# Applying Machine Learning Technique to Detect Failures in Hard Disk Drive Test Process

## ABSTRACT

This paper presents machine learning techniques to detect the failure in hard disk drive manufacturing test process. The data is high dimensionality and highly imbalance. Feature selection technique with filter method and embedded method with light gradient boost are used to reduce the dimension of data. We apply three techniques: SMOTE, Different Cost and SMOTE with Different Cost to handle imbalance data. Several machine learning methods are compared. The XGBoost with SMOTE and XGBoost with Different Cost (XGB DC) give the best performance with 91% ROC AUC and 73% PRC AUC. The SVM algorithm shows good performance on ROC AUC while low performance on PRC AUC. The XGBoost algorithm shows good performance of both ROC AUC and PRC AUC.

*Keywords*: SVM, imbalance data, high dimensionality, feature selection, hard disk drive

## 1 INTRODUCTION

Machine learning technique is widely used in manufacturing process. In hard disk drive manufacturing machine learning is used to improve the productivity, such as parameter improvement, anomaly detection and failure detection. There are several processes in hard disk drive manufacturing. Failure should be detected as early as possible in manufacturing process to reduce waste time and reduce cost. Each process has measured parameters and the data stored in database. The failure detection with machine learning technique is an opportunity to reduce the cost and increase productivity.

## 2 LITERATURE REVIEW

Support vector machine is one of popular machine learning algorithm for classification on two classes data. SVM is used to detect the failure in hard disk drive assembly process by using voice coil motor current to train the model (Simongyi & Chongstitvatana, 2018). The data set is imbalance, Fail case is only 3%. SVM algorithm is able to classify with 100% accuracy. Contrast to this previous work, this paper studies the failure during test process. Data set are the parameters collected from assembly and servo track write processes.

There are several research applying SVM and SMOTE technique to handle imbalance data such as Akbani et al. (2004). They applied SVM with SMOTE technique and used different error cost, called SDC (SMOTE with Different Cost). The SDC method gives the best performance compared to SVM and SVM with SMOTE.

SVM with different method to handle imbalance data is used (Tang et al., 2009). Several methods are studied: SVM-weight is a cost sensitive learning method, SVM-SMOTE oversampling minority class, SVM-RANDU under sampling and GSVM-RU random under sampling only the data which no vector supported. Comparing performance with 4 metrics, the GSVM-RU gives the best performance.

Tong and Daphne (2001) uses SVM with imbalance data by applying active learning to reduce the train data size and obtain better performance. Chakravarthy et al. (2019) studied C50, KNN, NN, RF and SVM with Random oversampling (ROSE) and SMOTE oversampling with different ratio. The SMOTE with 1:3 oversampling ratio gives best performance in all models. A fuzzy support vector machines is introduced (Ma

42 et al., 2011) for class imbalance learning. There are several studies of the problem of class imbalance such
43 as Guo et al., (2008), Chakravarthy et al., (2019), Wang and Japkowicz (2009).

44 Synthetic Minority Over-sampling Technique (SMOTE) (Chawla et al., 2002) is one of popular method to
45 oversampling the minor class. Borderline SMOTE (Han et al., 2005) is another sampling method that develop
46 from SMOTE to sample only the borderline of minor class. Sharma et al., (2018) studied oversampling
47 minority class with consider on major class. Sampling with Majority (SWIM) creates minor class with similar
48 Mahalanbois distance of the majority class. In Khire et al., (2019), Chandrashekar and Sahin (2014) the
49 features input to the model is studied and the result showed that they are important to the performance of the
50 model. The work in Zhang et al., (2017) uses SSVM-FS to select features. This method focuses on the
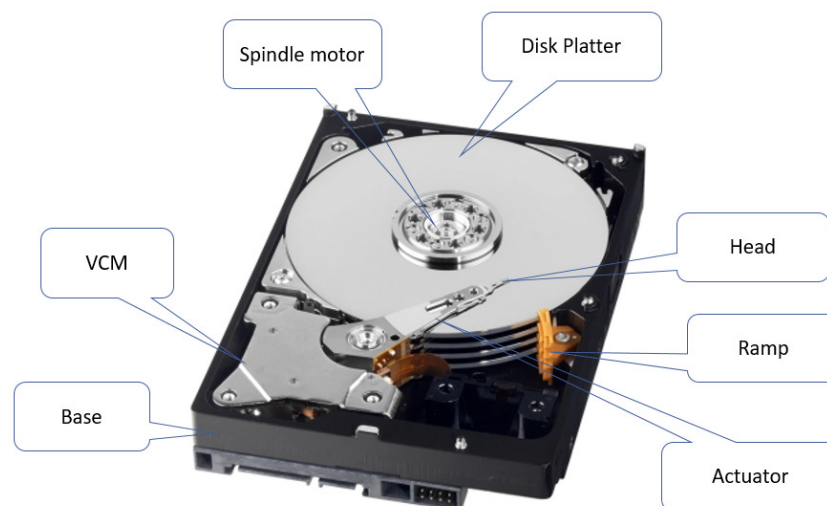51 imbalance class. The weight of SVM indicated the important features.

52 The Extreme Gradient Boost (XGBoost) is another machine learning algorithm which is a scalable tree
53 boosting system. It is widely used in machine learning competition to achieve state-of-art result. It is an
54 implement of gradient boost. XGBoost able to run beyond billions of examples on few resources (Chen &
55 Guestrin, 2016).

56 Using accuracy to measurement the performance of machine learning algorithm on imbalance data is not
57 appropriate. The model performance will show high accuracy when the model fails to predict the minor class.
58 The performance measurement with area under curve (AUC) is not affected by the ratio of class (skew),
59 while the accuracy, F1-score, Cohen's kappa and Krippendorf's are affected by skew (Jeni et al., 2013). The
60 Receive Operation Characteristic (ROC) plot gives the overview performance. Precision-Recall plot (PRC)
61 gives accurate prediction performance (Saito & Rehmsmeier, 2015).

62 This paper studies SVM and XGBoost algorithm to detect the failure in hard disk drives test process. We
63 also applying SMOTE, Different Cost and SMOTE with Different Cost techniques to train the model. We
64 measure the performance with area under curve of Receiver Operating Characteristics (ROC AUC) and area
65 under curve of associated Precision/Recall (PRC AUC) values. The rest of the paper is as follows Section 3
66 explained hard disk drive manufacturing process. Section 4 described failure detection with machine learning
67 techniques. Section 5 reported the experimental result. Conclusion is given in Section 6.

## 68 3  HARD DISK DRIVE MANUFACUTRING PROCESS

69 Hard disk drive process started in clean room to assembly all components such as base, spindle motor, head
70 and disk together (Fig.1). The complete unit is sent out from clean room to write servo track pattern and test
71 process. The servo track write is the process to write reference position signal on disk. The reference signal
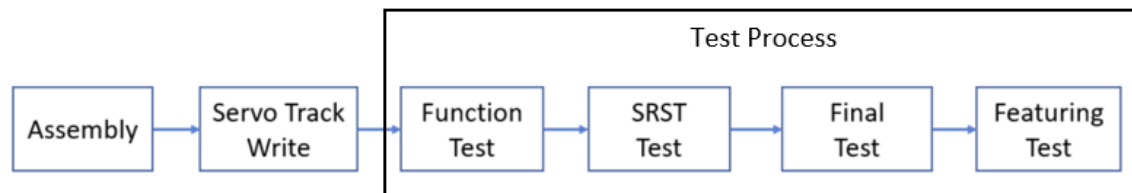72 is written on control circular track and it is consistency space of track.

73

74

Figure 1.   Hard disk drive components.

75 The test process consists of 4 main processes: Function test, Self Run Stress test (SRST), Final test and
76 Featuring test. Function test starts with micro code download to functional the hard disk drive, measures and
77 adjust parameters for best seek, read and write performance. SRST is the test to measure and analyze disk
78 surface then records the defect location. This test is performed at high temperature. Also Re-adjust parameters
79 and test customer functionality. It is the longest test process. Final Test is also called performance test. It
80 tests the whole surface read write performance. Featuring test is the test to adjust parameters to meet customer
81 requirement (Fig. 2).

82



83 Figure 2. Hard disk drive manufacturing process.

## 4 FAILURE DETECTION USING MACHINE LEARNING

84

### 4.1 Data Collection

85

86 Data set are collected from hard disk drives manufacturing process which is separated in two parts. The first
87 part is parameters collected from Assembly and Servo track write processes. Parameters are the measurement
88 during assembly parts and write servo signal, such as motor current, head resistance, servo signal quality,
89 distance between head and disk total 359 parameters. The second part is the target output collected from test
90 process. The target output are 2 classes, negative class (Bad) and positive class (Good). Negative class data
91 are very small compare to positive class. The ratio of failure to passer is 1:100. Data are collected with
92 sampling passer and whole failure. Total data is 84325 rows. Each row represents for one hard disk drive
93 (HDD) data. There are 79448 rows of passer and 4877 rows of failure. 70% of data is used for training and
94 30% for validation (Table 1.).

95 Table 1: Data set.

| | Input Data | Train Data | Validation Data | Unit |
|---|---|---|---|---|
| **Total Data** | 84325 | 59027 | 25298 | HDD |
| **Pass** | 79448 | 56613 | 22835 | HDD |
| **Fail** | 4877 | 3414 | 1463 | HDD |

### 4.2 Data Pre-processing

96

97 Data is pre-processed by eliminating parameters with missing value that are more than 40% and fill the
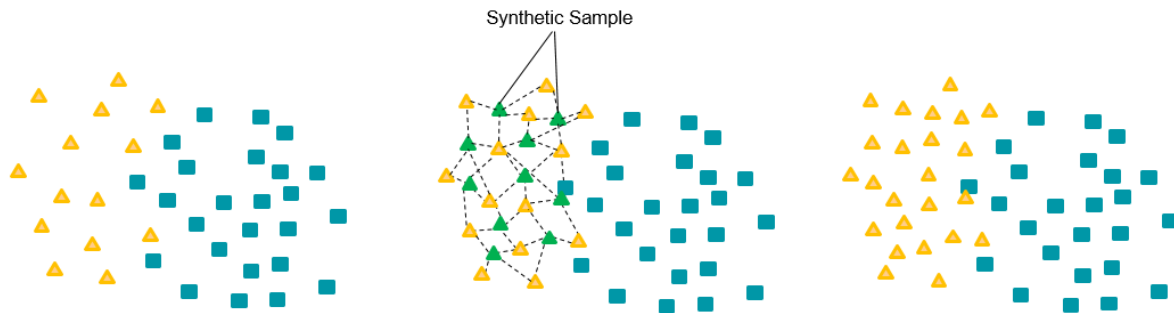98 missing data with mean value. Data is normalized with z-score.

### 4.3 Feature Selection

99

100 Feature selection is one of the methods to improve machine learning performance by reducing the redundant
101 and unnecessary features. Filter method is the method to eliminate the feature by ranking according to
102 importance with statistical measure such as Chi-square, ANOVA and correlation coefficient. Wrapper
103 method selects subset of feature that give high performance on machine learning model such as Forward
104 Selection, Backward Elimination and Recursive Feature Elimination. The wrapper method costs high
105 computation time (Khaire et al., 2019).

106 There are 359 features in the data set which is very high dimension. To input all data to train the machine
107 learning model will cost high computation time and low performance. The filter method is used to remove
108 the constant features, duplicated feature and correlated features. The group of correlate features defined by
109 Pearson's correlation coefficient more than 95% and select only one feature from each group with best AUC
110 on random forest model. The embedded method with light gradient boost to rank the important of feature is
111 used. To reduce the variation, feature important is rank with 10 iterations of accumulate values. 289
112 parameters are eliminated by feature selection. 70 parameters are input to train the model.

## 113 4.4 Imbalance Data Handling

114 The data set shows very small numbers of bad compare to good. Data set is highly imbalance. The SMOTE
115 (Chawla et al., 2002) is one of the methods to handle imbalance data by oversampling the minority class. It
116 uses the actual data to generate the synthetic sample of the minor class (Fig. 3). Another approach is assigning
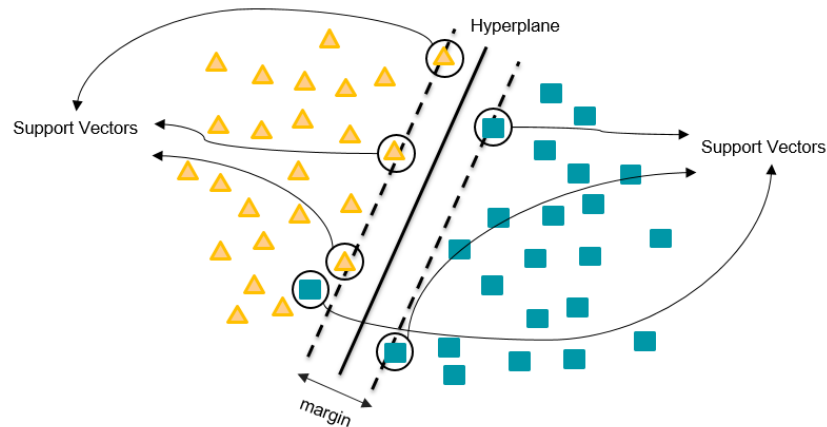117 different cost of the learning class.



118
119 Figure 3. SMOTE: Synthetic Minority Over-sampling Technique

## 120 4.5 Training Model

121 Support Vector Machine (SVM) is a supervised learning algorithm. It is one of popular classification method.
122 SVM algorithm finds the hyperplane that can separate the class of data with maximum margin. The closest
123 data to hyperplane of each class is called support vector. The margin measures from support vector to
124 hyperplane (Fig. 4).
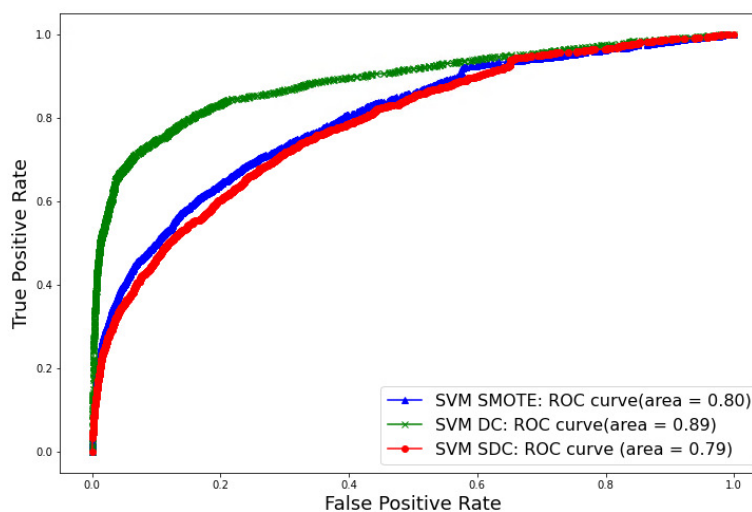


125
126 Figure 4. Support Vector Machine

127 Extreme Gradient Boosting (XGBoost) is an ensemble learning algorithm. It is a scalable machine learning
128 for tree boosting (Chen & Guestrin,2016). The boosting is sequential learner algorithm and use the error from
129 previous learner to improve the accuracy of the next learner. The model is added and learn until the accuracy

130  is not improved. XGBoost is a gradient boosting method which is fast and uses minimal resources. It is widely
131  used algorithm and achieve state-of -art result.

132  We compare 6 methods. There are 3 methods training with SVM algorithm and another 3 methods training
133  with XGBoost algorithm: 1.) SVM with SMOTE (SVM SMOTE), 2.) SVM with Different Cost (SVM DC),
134  3.) SVM with SMOTE and Different Cost (SVM SDC), 4.) XGBoost with SMOTE (XGB SMOTE), 5.)
135  XGBoost with Different Dost (XGB DC) and 6.) XGBoost with SMOTE and Different Cost (XGB SDC).
136  The SVM using 'rbf' kernel, class_weight equal to ratio of passer and failure and others hyper parameter
137  setting are as follows. Degree = 3, gamma = 'scale' and max_iter = 1. The XGBoost using booster = 'gbtree',
138  eta = 0.3, gamma = 0, max_depth = 6 and scale_pos_weight equal to ratio of passer and failure.
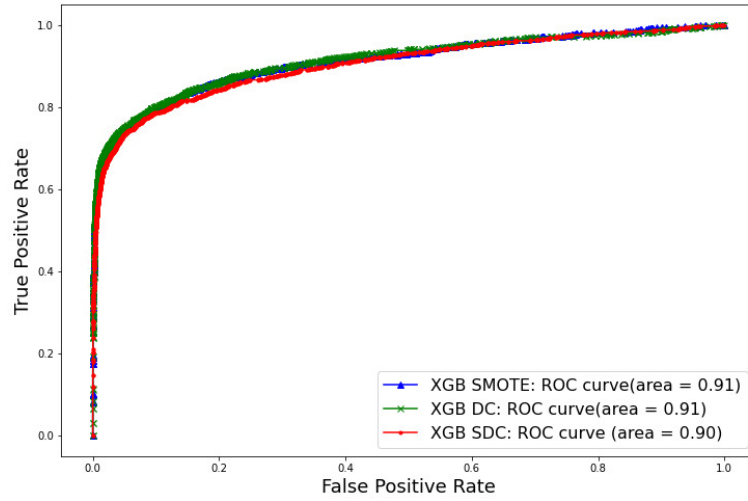
## 5  EXPERIMENT RESULT

140  The result of classifying the validation data of 25298 rows are as follows. The performance measurement
141  with ROC AUC: SVM with SMOTE gives 80%, SVM with Different Cost (SVM DC) gives 89% and SVM
142  with SMOTE and Different Cost (SVM SDC) gives 79% (Fig. 5). XGBoost with SMOTE (XGB SMOTE)
143  and XGBoost with Different Cost (XGB DC) give the best performance at 91% and XGBoost with SMOTE
144  and Different Cost (XGB SDC) gives 90% (Fig. 6). The performance measurement with PRC AUC: SVM
145  SMOTE 32%, SVM DC 59% PRC AUC and SVM SDC has the lowest performance at 29% (Fig. 7). XGB
146  SMOTE and XGB DC give the best performance at 73% while XGB SDC gives 71% (Fig. 8). The model
147  with Different Cost gives the best performance on both SVM and XGBoost algorithm. The XGBoost
148  algorithm in all methods give better performance than SVM (see Table 2).
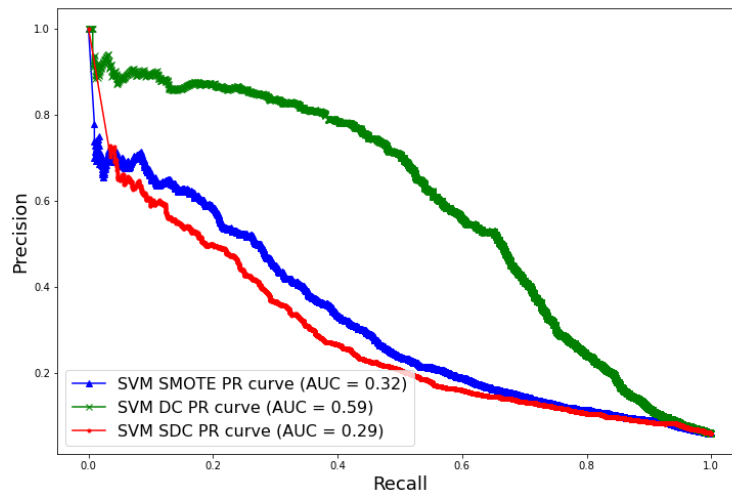
149



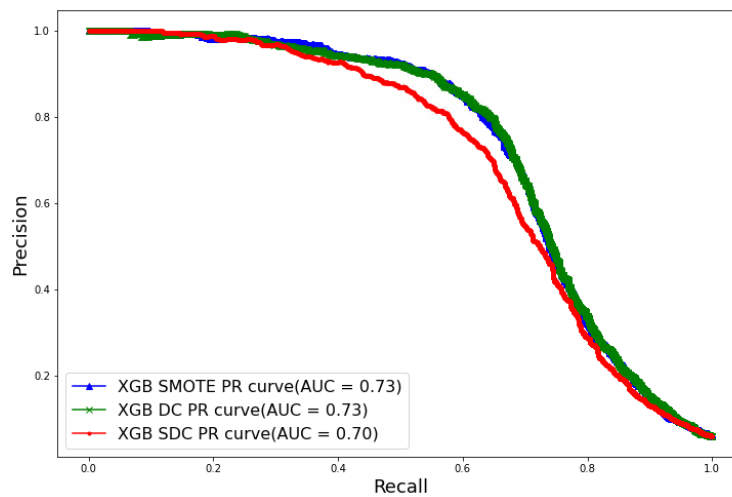150  Figure 5.   SVM ROC AUC plot

151

152
Figure 6.   XGBoost ROC AUC plot



153

154
Figure 7.   SVM PRC AUC plot



155

156
Figure 8.   XGboost PRC AUC plot

157     Table 2: Result

| Method | ROC AUC | PRC AUC |
|--------|---------|---------|
| **SVM SMOTE** | 80% | 32% |
| **SVM DC** | 89% | 59% |
| **SVM SDC** | 79% | 29% |
| **XGB SMOTE** | 91% | 73% |
| **XGB DC** | 91% | 73% |
| **XGB SDC** | 90% | 70% |

158

# 6   CONCLUSION

160 This study presents the method of failure detection with SVM and XGBoost algorithm. The proposed method
161 employs feature selection and data imbalance handling. The experiment performs on real hard disk drive
162 manufacturing data. The feature selection method are filter and embedded algorithm. The parameters are
163 reduced from 359 to 70 to be input to the model. The model training with SVM and XGBoost algorithm with
164 3 different data imbalance handling methods: 1) SMOTE, 2) Different Cost and 3) SMOTE with Different
165 Cost. The XGBoost algorithm has the better performance than SVM. The XGBoost with SMOTE and
166 XGBoost with DC give the best performance with 91% ROC AUC and 73% PRC AUC. The SVM algorithm
167 shows high performance on ROC AUC measurement while low performance on PRC AUC. XGBoost shows
168 good performance on both of two measurements.

# 7   REFERENCES

170 Akbani, R., Kwek, S., & Japkowicz, N. (2004). Applying support vector machines to Imbalanced datasets.
171     Machine Learning: ECML 2004, 39-50. https://doi.org/10.1007/978-3-540-30115-8_7
172 Chakravarthy, A. D., Bonthu, S., Chen, Z., & Zhu, Q. (2019). Predictive models with Resampling: A
173     comparative study of machine learning algorithms and their performances on handling Imbalanced
174     datasets. 2019 18th IEEE International Conference On Machine Learning And Applications
175     (ICMLA). https://doi.org/10.1109/icmla.2019.00245
176 Chandrashekar, G., & Sahin, F. (2014). A survey on feature selection methods. Computers & Electrical
177     Engineering, 40(1), 16-28. https://doi.org/10.1016/j.compeleceng.2013.11.024
178 Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority
179     over-sampling technique. Journal of Artificial Intelligence Research, 16, 321-357.
180     https://doi.org/10.1613/jair.953
181 Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. Proceedings of the 22nd
182     ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794.
183     doi:10.1145/2939672.2939785
184 Guo, x., Yin, Y., Dong, C., Yang, G., & Zhou, G. (2008). On the Class Imbalance Problem. 2008 Fourth
185     International Conference on Natural Computation, 192-201.
186     https://doi.org/10.1109/ICNC.2008.871
187 Han, H., Wang, W., & Mao, B. (2005). Borderline-SMOTE: A new over-sampling method in Imbalanced
188     data sets learning. Lecture Notes in Computer Science, 878-887.
189     https://doi.org/10.1007/11538059_91
190 Jeni, L. A., Cohn, J. F., & De La Torre, F. (2013). Facing Imbalanced data--recommendations for the use of
191     performance metrics. 2013 Humaine Association Conference on Affective Computing and
192     Intelligent Interaction. doi:10.1109/acii.2013.47
193 Khaire, U. M., & Dhanalakshmi, R. (2019). Stability of feature selection algorithm: A review. Journal of
194     King Saud University - Computer and Information Sciences.
195     https://doi.org/10.1016/j.jksuci.2019.06.012
196 Ma, H., Wang, L., & Shen, B. (2011). A new fuzzy support vector machines for class imbalance learning.
197     2011 International Conference on Electrical and Control Engineering, 3781-3784.
198     https://doi.org/10.1109/iceceng.2011.6056838

199  Saito, T., & Rehmsmeier, M. (2015). The precision-recall plot is more informative than the ROC plot when
200      evaluating binary classifiers on Imbalanced datasets. PLOS ONE, 10(3), e0118432.
201      doi:10.1371/journal.pone.0118432
202  Sharma, S., Bellinger, C., Krawczyk, B., Zaiane, O., & Japkowicz, N. (2018). Synthetic Oversampling with
203      the majority Class: A new perspective on handling extreme imbalance. 2018 IEEE International
204      Conference on Data Mining (ICDM). https://doi.org/10.1109/icdm.2018.00060
205  Simongyi, M., & Chongstitvatana, P. (2018). Abnormality detection in hard disk drive assembly process
206      using support vector machine. 2018 15th International Conference on Electrical
207      Engineering/Electronics, Computer, Telecommunications and Information Technology (ECTI-
208      CON), 612-615. https://doi.org/10.1109/ecticon.2018.8619935
209  Tong, S., & Daphne, K. (2001). Support Vector Machine Active Learning With Applications To Text
210      Classification. The Journal of Machine Learning Research, 2(1), 45-66.
211      https://doi.org/10.1162/153244302760185243
212  Wang, B. X., & Japkowicz, N. (2009). Boosting support vector machines for imbalanced data sets.
213      Knowledge and Information Systems, 25(1), 1-20. https://doi.org/10.1007/s10115-009-0198-y
214  Yuchun Tang, Yan-Qing Zhang, Chawla, N., & Krasser, S. (2009). SVMs modeling for highly Imbalanced
215      classification. IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics), 39(1),
216      281-288. https://doi.org/10.1109/tsmcb.2008.2002909
217  Zhang, C., Wang, G., Zhou, Y., Yao, L., Jiang, Z. L., Liao, Q., & Wang, X. (2017). Feature selection for
218      high dimensional imbalanced class data based on F-measure optimization. 2017 International
219      Conference on Security, Pattern Analysis, and Cybernetics (SPAC).
220      https://doi.org/10.1109/spac.2017.8304290
221