure 2: Data Cleaning of WAFER Dataset

3.2 Dimensionality Reduction on both data-sets

3.2.1 UCI SEMCOM Dataset

Feature extraction was performed using principle component analysis (PCA) technique to extract top components explaining 80% variance of the data. PCA was applied to extract 250 components from 447 attributes. Then after, variance ratio was calculated and plotted for principal components which led to the selection of top 100 components as they explained more that 80% variance of data. Figure 3 represents the code for implementation of PCA on UCI SEMCOM data for extracting principle components.

Applying PCA

```
1  ## Applying PCA to get top Principle components representing maximum variance in data
2  pca = PCA(n_components=250)
3  principalComponents = pca.fit_transform(x)
4  columns = ['pca_%i' % i for i in range(250)]
5  principalDf = pd.DataFrame(data = principalComponents, columns = columns)
6
7  ## Identifying number of PCA's that explain maximum variance of data (we are attempting to find for 98-99% or more)
8  pca.fit(x)
9  variance = pca.explained_variance_ratio_ #calculate variance ratios
10 var=np.cumsum(np.round(pca.explained_variance_ratio_, decimals=3)*100)
11 var #cumulative sum of variance explained with [n] features
12
13 ##Plotting PCA's against variance to identify the PCA's to be selected
14 fig= plt.figure(figsize=(8,5))
15 plt.ylabel('% Variance Explained')
16 plt.xlabel("Number of PCA's")
17 plt.title('PCA Analysis')
18 plt.ylim(5,100.5)
19 plt.style.context('seaborn-blackgrid')
20 plt.grid(True)
21 plt.plot(var)

1  ##Thus from th previous plot, selecting first 100 PCA for this research as they explain more than 80% variance of data
2  principalDf.drop(principalDf.iloc[:,100:250], axis = 1, inplace = True)
3
4  ##Concatinating Dependent Variable to the dataframe
5  FinalDf = pd.concat([principalDf, y], axis = 1)
6
7  ##Renumbering last column count
8  FinalDf.rename(columns={591: "Pass/Fail"}, inplace= True)
9
10 ##Shifting last column up by 1 row
11 FinalDf['Pass/Fail'] = FinalDf['Pass/Fail'].shift(-1)
12
13 ##Dropping the last row
14 FinalDf = FinalDf[:-1]
```

Figure 3: Feature Extraction using PCA on UCI SEMCOM Data

Feature selection was performed using Analysis of Variance (ANOVA) technique to select top features explaining more that 80% variance in data. Initially number of features were gradually reduced to identify feature count for which models provide optimum performance, however different models provided optimum performance for different feature count. Then after, top 100 features were selected from the data to compare model's performance with that of feature extraction technique. Code for feature selection using ANOVA is highlighted in Figure 4.

Applying ANOVA

```
1  ##Select Features With Best ANOVA F-Values
2  ## Create an SelectKBest object to select features with two best ANOVA F-Values
3  fvalue_selector = SelectKBest(f_classif, k=100)
4
5  ## Apply the SelectKBest object to the features and target
6  X_kbest = fvalue_selector.fit_transform(x, y)
7  FS_1Df = pd.DataFrame(X_kbest)
8
9  ##Concatinating Dependent Variable to the dataframe
10 FS_1Df = pd.concat([FS_1Df, y], axis = 1)
11
12 ##Renumbering last column count
13 FS_1Df.rename(columns={591: "Pass/Fail"}, inplace= True)
14
15 ##Shifting last column up by 1 row
16 FS_1Df['Pass/Fail'] = FS_1Df['Pass/Fail'].shift(-1)
17
18 ##Dropping the last row
19 FS_1Df = FS_1Df[:-1]
20 FS_1Df.head()
```

Figure 4: Feature Selection using ANOVA on UCI SEMCOM Data

3.2.2 WAFER Dataset

Feature selection was performed using Analysis of Variance (ANOVA) technique to select top features explaining more than 80% variance in data. Initially number of features were gradually reduced to identify feature count for which models provide optimum performance, however different models provided optimum performance for different feature count. Then after, top 100 features were selected from the data to compare model's performance with that of feature extraction technique. Code for feature selection using ANOVA is highlighted in Figure 5.

Applying ANOVA

```
1  ##Select Features With Best ANOVA F-Values
2  ## Create an SelectKBest object to select features with two best ANOVA F-Values
3  fvalue_selector = SelectKBest(f_classif, k=100)
4
5  ## Apply the SelectKBest object to the features and target
6  X2_kbest = fvalue_selector.fit_transform(x2, y2)
7  FS_2Df = pd.DataFrame(X2_kbest)
8
9  ##Concatinating Dependent Variable to the dataframe
10 FS_2Df = pd.concat([FS_2Df, y2], axis = 1)
11
12 ##Renumbering last column count
13 FS_2Df.rename(columns={153: "Pass/Fail"}, inplace= True)
14
15 ##Shifting last column up by 1 row
16 FS_2Df['Pass/Fail'] = FS_2Df['Pass/Fail'].shift(-1)
17
18 ##Dropping the last row
19 FS_2Df = FS_2Df[:-1]
20 FS_2Df.head()
```

Figure 5: Feature Selection using ANOVA on WAFER Data

Feature extraction was performed using principle component analysis (PCA) technique to

extract top components explaining 80% variance of the data. PCA was applied to extract 150 components from attributes. Then after, variance ratio was calculated and plotted for principal components which led to the selection of top 100 components as they explained more than 80% variance of data. Figure 6 represents the code for implementation of PCA on UCI SEMCOM data for extracting principle components.

Applying PCA

```

1  ## Applying PCA to get top Principle components representing maximum variance in data
2  pca2 = PCA(n_components=150)
3  principalComponents2 = pca2.fit_transform(x2)
4  columns = ['pca_%i' % i for i in range(150)]
5  principal2_Df = pd.DataFrame(data = principalComponents2, columns = columns)
6
7  ## Identifying number of PCA's that explain maximum variance of data (we are attempting to find for 98-99% or more)
8  pca2.fit(x2)
9  variance2 = pca2.explained_variance_ratio_ #calculate variance ratios
10 var2=np.cumsum(np.round(pca2.explained_variance_ratio_, decimals=3)*100)
11 var2 #cumulative sum of variance explained with [n] features
12
13 ##Plotting PCA's against variance to identify the PCA's to be selected
14 fig= plt.figure(figsize=(8,5))
15 plt.ylabel('% Variance Explained')
16 plt.xlabel("Number of PCA's")
17 plt.title('PCA Analysis')
18 plt.ylim(5,100.5)
19 plt.style.context('seaborn-blackgrid')
20 plt.grid(True)
21 plt.plot(var2)

```

```

1  ##Thus from th previous plot, selecting first 100 PCA for this research
2  principal2_Df.drop(principal2_Df.iloc[:,100:150], axis = 1, inplace = True)
3
4  ##Concatinating Dependent Variable to the dataframe
5  Final2_Df = pd.concat([principal2_Df, y2], axis = 1)
6
7  ##Renumbering Last column count
8  Final2_Df.rename(columns={153: "Pass/Fail"}, inplace= True)
9
10 ##Shifting Last column up by 1 row
11 Final2_Df['Pass/Fail'] = Final2_Df['Pass/Fail'].shift(-1)
12
13 ##Dropping the last row
14 Final2_Df = Final2_Df[:-1]

```

Figure 6: Feature Extraction using PCA on WAFER Data

3.3 Merging of data

3.3.1 Merging of feature extracted data frames

Two data frames are created of principle components extracted from both data-sets which are then merged as both had same number of columns. The pass category which was initially assigned '-1' label was reassigned with label '0' and descriptive analysis was performed on final merged dataset. Code in Figure 7 shows how feature extracted PCA data frames from both data-sets were merged together.

3.3.2 Merging of feature selected data frames

Two data frames are created of feature's selected from both data-sets which are then merged as both had same number of columns. The pass category which was initially assigned '-1' label was reassigned with label '0'. Code in Figure 8 shows how feature selected data frames from both data-sets were merged together.

Merging Two PCA Data Frames

```
1 ##Merging data frames
2 Wafer = pd.concat([FinalDf, Final2_Df], ignore_index=True)
3
4 ## Changing -1 to 0
5 Wafer['Pass/Fail'].replace({-1.0: 0.0}, inplace=True)
6
7 ##Statistical Values of Each Column
8 Wafer_des = Wafer.describe()
9 Wafer_des
```

Figure 7: Merging two PCA Data Frames

Merging Two Feature Selected Data Frames

```
1 ##Merging data frames
2 FS_Wafer = pd.concat([FS_1Df, FS_2Df], ignore_index=True)
3
4 ## Changing -1 to 0
5 FS_Wafer['Pass/Fail'].replace({-1.0: 0.0}, inplace=True)
```

Figure 8: Merging two Feature Selected Data Frames

3.4 Splitting the data into Train and Test set

After merging, both the data-set were split into train and test part in 75:25 ratio respectively. Models were trained on train set and evaluated on test set. Their performance was cross validated using stratified K-fold validation technique. Figure 9 illustrates the code for train test split of final dataset.

```
1 ##Again separating the dependent and independent variables from FinalDf
2 x = Wafer.iloc[:,100]
3 y = Wafer.iloc[:, 100]
4
5 #Getting the shapes of new data sets x and y
6 print("Shape of x:", x.shape)
7 print("Shape of y:", y.shape)
8
9 ##Splitting the data into train and test sets
10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.25, random_state = 0)
11
12 # getting the shapes
13 print("Shape of x_train: ", x_train.shape)
14 print("Shape of x_test: ", x_test.shape)
15 print("Shape of y_train: ", y_train.shape)
16 print("Shape of y_test: ", y_test.shape)
```

```
Shape of x: (8731, 100)
Shape of y: (8731,)
Shape of x_train: (6548, 100)
Shape of x_test: (2183, 100)
Shape of y_train: (6548,)
Shape of y_test: (2183,)
```

Figure 9: Train Test Split of Merged Data

3.5 Addressing Class Imbalance

3.5.1 Oversampling of Feature Selected and Feature Extracted Data

After splitting the data, major class imbalance was observed in train set with fail class contributing approximately 10% of entire data. This was then addressed using Synthetic Minority Over-Sampling Technique (SMOTE) wherein the minority class was over-sampled to 50% to that of majority class in both features extracted and features selected data. Figure 10 represents the code for oversampling of minority class using SMOTE.

Oversampling using SMOTE of the Fail Cases

```
1 SM = SMOTE(sampling_strategy= 0.50, random_state= None)
2 x_train_os, y_train_os = SM.fit_sample(x_train, y_train)
3
4 print("Shape of x_train_os: ", x_train_os.shape)
5 print("Shape of y_train_os: ", y_train_os.shape)
```

Figure 10: Oversampling of Train Set using SMOTE

3.5.2 Random Sampling of Feature Selected Data

In another experiment, class imbalance of feature selected data was address by random oversampling of minority class along with random under-sampling of majority class. 3 different rations of oversampling and under-sampling respectively were experimented viz. 40:60, 45:55 and 50:50. Figure 11 represents the code for sampling of majority and minority class using random sampling.

40:60 Sampling Ratio

```
1 ## Oversample
2 oversample = RandomOverSampler(sampling_strategy=0.40)
3 x_train_ros_1, y_train_ros_1 = oversample.fit_sample(x_train, y_train)
4
5 ## Undersample
6 undersample = RandomUnderSampler(sampling_strategy=0.60)
7 x_train_bd, y_train_bd= undersample.fit_sample(x_train_ros_1, y_train_ros_1)
```

45:55 Sampling Ratio

```
1 ## Oversample
2 oversample = RandomOverSampler(sampling_strategy=0.45)
3 x_train_ros_1, y_train_ros_1 = oversample.fit_sample(x_train, y_train)
4
5 ## Undersample
6 undersample = RandomUnderSampler(sampling_strategy=0.55)
7 x_train_bd, y_train_bd= undersample.fit_sample(x_train_ros_1, y_train_ros_1)
```

50:50 Sampling Ratio

```
1 ## Oversample
2 oversample = RandomOverSampler(sampling_strategy=0.50)
3 x_train_ros_1, y_train_ros_1 = oversample.fit_sample(x_train, y_train)
4
5 ## Undersample
6 undersample = RandomUnderSampler(sampling_strategy=0.50)
7 x_train_bd, y_train_bd= undersample.fit_sample(x_train_ros_1, y_train_ros_1)
```

Figure 11: Random Sampling of Train Set in Feature Selected Data

3.6 Model Implementation and Cross Validation

Various classification models viz. Decision Tree, Logistic Regression, XGBoost, Random Forest, SVM-Linear, SVM-RBF, Naïve Bayes, KNN and basic Neural Network were implemented on pre-processed and feature engineered data. Their performance was evaluated for precision and accuracy. The accuracy of each model was further cross validated using stratified K-fold validation.

3.6.1 Decision Tree

Braha and Shmilovici (2002) used Decision Tree (DT) in their research and achieved an accuracy of 77%. DT was implemented using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for DT is illustrated in Figure 12 whereas Figure 13 represents the K-Fold validation of DT.

Decision Tree and its Confusion Matrix

```
1  ## Decision Tree
2  DT = DecisionTreeClassifier()
3
4  ## Training DT
5  DT = DT.fit(x_train_fos, y_train_fos)
6
7  ## Predicting response on Test
8  y_pred = DT.predict(x_test)
9
10 DT_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 DT_R2 = metrics.recall_score(y_test, y_pred)*100
12 DT_P2 = metrics.precision_score(y_test, y_pred)*100
13 DT_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print('Accuracy: %.2f%% ' % DT_A2)
16 print("Recall_Accuracy: %.2f%%" % DT_R2)
17 print("Precision_Accuracy: %.2f%%" % DT_P2)
18 print("F1 Score: %.2f%%" % DT_F2 )
```

```
Accuracy: 96.98%
Recall_Accuracy: 87.50%
Precision_Accuracy: 83.76%
F1 Score: 85.59%
```

```
1  cm = confusion_matrix(y_test, y_pred)
2  TP_DT_2 = cm[1][1]
3
4  plt.rcParams['figure.figsize'] = (5, 5)
5  #sns.set(style = 'dark', font_scale = 1.4)
6  sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8  print(confusion_matrix(y_test, y_pred))
```

```
[[1921  38]
 [ 28 196]]
```

Figure 12: Code for Decision Tree

Decision Tree with Stratified K-Fold

```
1  ## Decision Tree for k = 10
2  skfold = StratifiedKFold(n_splits = 10,random_state=None)
3  model_skfold = DecisionTreeClassifier()
4  results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5  DT_1 = (results_skfold.mean()*100)
6  print("Accuracy when k is 10 : %.2f%%" % DT_1)
7
8  ## Decision Tree for k = 20
9  skfold = StratifiedKFold(n_splits = 20,random_state=None)
10 model_skfold = DecisionTreeClassifier()
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 DT_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % DT_2)
14
15 ## Decision Tree for K = 30
16 skfold = StratifiedKFold(n_splits = 30,random_state=None)
17 model_skfold = DecisionTreeClassifier()
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 DT_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % DT_3)
21
22 ## Decision Tree for K = 50
23 skfold = StratifiedKFold(n_splits = 50,random_state=None)
24 model_skfold = DecisionTreeClassifier()
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 DT_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % DT_4)
```

Accuracy when k is 10 : 92.53%
Accuracy when k is 20 : 95.41%
Accuracy when k is 30 : 95.55%
Accuracy when k is 50 : 96.47%

Figure 13: K-Fold Validation of Decision Tree

3.6.2 Logistic Regression

Logistic Regression was implemented using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for Logistic Regression is illustrated in Figure 14 whereas Figure 15 represents its K-Fold validation.

Logistic Regression & its Confusion Matrix

```
1  ## L0gistic Regression
2  logreg = LogisticRegression(random_state= 0)
3
4  ## Training Model
5  LogReg = logreg.fit(x_train_fos, y_train_fos)
6
7  ## Predicting response
8  y_predLog = LogReg.predict(x_test)
9
10 LogReg_A2 = metrics.accuracy_score(y_test, y_predLog)*100
11 LogReg_R2 = metrics.recall_score(y_test, y_predLog)*100
12 LogReg_P2 = metrics.precision_score(y_test, y_predLog)*100
13 LogReg_F2 = metrics.f1_score(y_test, y_predLog)*100
14
15 print("Accuracy: %.2f%%" % LogReg_A2)
16 print("Recall_Accuracy: %.2f%%" % LogReg_R2)
17 print("Precision_Accuracy: %.2f%%" % LogReg_P2)
18 print("F1 Score: %.2f%%" % LogReg_F2)
```

```
1  cm = confusion_matrix(y_test, y_predLog)
2  TP_LogReg_2 = cm[1][1]
3
4  plt.rcParams['figure.figsize'] = (5, 5)
5  #sns.set(style = 'dark', font_scale = 1.4)
6  sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8  print(confusion_matrix(y_test, y_predLog))
```

Figure 14: Code for Logistic Regression

Logistic Regression with Stratified K-Fold

```
1  ## Logistic Regression for k = 10
2  skfold = StratifiedKFold(n_splits = 10, random_state=None)
3  model_skfold = LogisticRegression()
4  results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5  LogReg_1 = (results_skfold.mean()*100)
6  print("Accuracy when k is 10 : %.2f%%" % LogReg_1)
7
8  ## Logistic Regression for k = 20
9  skfold = StratifiedKFold(n_splits = 20, random_state=None)
10 model_skfold = LogisticRegression()
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 LogReg_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % LogReg_2)
14
15 ## Logistic Regression for K = 30
16 skfold = StratifiedKFold(n_splits = 30, random_state=None)
17 model_skfold = LogisticRegression()
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 LogReg_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % LogReg_3)
21
22 ## Logistic Regression for K = 50
23 skfold = StratifiedKFold(n_splits = 50, random_state=None)
24 model_skfold = LogisticRegression()
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 LogReg_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % LogReg_4)

```

C:\Users\Omkar Doke\Anaconda3\lib\site-packages\sklearn\linear_model_logistic.py:940: ConvergenceWarning: lbfgs failed to converge (status=1):

Figure 15: K-Fold Validation of Logistic Regression

Model's precision failure was further studied for its threshold for classification of its probabilities. It was then adjusted after plotting the histogram plot and the model was re-implemented which saw further reduction in precision. Figure 16 shows code for experiment with logistic regression.

Understanding why True Positive is less and how to adjust the threshold of classification for that

```
1  # print the first 25 true and predicted responses
2  print('True:', y_test.values[0:100])
3  print('Predicted:', y_predLog[0:100])

```

```
1  # print the first 10 predicted probabilities of class membership
2  LogReg.predict_proba(x_test)[79:100]
3
4  # print the first 10 predicted probabilities for class 1
5  LogReg.predict_proba(x_test)[0:10, 1]
6
7  # store the predicted probabilities for class 1
8  y_pred_prob = LogReg.predict_proba(x_test)[: , 1]
9  y_pred_prob_df = pd.DataFrame(y_pred_prob)

```

```
1  # histogram of predicted probabilities
2  # 8 bins
3  plt.hist(y_pred_prob, bins=5)
4
5  # x-axis limit from 0 to 1
6  plt.xlim(0,1)
7  plt.title('Histogram of predicted probabilities')
8  plt.xlabel('Predicted probability of wafer test')
9  plt.ylabel('Frequency')

```

```
1  # predict wafer fail if the predicted probability is greater than 0.4, it will return 1 for all values above 0.2 and 0 other
2  # results are 2D so we slice out the first column
3
4  y_pred_class = pd.DataFrame(binimize(y_pred_prob_df, 0.2))
5
6  LogReg_B_A1 = metrics.accuracy_score(y_test, y_pred_class)*100
7  LogReg_B_R1 = metrics.recall_score(y_test, y_pred_class)*100
8  LogReg_B_P1 = metrics.precision_score(y_test, y_pred_class)*100
9  LogReg_B_F1 = metrics.f1_score(y_test, y_pred_class)*100
10
11 print("Accuracy: %.2f%%" % LogReg_B_A1)
12 print("Recall_Accuracy: %.2f%%" % LogReg_B_R1)
13 print("Precision_Accuracy: %.2f%%" % LogReg_B_P1)
14 print("F1 Score: %.2f%%" % LogReg_B_F1)

```

```
1  m = confusion_matrix(y_test, y_pred_class)
2  TP_LogReg_B = cm[1][1]
3
4  plt.rcParams['figure.figsize'] = (5, 5)
5  sns.set(style = 'dark', font_scale = 1.4)
6  sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8  print(confusion_matrix(y_test, y_pred_class))

```

Figure 16: Understanding Poor Performance of Logistic Regression

3.6.3 XGBoost

XGBoost was implemented using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for XGBoost is illustrated in Figure 17 whereas Figure 18 represents its K-Fold validation.

XGB Classifier and its CM

```
1  ## XGB Boost
2  XGB = XGBClassifier()
3
4  ## Training Model
5  XGB = XGB.fit(x_train_fos, y_train_fos)
6
7  ## Predicting response on Test
8  y_pred = XGB.predict(x_test)
9
10 XGB_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 XGB_R2 = metrics.recall_score(y_test, y_pred)*100
12 XGB_P2 = metrics.precision_score(y_test, y_pred)*100
13 XGB_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print("Accuracy: %.2f%%" % XGB_A2)
16 print("Recall_Accuracy: %.2f%%" % XGB_R2)
17 print("Precision_Accuracy: %.2f%%" % XGB_P2)
18 print("F1 Score: %.2f%%" %XGB_F2)
```

```
Accuracy: 98.44%
Recall_Accuracy: 88.84%
Precision_Accuracy: 95.67%
F1 Score: 92.13%
```

```
1  cm = confusion_matrix(y_test, y_pred)
2  TP_XGB_2 = cm[1][1]
3
4  plt.rcParams['figure.figsize'] = (5, 5)
5  #sns.set(style = 'dark', font_scale = 1.4)
6  sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8  print(confusion_matrix(y_test, y_pred))
```

```
[[1950  9]
 [ 25 199]]
```

Figure 17: Code for XGBoost Classifier

XGB with Stratified K-Fold

```
1  ## XGB for k = 10
2  skfold = StratifiedKFold(n_splits = 10,random_state=None)
3  model_skfold = XGBClassifier()
4  results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5  XGB_1 = (results_skfold.mean()*100)
6  print("Accuracy when k is 10 : %.2f%%" % XGB_1)
7
8  ## XGB for k = 20
9  skfold = StratifiedKFold(n_splits = 20,random_state=None)
10 model_skfold = XGBClassifier()
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 XGB_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % XGB_2)
14
15 ## XGB for K = 30
16 skfold = StratifiedKFold(n_splits = 30,random_state=None)
17 model_skfold = XGBClassifier()
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 XGB_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % XGB_3)
21
22 ## XGB for K = 50
23 skfold = StratifiedKFold(n_splits = 50,random_state=None)
24 model_skfold = XGBClassifier()
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 XGB_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % XGB_4)
```

```
Accuracy when k is 10 : 92.78%
Accuracy when k is 20 : 97.17%
Accuracy when k is 30 : 98.05%
Accuracy when k is 50 : 98.27%
```

Figure 18: K-Fold Validation of XGBoost

3.6.4 Random Forest

Random Forest was implemented using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for RF is illustrated in Figure 19 whereas Figure 20 represents its K-Fold validation.

Random Forest and its CM

```
1 ## Random Forest
2 RF = RandomForestClassifier()
3
4 ## Training Model
5 RF = RF.fit(x_train_fos, y_train_fos)
6
7 ## Predicting response on Test
8 y_pred = RF.predict(x_test)
9
10 RF_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 RF_R2 = metrics.recall_score(y_test, y_pred)*100
12 RF_P2 = metrics.precision_score(y_test, y_pred)*100
13 RF_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print("Accuracy: %.2f%%" % RF_A2)
16 print("Recall_Accuracy: %.2f%%" % RF_R2)
17 print("Precision_Accuracy: %.2f%%" % RF_P2)
18 print("F1 Score: %.2f%%" % RF_F2)
```

```
Accuracy: 98.40%
Recall_Accuracy: 87.50%
Precision_Accuracy: 96.55%
F1 Score: 91.80%
```

```
1 cm = confusion_matrix(y_test, y_pred)
2 TP_RF_2 = cm[1][1]
3
4 plt.rcParams['figure.figsize'] = (5, 5)
5 #sns.set(style = 'dark', font_scale = 1.4)
6 sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8 print(confusion_matrix(y_test, y_pred))
```

```
[[1952  7]
 [ 28 196]]
```

Figure 19: Code for Random Forest

Random Forest with Stratified K-Fold

```
1 ## Random Forest for k = 10
2 skfold = StratifiedKFold(n_splits = 10, random_state=None)
3 model_skfold = RandomForestClassifier()
4 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5 RF_1 = (results_skfold.mean()*100)
6 print("Accuracy when k is 10 : %.2f%%" % RF_1)
7
8 ## Random Forest for k = 20
9 skfold = StratifiedKFold(n_splits = 20, random_state=None)
10 model_skfold = RandomForestClassifier()
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 RF_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % RF_2)
14
15 ## Random Forest for K = 30
16 skfold = StratifiedKFold(n_splits = 30, random_state=None)
17 model_skfold = RandomForestClassifier()
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 RF_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % RF_3)
21
22 ## Random Forest for K = 50
23 skfold = StratifiedKFold(n_splits = 50, random_state=None)
24 model_skfold = RandomForestClassifier()
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 RF_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % RF_4)
```

```
Accuracy when k is 10 : 94.01%
Accuracy when k is 20 : 97.50%
Accuracy when k is 30 : 98.14%
Accuracy when k is 50 : 98.28%
```

Figure 20: K-Fold Validation of Random Forest

3.6.5 SVM-Linear

Yu et al. (2017) used SVM-Linear in their research and achieved a F1 Score of 90%. SVM was implemented with 'Linear' Kernel using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for SVM-Linear is illustrated in Figure 21 whereas Figure 22 represents its K-Fold validation.

SVM-Linear and its Confusion Matrix

```
1  ## SVM Linear
2  SVM= svm.SVC(kernel = 'linear')
3
4  ## Training Model
5  SVM = SVM.fit(x_train_fos, y_train_fos)
6
7  ## Predicting response on Test
8  y_pred = SVM.predict(x_test)
9
10 SVM_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 SVM_R2 = metrics.recall_score(y_test, y_pred)*100
12 SVM_P2 = metrics.precision_score(y_test, y_pred)*100
13 SVM_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print("Accuracy: %.2f%%" % SVM_A2)
16 print("Recall_Accuracy: %.2f%%" % SVM_R2)
17 print("Precision_Accuracy: %.2f%%" % SVM_P2)
18 print("F1 Score: %.2f%%" % SVM_F2)
```

```
Accuracy: 92.81%
Recall_Accuracy: 58.93%
Precision_Accuracy: 67.01%
F1 Score: 62.71%
```

```
1  cm = confusion_matrix(y_test, y_pred)
2  TP_SVM_2 = cm[1][1]
3
4  plt.rcParams['figure.figsize'] = (5, 5)
5  #sns.set(style = 'dark', font_scale = 1.4)
6  sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8  print(confusion_matrix(y_test, y_pred))
```

```
[[1894  65]
 [ 92 132]]
```

Figure 21: Code for SVM-Linear

SVM Linear with Stratified K-Fold

```
1  ## SVM Linear for k = 10
2  skfold = StratifiedKFold(n_splits = 10,random_state=None)
3  model_skfold = svm.SVC(kernel= 'linear')
4  results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5  SVM_1 = (results_skfold.mean()*100)
6  print("Accuracy when k is 10 : %.2f%%" % SVM_1)
7
8  ## SVM Linear for k = 20
9  skfold = StratifiedKFold(n_splits = 20,random_state=None)
10 model_skfold = svm.SVC(kernel= 'linear')
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 SVM_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % SVM_2)
14
15 ## SVM Linear for K = 30
16 skfold = StratifiedKFold(n_splits = 30,random_state=None)
17 model_skfold = svm.SVC(kernel= 'linear')
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 SVM_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % SVM_3)
21
22 ## SVM Linear for K = 50
23 skfold = StratifiedKFold(n_splits = 50,random_state=None)
24 model_skfold = svm.SVC(kernel= 'linear')
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 SVM_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % SVM_4)
```

```
Accuracy when k is 10 : 77.85%
Accuracy when k is 20 : 82.89%
Accuracy when k is 30 : 83.95%
Accuracy when k is 50 : 84.39%
```

Figure 22: K-Fold Validation of SVM-Linear

3.6.6 SVM-RBF

Adly et al. (2015) used SVM-RBF in their research and achieved an accuracy of 87.5%. SVM was implemented with 'Radial Basis Function' Kernel using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for SVM-RBF is illustrated in Figure 23 whereas Figure 24 represents its K-Fold validation.

SVM-RBF and its Confusion Matrix

```
1 ## SVM RBF
2 SVM= svm.SVC(kernel = 'rbf')
3
4 ## Training Model
5 SVM = SVM.fit(x_train_fos, y_train_fos)
6
7 ## Predicting response on Test
8 y_pred = SVM.predict(x_test)
9
10 SVM_RBF_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 SVM_RBF_R2 = metrics.recall_score(y_test, y_pred)*100
12 SVM_RBF_P2 = metrics.precision_score(y_test, y_pred)*100
13 SVM_RBF_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print("Accuracy: %.2f%%" % SVM_RBF_A2)
16 print("Recall_Accuracy: %.2f%%" % SVM_RBF_R2)
17 print("Precision_Accuracy: %.2f%%" % SVM_RBF_P2)
18 print("F1 Score: %.2f%%" % SVM_RBF_F2)
```

```
Accuracy: 97.76%
Recall_Accuracy: 87.50%
Precision_Accuracy: 90.32%
F1 Score: 88.89%
```

```
1 cm = confusion_matrix(y_test, y_pred)
2 TP_SVM_RBF_2 = cm[1][1]
3
4 plt.rcParams['figure.figsize'] = (5, 5)
5 #sns.set(style = 'dark', font_scale = 1.4)
6 sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8 print(confusion_matrix(y_test, y_pred))
```

```
[[1938  21]
 [ 28 196]]
```

Figure 23: Code for SVM-RBF

SVM RBF with Stratified K-Fold

```
1 ## SVM RBF for k = 10
2 skfold = StratifiedKFold(n_splits = 10,random_state=None)
3 model_skfold = svm.SVC(kernel= 'rbf')
4 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5 SVM_RBF_1 = (results_skfold.mean()*100)
6 print("Accuracy when k is 10 : %.2f%%" % SVM_RBF_1)
7
8 ## SVM RBF for k = 20
9 skfold = StratifiedKFold(n_splits = 20,random_state=None)
10 model_skfold = svm.SVC(kernel= 'rbf')
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 SVM_RBF_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % SVM_RBF_2)
14
15 ## SVM RBF for K = 30
16 skfold = StratifiedKFold(n_splits = 30,random_state=None)
17 model_skfold = svm.SVC(kernel= 'rbf')
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 SVM_RBF_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % SVM_RBF_3)
21
22 ## SVM RBF for K = 50
23 skfold = StratifiedKFold(n_splits = 50,random_state=None)
24 model_skfold = svm.SVC(kernel= 'rbf')
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 SVM_RBF_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % SVM_RBF_4)
```

```
Accuracy when k is 10 : 88.88%
Accuracy when k is 20 : 93.48%
Accuracy when k is 30 : 94.63%
Accuracy when k is 50 : 95.11%
```

Figure 24: K-Fold Validation of SVM-RBF

3.6.7 Naive Bayes

Naïve Bayes was implemented using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for NB is illustrated in Figure 25 whereas Figure 26 represents its K-Fold validation.

Naive Bayes and its Confusion Matrix

```
1 ## Naive Bayes
2 NB = GaussianNB()
3
4 ## Training Model
5 NB = NB.fit(x_train_fos, y_train_fos)
6
7 ## Predicting response on Test
8 y_pred = NB.predict(x_test)
9
10 NB_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 NB_R2 = metrics.recall_score(y_test, y_pred)*100
12 NB_P2 = metrics.precision_score(y_test, y_pred)*100
13 NB_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print("Accuracy: %.2f%%" % NB_A2)
16 print("Recall_Accuracy: %.2f%%" % NB_R2)
17 print("Precision_Accuracy: %.2f%%" % NB_P2)
18 print("F1 Score: %.2f%%" % NB_F2)
```

```
Accuracy: 72.47%
Recall_Accuracy: 58.04%
Precision_Accuracy: 20.41%
F1 Score: 30.20%
```

```
1 cm = confusion_matrix(y_test, y_pred)
2 TP_NB_2 = cm[1][1]
3
4 plt.rcParams['figure.figsize'] = (5, 5)
5 #sns.set(style = 'dark', font_scale = 1.4)
6 sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8 print(confusion_matrix(y_test, y_pred))
```

```
[[1452  507]
 [  94  130]]
```

Figure 25: Code for Naïve Bayes

Naive Bayes with Stratified K-Fold

```
1 ## Naive Bayes for k = 10
2 skfold = StratifiedKFold(n_splits = 10,random_state=None)
3 model_skfold = GaussianNB()
4 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5 NB_1 = (results_skfold.mean()*100)
6 print("Accuracy when k is 10 : %.2f%%" % NB_1)
7
8 ## Naive Bayes for k = 20
9 skfold = StratifiedKFold(n_splits = 20,random_state=None)
10 model_skfold = GaussianNB()
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 NB_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % NB_2)
14
15 ## Naive Bayes for K = 30
16 skfold = StratifiedKFold(n_splits = 30,random_state=None)
17 model_skfold = GaussianNB()
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 NB_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % NB_3)
21
22 ## Naive Bayes for K = 50
23 skfold = StratifiedKFold(n_splits = 50,random_state=None)
24 model_skfold = GaussianNB()
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 NB_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % NB_4)
```

```
Accuracy when k is 10 : 58.03%
Accuracy when k is 20 : 61.88%
Accuracy when k is 30 : 65.34%
Accuracy when k is 50 : 67.67%
```

Figure 26: K-Fold Validation of Naïve Bayes

3.6.8 KNN

Chien et al. (2012) used KNN in their research and achieved an accuracy of 75%. KNN was implemented using default parameters and was cross validated using K-fold validation with folds ranging from 10 to 50. Code for KNN is illustrated in Figure 27 whereas Figure 28 represents its K-Fold validation.

KNN and its Confusion Matrix

```
1 ## KNN
2 KNN = KNeighborsClassifier()
3
4 ## Training Model
5 KNN = KNN.fit(x_train_fos, y_train_fos)
6
7 ## Predicting response on Test
8 y_pred = KNN.predict(x_test)
9
10 KNN_A2 = metrics.accuracy_score(y_test, y_pred)*100
11 KNN_R2 = metrics.recall_score(y_test, y_pred)*100
12 KNN_P2 = metrics.precision_score(y_test, y_pred)*100
13 KNN_F2 = metrics.f1_score(y_test, y_pred)*100
14
15 print("Accuracy: %.2f%%" % KNN_A2)
16 print("Recall_Accuracy: %.2f%%" % KNN_R2)
17 print("Precision_Accuracy: %.2f%%" % KNN_P2)
18 print("F1 Score: %.2f%%" % KNN_F2)
```

```
Accuracy: 94.96%
Recall_Accuracy: 89.29%
Precision_Accuracy: 69.93%
F1 Score: 78.43%
```

```
1 cm = confusion_matrix(y_test, y_pred)
2 TP_KNN_2 = cm[1][1]
3
4 plt.rcParams['figure.figsize'] = (5, 5)
5 #sns.set(style = 'dark', font_scale = 1.4)
6 sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7
8 print(confusion_matrix(y_test, y_pred))
```

```
[[1873  86]
 [ 24 200]]
```

Figure 27: Code for KNN

KNN with Stratified K-Fold

```
1 ## KNN for k = 10
2 skfold = StratifiedKFold(n_splits = 10, random_state=None)
3 model_skfold = KNeighborsClassifier()
4 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
5 KNN_1 = (results_skfold.mean()*100)
6 print("Accuracy when k is 10 : %.2f%%" % KNN_1)
7
8 ## KNN for k = 20
9 skfold = StratifiedKFold(n_splits = 20, random_state=None)
10 model_skfold = KNeighborsClassifier()
11 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
12 KNN_2 = (results_skfold.mean()*100)
13 print("Accuracy when k is 20 : %.2f%%" % KNN_2)
14
15 ## KNN for K = 30
16 skfold = StratifiedKFold(n_splits = 30, random_state=None)
17 model_skfold = KNeighborsClassifier()
18 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
19 KNN_3 = (results_skfold.mean()*100)
20 print("Accuracy when k is 30 : %.2f%%" % KNN_3)
21
22 ## KNN for K = 50
23 skfold = StratifiedKFold(n_splits = 50, random_state=None)
24 model_skfold = KNeighborsClassifier()
25 results_skfold = cross_val_score(model_skfold, a_fos, b_fos, cv = skfold)
26 KNN_4 = (results_skfold.mean()*100)
27 print("Accuracy when k is 50 : %.2f%%" % KNN_4)
```

```
Accuracy when k is 10 : 91.12%
Accuracy when k is 20 : 93.89%
Accuracy when k is 30 : 94.84%
Accuracy when k is 50 : 95.17%
```

Figure 28: K-Fold Validation of KNN

3.6.9 Neural Network

Fernandes et al. (2020) used KNN in their research and achieved an accuracy of 89.64%. Basic Neural Network was designed and implemented. It was tested for epochs 25 and 50 with constant batch size of 60. Code for design, training and implementation of NN is illustrated in Figure 29.

Neural Network & its Confusion Matrix

```
1 ##Neural Network
2 NN = Sequential()
3 NN.add(Dense(51, input_dim = 100, activation = 'relu'))
4 NN.add(Dense(27, activation = 'relu'))
5 NN.add(Dense(15, activation = 'relu'))
6 NN.add(Dense(9, activation = 'relu'))
7 NN.add(Dense(6, activation = 'relu'))
8 NN.add(Dense(2, activation = 'sigmoid'))
9 NN.compile(loss='sparse_categorical_crossentropy', optimizer='adam', metrics=['accuracy'])

1 ## Training Neural Network
2 Neural_N3 = NN.fit(x_train_fos, y_train_fos, epochs=25, batch_size= 60)

1 ## Testing Neural Network
2 y_pred = NN.predict(x_test)
3
4 #Converting predictions to Label
5 pred = list()
6 for i in range(len(y_pred)):
7     pred.append(np.argmax(y_pred[i]))
8 NN_3 = accuracy_score(pred,y_test)*100
9 print('Accuracy is: %.2f%%' % NN_3)

1 ## Confusion Matrix for Neural Network
2 cm = confusion_matrix(y_test, pred)
3 TP_NN_3 = cm[1][1]
4 plt.rcParams['figure.figsize'] = (5, 5)
5 #sns.set(style = 'dark', font_scale = 1.4)
6 sns.heatmap(cm, annot = True, annot_kws = {"size": 15})
7 print(confusion_matrix(y_test, pred))
```

Figure 29: Code for Neural Network

References

- Adly, F., Yoo, P. D., Muhaidat, S., Al-Hammadi, Y., Lee, U. and Ismail, M. (2015). Randomized general regression network for identification of defect patterns in semiconductor wafer maps, *IEEE Transactions on Semiconductor Manufacturing* **28**(2): 145–152.
- Braha, D. and Shmilovici, A. (2002). Data mining for improving a cleaning process in the semiconductor industry, *IEEE Transactions on Semiconductor Manufacturing* **15**(1): 91–101.

- Chien, C.-F., Hsu, C.-Y. and Chen, P.-N. (2012). Semiconductor fault detection and classification for yield enhancement and manufacturing intelligence, *Flexible Services and Manufacturing Journal* **25**.
- Fernandes, S., Antunes, M., Santiago, A., Barraca, J., Gomes, D. and Aguiar, R. (2020). Forecasting appliances failures: A machine-learning approach to predictive maintenance, *MDPI Journals* .
- Yu, C., Chien, C. and Kuo, C. (2017). Exploit the value of production data to discover opportunities for saving power consumption of production tools, *IEEE Transactions on Semiconductor Manufacturing* **30**(4): 345–350.