

Representation Learning and Predictive Classification: Application with an Electric Arc Furnace

L.D. Rippon^{1a}, I. Yousef^a, B. Hosseini^b, A. Bouchoucha^b, J.F. Beaulieu^b, C. Prévost^b, M. Ruel^b, S.L. Shah^c, R.B. Gopaluni^a

^a*University of British Columbia, Vancouver, B.C., Canada*

^b*BBA Engineering Consultants, Mont-Saint-Hilaire, Q.C., Canada*

^c*University of Alberta, Edmonton, A.B., Canada*

Abstract

Data-driven disciplines such as biostatistics and chemometrics are undergoing a period of transformation propelled by powerful advances in computational hardware, parallel processing and algorithmic efficiency. Process systems engineering is positioned for concurrent advances in data-driven sub-disciplines such as modeling, optimization, control, fault detection and diagnosis. This work embodies this transformation as it addresses a novel industrial fault detection problem from both traditional and contemporary approaches to process analytics. Traditional approaches such as partial least squares are compared with powerful new techniques inspired by deep representation learning such as convolutional neural networks. Novel contributions include the formulation and introduction of a novel industrial predictive classification problem, the design and implementation of a comprehensive machine learning workflow that converts raw industrial data into critical operational insights, and the presentation of a robust comparative analysis between traditional and contemporary approaches to representation learning and binary classification. Specifically, this work addresses the unexpected loss of plasma arc in the electric arc furnace of a large-scale metallurgical process. The objective is to learn an efficient and informative representation from the raw industrial data that enables the prediction of an arc loss event such that operators can take corrective actions. A comprehensive representation learning and predictive classification framework is presented for development of the inferential sensor from large quantities of historical industrial process data.

Keywords: Process Monitoring and Diagnostics, Machine Learning, Deep Learning, Process Control, Inferential Sensor, Applications

1. Introduction

Representation learning is described as a subset of machine learning (ML) and a superset of deep learning. Classical ML is distinguished from representation learning through the selection of features. In classical ML features are hand-designed whereas in representation learning the features are learned from the data. Moreover, in deep representation learning there exist additional layers of abstraction between simple learned features and more complicated features that may further improve representation [1]. Representation learning is defined as the means by which an efficient and informative representation can be learned that extracts useful information to improve the performance of classification, regression or prediction models [2]. The labor intensive procedure of engineering features is excellent at leveraging application specific domain knowledge but this approach lacks efficiency and ease of applicability across various domains [3]. With representation learning the important discriminatory features can be learned from the data in a systematic fashion allowing for much faster deployment of effective ML models for problems from a variety of domains [2].

The use of deep representation learning algorithms for process systems engineering (PSE) applications (e.g., control, process monitoring, and fault detection) is a relatively new, but highly active research area [4]. The Tennessee Eastman (TE) process has been used as a fault detection and diagnosis (FDD) benchmark to validate many advanced neural architectures including stacked sparse auto-encoders [5], deep belief networks [6], and deep convolutional neural networks (CNNs) [7]. A limitation of these studies is that they do not extend the validation to industrial case studies which include many added complexities. The design of the three phase flow facility by Cranfield University has improved this predicament by providing a benchmark for comparative process monitoring studies based on real experimental data [8, 9]. This work offers a different form of contribution as it validates and compares traditional and advanced process monitoring methods on historical data taken from a real, large-scale industrial process. Another limitation of these studies is that they focus on detecting and diagnosing faults that have already occurred. Alternatively, this work is distinguished by the fact that we aim to predict our fault approximately five minutes before it occurs.

In this paper we develop and introduce an end-to-end workflow for industrial fault prediction that includes data pre-processing, representation learning and predictive classification. Traditional and contemporary approaches to representation

36 learning and binary classification are compared in a comprehensive analysis for
 37 their ability to predict the fault using large quantities of historical industrial data.
 38 Extensive cross-validation and hyper-parameter optimization trials are performed
 39 to obtain rigorous empirical results and maximize the accuracy of fault prediction.
 40 The fault that we aim to detect and predict is the unexpected loss of plasma arc in
 41 an industrial direct current (DC) electric arc furnace (EAF) that serves as a smelter
 42 to refine ores into base metals.

43 Ore is transported from the mine and passed through a series of hammer mill
 44 grinders, flash dryers, preheater cyclones, calciner combustion chambers and flu-
 45 idized bed reducers. This processing provides a fine particulate feed that is dried,
 46 heated and reduced to maximize the efficiency of the energy intensive twin elec-
 47 trode DC EAF, illustrated in Figure 1. An open plasma arc spans from the graphite
 48 electrodes (i.e., the cathode) to the surface of the molten slag (i.e., the anode) pro-
 49 viding energy required to maintain the slag and alloy at target temperatures over
 50 1400°C [10]. The roof and side walls of the furnace are water cooled whereas the
 51 bottom anode is air cooled to maintain safe structural temperatures [11]. Hot off-
 52 gas released from the furnace is recycled to provide upstream preheating. The feed
 53 enters from multiple ports along the roof whereas the slag and alloy are tapped in-
 54 termittently from launders [12]. This work is directly relevant to a variety of EAF
 operations including nine in the Canadian steel-making industry [13].

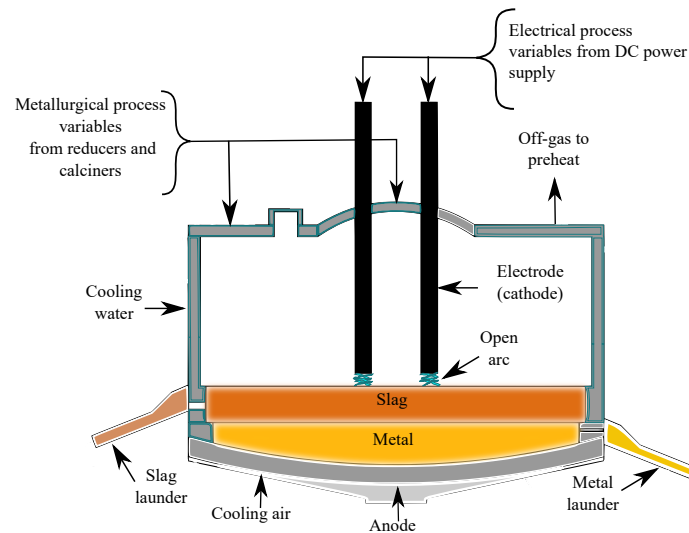


Figure 1: An illustration of a direct current electric arc furnace [14].

55 Stable EAF operation is critical for maximizing production efficiency and
 56 profitability. Unexpected loss of the plasma arc is a recurring and unresolved
 57 fault that significantly impacts the production rate and the electrical efficiency of
 58 the furnace. There are three primary categories of suspected arc loss mechanisms,
 59 i.e., electrical disturbances from the DC power supply, feed disturbances from
 60 the upstream metallurgical processes and the operation of the EAF. Therefore, a
 61 broad process aspect ratio is considered in the representation learning analysis that
 62 includes dozens of measured variables (MVs) from the power supply, numerous
 63 upstream unit operations, and the EAF. Moreover, an entire year of high frequency
 64 operating data is collected and analyzed to develop the arc loss predictor.

65 The goal of the fault predictor is to provide operators with a warning five to ten
 66 minutes in advance of an event with a 75% or higher probability of inducing arc
 67 loss such that operators can take preventative measures. Operators require at least
 68 two minutes prior to the arc loss event in order to take the necessary corrective
 69 actions. Figure 2 illustrates the entire ML workflow including the data prepro-
 70 cessing tasks resulting in segmented datasets ready for representation learning
 71 and predictive classification. Note, some classification methods bypass explicit
 72 representation learning and instead learn from the raw features. There are also
 73 hidden feedback connections between the modules as the workflow progresses in
 74 a largely iterative fashion. Initially this alarm will serve as a tool for engineers
 75 and operators but ultimately the goal is to implement an advanced controller that
 76 can automatically take corrective action.

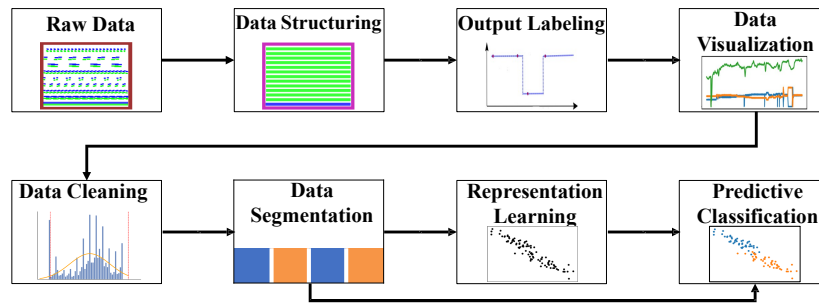


Figure 2: Flowchart illustrating the overall data analytics workflow.

77 The novel contributions presented in this work include the introduction of the

78 DC EAF arc loss problem and the formulation of this problem as a supervised
79 ML problem. Successful problem formulation is a significant contribution that
80 includes transforming a year of raw industrial operating data into cleaned, struc-
81 tured, labeled, and segmented datasets that are amenable to further statistical ML
82 analysis. Labeling the data requires the introduction of rigorous quantitative con-
83 ditions to detect the arc loss. Given a precise problem formulation and procured
84 training data, the remaining contribution is the development of the arc loss predic-
85 tion inferential sensor. This contribution also includes a comprehensive validation
86 and comparison of traditional and advanced approaches to representation learning
87 and predictive classification on real industrial operating data.

88 **2. Data preprocessing and visualization**

89 Data preprocessing produces the datasets that are used to train, validate, and
90 test the predictive models. Therefore, the preprocessing methodology is a signifi-
91 cant factor for the generalization performance of a supervised ML algorithm [15].
92 The goal of the preprocessing module is to transform the raw process historian
93 data into a form that is amenable for statistical ML algorithms. Moreover, this
94 transformation should maximally retain information from the raw data, minimize
95 extraneous information injected during preprocessing and remove redundant data.
96 Finally, the preprocessing procedure should maintain generality and flexibility for
97 efficient deployment to other PSE applications. The remainder of this section in-
98 cludes a description of the raw data followed by an overview of the methods used
99 to structure, visualize, clean and segment the data.

100 **2.1. Wrangling big data - size and quality**

101 Big data is a contentious term primarily because the meaning of *big* in data
102 science has undergone rapid semantic changes as the standard size of data-sets
103 across various disciplines and domains has grown rapidly. Not only is the amount
104 or volume of data context dependent but the velocity, variety and veracity of typi-
105 cal data-sets vary widely across both time and disciplines [16]. This work qualifies
106 as big data analysis from both contextual and pragmatic perspectives as most in-
107 ferential sensing literature in PSE relies on significantly less data (e.g., Tennessee
108 Eastman and penicillin fermentation benchmarks) and there is too much data for
109 straightforward processing on most consumer-grade hardware.

110 The raw data used in this work involves one year of daily exports from a
111 real industrial process historian. Our scope encompasses approximately all of the
112 process variable data collected from a metallurgical process from the milling of
113 crushed ore to the refining of base metals. Each day of operation is captured and
114 stored as a comma separated value (CSV) file with approximately 228 columns
115 and thirty thousand rows. Half of the 228 columns are process variables (PVs) and
116 the other half are corresponding timestamps. In total, the entire year of operating
117 data has an uncompressed size of 17.4 gigabytes (GB) and thus requires more than
118 16 GB of random access memory (RAM) to simply load the data into a data-frame.

119 **2.2. Data structuring and output labeling**

120 The columns of each daily export have a varying number of rows with more
121 densely sampled PVs having up to thirty thousand rows and others having as few
122 as ten samples. The raw data contains asynchronous data with both numerical PVs
123 such as furnace temperature and categorical PVs such as valve positions. The raw
124 data contains errors such as missing values, bad inputs and not a number (NaN)
125 values. Systematically structuring the raw data and removing the corrupted data
126 is one of the first stages of preprocessing.

127 Each of the 365 daily CSV file exports is loaded as a data-frame to replace
128 non-numeric inputs (e.g., ‘tag not found’) with NaN values and remove rows and
129 columns with overwhelming NaN values (e.g., rows with less than three non NaN
130 values). The illustration in Figure 3 shows the preprocessing of three consecutive
131 days and represents time horizontally. Each MV is represented by a green row and
132 the accompanying timestamp is represented by a blue row with the dashed rows
133 representing the differing frequency of measurements. The top of Figure 3 shows
134 the structured dataset that has an equivalent number of samples for each MV, one
135 unified timestamp and no NaN values. To preserve information the most densely
136 sampled variable from each day is identified and the corresponding timestamp is
137 used as a unified timestamp (blue bars at the top of Figure 3) for all PVs. To
138 minimize insertion of synthetic data the less frequently sampled variables are re-
139 sampled using a simple forward fill or zero order hold operation.

140 Once the data is cleaned and structured it amounts to only 5.6 GB of uncom-
141 pressed memory with each file having approximately 110 columns and twenty
142 eight thousand rows per column. The data is now suitable for generating the arc
143 loss labels. To ensure the labels are robust, all three conditions in Figure 4 regard-

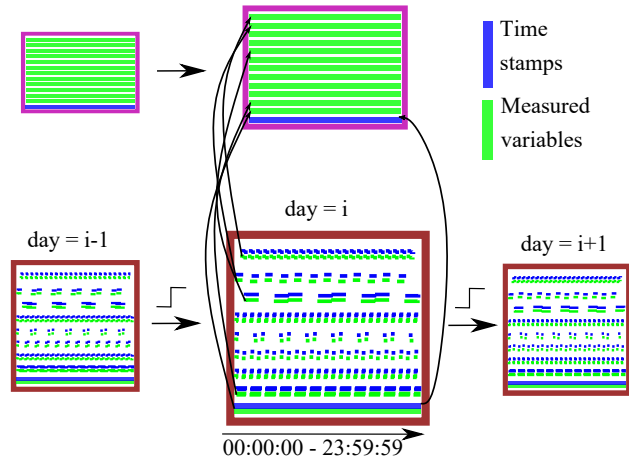


Figure 3: Preprocessing consecutive days of historical data.

144 ing the power of an electrode must be met in order to constitute a loss of arc in
 145 that electrode. Specifically, the power must be stable within a standard deviation
 146 of 2 MW for approximately 11.5 minutes, there must be a power drop greater than
 147 10 MW within the past 36s, and the power must recover to within 5 MW of the
 148 original stable power within a period of approximately 10 minutes. These condi-
 149 tions are applied to each sample for both electrodes to generate output labels that
 150 are binary indicators of arc loss in the respective electrode.

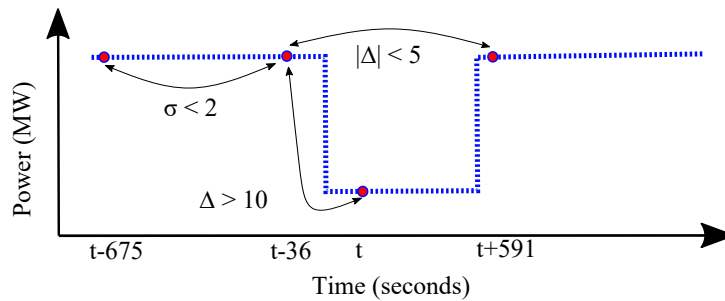


Figure 4: Illustration and quantitative definition of conditions constituting arc loss.

151 2.3. Data visualization

152 Data visualization provides key insights into the frequency of the faults and
 153 the severity of the arc loss on overall production efficiency. Visualization also as-

154 sists in troubleshooting and validation of data pre-processing and output labeling,
155 respectively. A sanity check is performed on the binary arc loss labels, such as
156 that shown in Figure 5, to ensure that they correspond to a representative power
157 drop. Three discrete arc loss labels are shown in the top plot of Figure 5 and the
158 corresponding drops in power of arc A are shown in the bottom plot of Figure 5.

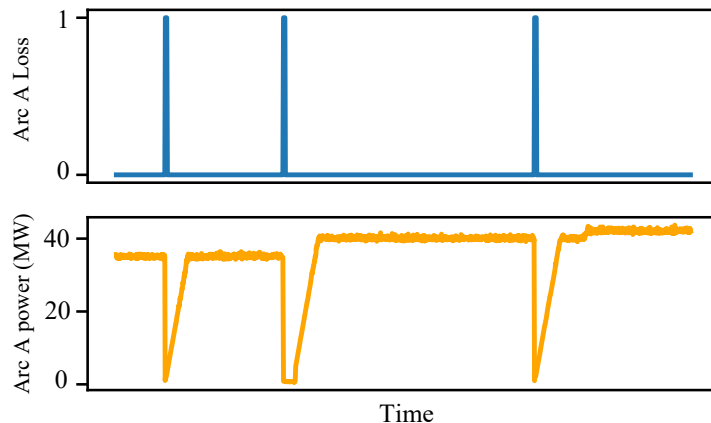


Figure 5: Visual validation of the arc loss labels.

159 The frequency of arc loss events is clear from Figure 6 which shows the num-
160 ber of arc loss events per day for each plasma arc throughout a year of operation.
161 Arc loss is a significant problem that can either occur as often as twenty five times
162 per day (indicating a chain of arc losses) or not at all for multiple consecutive
163 days. This distinction provides motivation to apply data-driven pattern recogni-
164 tion techniques to determine the difference in operation between arc loss cascades
165 and stable operation. Although the average sequence of positive arc loss indica-
166 tions is less than one minute in duration, the disruption to the EAF of a single loss
167 event can cause up to twenty minutes of lost production. This visualization not
168 only provides motivation but it also helps to recognize the class imbalance in our
169 output labels due to the short average duration of each arc loss indication. Class
170 imbalance is an important consideration addressed at the end of Section 2.

171 Finally, the severity of the arc loss fault on EAF operation is visualized in
172 Figure 7 by comparing a period of relatively stable operation (top) to a period
173 of faulty operation (bottom) using the power applied to each arc and the furnace
174 feed rate. The occurrence of arc loss has a significant impact on the furnace
175 feed rates and subsequently on the production rate of the EAF. Thus, it is imperative to

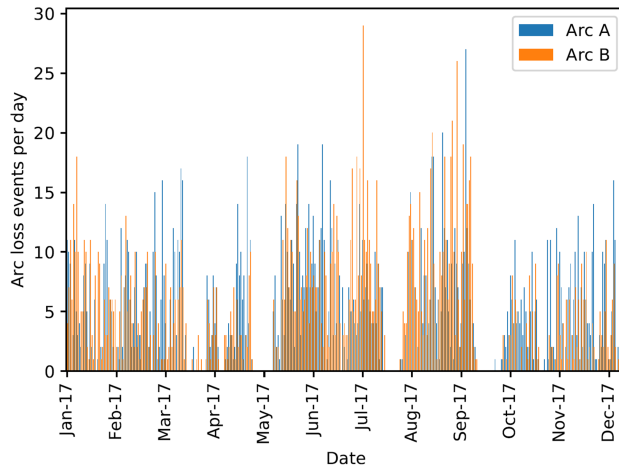


Figure 6: Daily arc loss events in each electrode over one year of operation.

176 prevent loss of the plasma arc in order to sustain economic viability of the process.

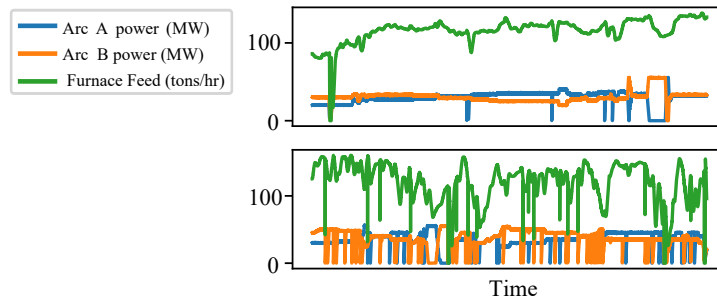


Figure 7: Visually comparing stable operation (top) to faulty operation (bottom).

177

178 2.4. Data cleaning

179 The quality of any ML model depends on the quality of the input it receives.
 180 Here, data cleaning involves setting PV limits using process knowledge to fil-
 181 ter out nonsensical values (e.g., negative feed rates), removing problematic PVs,
 182 and removing data from plant shut-downs. Erroneous process data and outliers
 183 can induce spurious correlations and increase the rate of misclassification for ML
 184 classifiers [17]. Removal of this data is accomplished through domain expertise

185 and consultations with our collaborators at BBA. A set of minimum and maximum
186 limits are agreed upon for each PV and measurements outside of these limits are
187 set to either the nearest limit or three standard deviations from the mean.

188 The left side of Figure 8 shows the power values for arc A as a histogram with
189 a normally distributed probability density function (PDF) and PV limits shown by
190 the vertical red lines. There are some negative power values that are subsequently
191 adjusted to zero during data cleaning. Using process knowledge to set PV limits
192 is not an infallible strategy. Sanity checks are necessary to ensure the PV limits
193 are correct as demonstrated by the right side of Figure 8 which shows the crucible
194 heat loss as a PDF with the original PV limits as vertical red lines. All of the
195 crucible heat loss data is outside the original PV limits but instead of cleaning this
196 data the PV limits are re-evaluated and it is deemed acceptable. This PV limit and
sanity check procedure was conducted for all of the PVs.

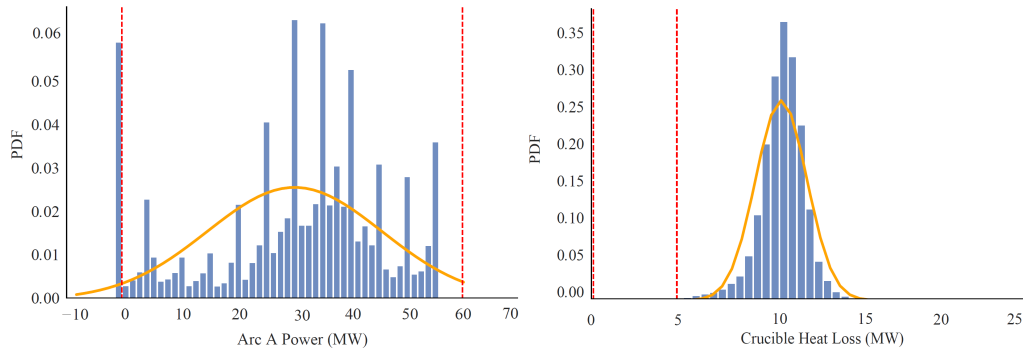


Figure 8: Setting PV limits with process knowledge to filter out erroneous data.

197

198 Outliers are often considered to be values that are greater than three standard
199 deviations away from the mean. Box plots are commonly used to show the dis-
200 tribution of a variable and indicate the number of outliers. Treatment of outliers
201 is application specific and modeling abnormal behavior requires retention of data
202 that may be statistically defined as outliers. Figure 9 shows the use of box plots
203 to visualize the number of furnace feed values greater than 1.5 times the inter-
204 quartile range (shown by the black circles). The statistical outliers at the top of
205 Figure 9 are considered feasible whereas the outliers with a negative value at the
206 bottom of Figure 9 are deemed irrelevant and are removed from the data.

207 Two final tasks remain for data cleaning, i.e., removing unhelpful PVs and

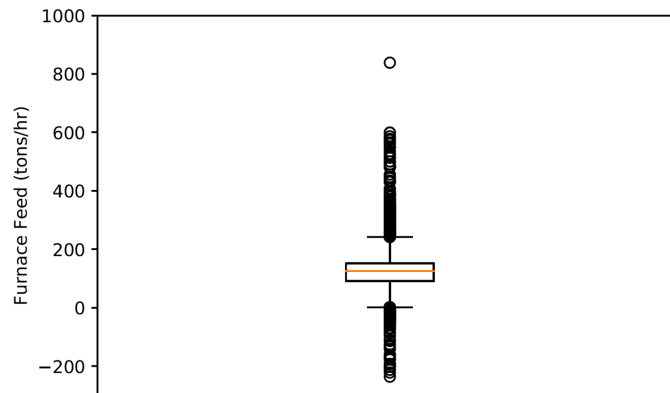


Figure 9: Box-plots illustrating outliers that represent abnormal process behavior.

208 removing shut-down data that is not representative of the process during operation.
 209 Five laboratory measurements were deemed to have too low of a sampling
 210 frequency for use as a fault predictor and were therefore removed entirely from
 211 the dataset. Seven PVs were removed based on prior knowledge that they had
 212 consistently faulty or inaccurate measurements (i.e., PVs that are not decommissioned
 213 from the historian correctly). A shut-down is a series of steps to take a chemical
 214 process from a normal state of operation to an idle state of operation
 215 until all required maintenance is complete. Two plant shutdowns are clearly visible
 216 in Figure 6 in May and early October. The data for these periods is carefully
 217 removed for all PVs during data cleaning to preserve useful information during
 218 the initial shutdown and initial startup phases. By carefully structuring and cleaning
 219 the data we obtain a significantly smaller set of data that has preserved useful
 220 process information and is more amenable to subsequent modeling.

221 2.5. Data segmentation

222 For binary classification problems with a large degree of class imbalance, the
 223 vast majority of instances fall into the majority class while significantly fewer instances
 224 fall into the minority class (i.e., the class of interest for fault detection).
 225 Most binary classification methods perform poorly on imbalanced datasets due to
 226 assuming the data are drawn from the same distribution and assigning equal error
 227 value to both classes. Classifiers aim to achieve the highest accuracy along the
 228 whole range of data and therefore tend to largely ignore the minority class which
 229 has relatively negligible impact [18]. Previous studies suggest techniques to ad-

230 dress class imbalance can be divided mainly into three categories: re-sampling,
 231 feature engineering, and classifier manipulation [19]. Artificially re-sampling the
 232 instances to balance class distributions can be performed by either under-sampling
 233 the majority class or over-sampling the minority class [20]. Under-sampling is at
 234 risk of discarding information from the majority class while over-sampling in-
 235 creases the likelihood of over-fitting by duplicating instances from the minority
 236 class [21]. More advanced methods include explicitly combining separate fea-
 237 tures from the minority and majority classes as well as manipulating the classifier
 238 weights internally [22].

239 The arc loss dataset is highly imbalanced with 99.67% of the samples labeled
 240 as the majority class (i.e., no arc loss) and only 0.33% of the instances labeled with
 241 arc loss. A random under-sampling approach is taken to address class imbalance
 242 by extracting a segment that contains 55 minutes worth of data in the 5-60 minute
 243 period before every arc loss. All 1526 arc loss events (taken from both arc A and
 244 arc B) are extracted to represent the minority class. The majority class is randomly
 245 under-sampled and only 1526 segments that correspond to 55 consecutive minutes
 246 taken 5 minutes prior to periods of extended stable operation are extracted. The
 247 data segmentation process is illustrated in Figure 10. The entire dataset containing
 248 3052 segments is further divided with 85% (or 2594 balanced segments) for cross-
 249 validated training and 15% (or 458 balanced segments) for testing. With the data
 250 finally procured to a suitable format it can be used to train the representation
 251 learning and predictive classification algorithms.

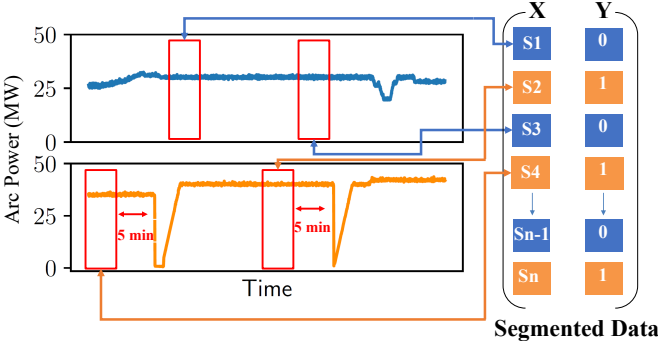


Figure 10: Illustration of data segmentation to create a balanced data-set.

252 3. Learned representations and predictive classifiers

253 This work focuses on studying and validating the benefits of using repre-
254 sentation learning (e.g., dimensionality reduction) and deep learning for predictive
255 classification with real industrial operating data. A comprehensive methods se-
256 lection is illustrated in Figure 11 with explicit representation learning algorithms
257 on the left and the predictive classifiers on the right. Partial least squares (PLS)
258 and principal components analysis (PCA) are compared for generating explicit
259 representations while logistic regression (LR), linear support vector classifier (L-
260 SVC), kernel SVC (K-SVC), artificial neural networks (ANNs), and CNNs are
261 all compared as predictive classification models. Note, the not applicable (N/A)
262 indicates the use of raw features instead of explicit representation, but the K-SVC,
263 ANN, and CNN methods have internal representations with kernels, hidden layers
264 and convolutions, respectively. Altogether, Figure 11 shows fifteen experimental
265 combinations with seven algorithms that are introduced in what follows.

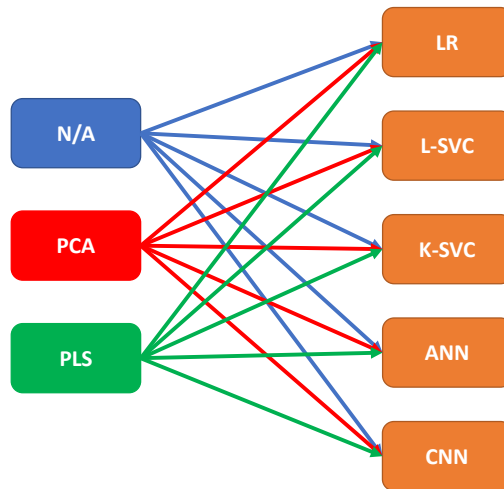


Figure 11: Experimental configurations with representations (left) and predictive classifiers (right).

266 3.1. Explicit representations with reduced dimensionality

267 Two traditional process analytics tools are applied to learn explicit dimension-
268 ally reduced representations from raw features, i.e., PCA and PLS.

269 3.1.1. Principal component analysis

270 The PCA statistical procedure was introduced in the early 20th century to de-
271 compose a multivariate dataset into a basis set of linearly uncorrelated orthogonal
272 variables called principal components [23]. It was subsequently developed for use
273 in multivariable quality control and has since been further extended and applied in
274 PSE where it is categorized along with PLS and ANNs as a quantitative process
275 history based method for FDD [24, 25, 26, 27, 28]. The convention for fault de-
276 tection is to calculate the Hotelling T^2 statistic with the largest singular values and
277 the Q statistic with the smallest singular values. The T^2 statistic defines normal
278 process behavior and any observation vectors that fall outside of the T^2 region
279 indicate that a fault has occurred. Alternatively, the Q statistic is used to define
280 a threshold that indicates whether or not the characteristics of the measurement
281 noise have changed significantly [29, 30].

282 Consider a pre-processed set of historical data that has been centered (i.e., col-
283 umn means subtracted), $X \in \mathbb{R}^{n \times d}$ where X includes the output label data as
284 additional columns. The covariance matrix of X is denoted $S \in \mathbb{R}^{d \times d}$ and is
285 given by $S = X^T X / (n - 1)$. The typical eigendecomposition of S is given by
286 $S = V \Lambda V^{-1}$ where the i -th column of $V \in \mathbb{R}^{d \times d}$ is the eigenvector v_i of S or
287 alternatively, the loading vectors or principal directions of the data X . The diag-
288 onal matrix $\Lambda \in \mathbb{R}^{d \times d}$ contains eigenvalues λ_i of decreasing magnitude. Given a
289 symmetric matrix S with distinct eigenvalues λ_i , the eigenvector columns of V are
290 orthogonal (i.e., $V^{-1} = V^T$) and the eigendecomposition becomes $S = V \Lambda V^T$.
291 The principal components or principal component scores can be calculated by
292 projecting the data onto the principal directions, i.e., $C = XV$, where the i -th
293 column of C is the i -th principal component of X [31]. Alternatively, PCA can
294 be conducted with singular value decomposition of the centered data matrix X
295 where singular values (σ) are related to eigenvalues, i.e., $\lambda_i = \sigma_i^2 / (n - 1)$ [30].

296 3.1.2. Partial least squares

297 As with PCA, PLS (also known as projection to latent structures) is a linear
298 representation learning method with a rich history of use in PSE. The PLS ap-
299 proach was first introduced by Herman Wold in the 1970s and has since been used
300 extensively in chemical process industries as a chemometrics method for applica-
301 tions such as FDD [32, 33]. One drawback of PCA is that although some principal
302 components may describe significant variance in X , those same principal com-

303 ponents might not be relevant for predicting the output labels. As a supervised
 304 learning method, PLS regression (PLSR) maximizes the covariance between the
 305 input data, X , and output data (or labels), $Y \in \mathbb{R}^{n \times d_y}$, in the latent space via the
 306 non-linear iterative partial least squares (NIPALS) algorithm [34].

The centered input matrix X and output matrix Y are each decomposed as,
 $X = LP^\top + E$ and $Y = MQ^\top + F$, where $L \in \mathbb{R}^{n \times a}$ and $M \in \mathbb{R}^{n \times a}$ are latent
 score matrices, $P \in \mathbb{R}^{d_x \times a}$ and $Q \in \mathbb{R}^{d_y \times a}$ are loading matrices, $E \in \mathbb{R}^{n \times d_x}$ and
 $F \in \mathbb{R}^{n \times d_y}$ are residual matrices and a is the PLS component or reduction order
 [35]. The iterative PLSR algorithm initializes $X_1 := X$ and $Y_1 := Y$ and then pro-
 ceeds to maximizing $l_i^\top m_i$ (for each iteration i) by initializing m_1 as one column
 of Y and solving the following set of equations until convergence is achieved:

$$w_1 = \frac{X_1^\top m_1}{\|X_1^\top m_1\|}, \quad l_1 = X_1 w_1, \quad q_1 = \frac{Y_1^\top l_1}{\|Y_1^\top l_1\|}, \quad \text{and} \quad m_1 = Y_1 q_1, \quad (1)$$

307 where $\|\cdot\|$ represents the Euclidean norm or ℓ_2 norm. The X-weights (w_1) are
 308 updated with the Y-scores (m_1) until the change in l_1 is negligible or below some
 309 specified error [32, 35]. The same procedure is repeated for the next iteration by
 310 replacing X and Y with the residual matrices, i.e., $X_{i+1} = E_i = X_i - l_i p_i^\top$ and
 311 $Y_{i+1} = F_i = Y_i - m_i q_i^\top$ where $p_i = X_i^\top l_i / \|l_i^\top l_i\|$.

312 **3.2. Predictive classification and implicit representations**

313 The right side of Figure 11 lists the five predictive classification methods that
 314 are trained and tested with the raw features, representations learned through PCA,
 315 and representations learned through PLS.

316 **3.2.1. Logistic regression**

Choosing LR for binary classification is natural as the standard logistic func-
 tion (i.e., the sigmoid function) given by

$$P(Z) = \frac{\exp(Z)}{1 + \exp(Z)} = \frac{1}{1 + \exp(-Z)} \quad (2)$$

317 provides a bounded output between zero and one that can be interpreted as the
 318 probability of a binary outcome and mapped to discrete classes (e.g., arc loss or
 319 no arc loss). The input, $Z = \alpha + \beta X$, to the logistic function illustrates the

320 connection with linear regression where X is the preprocessed data (or a learned
321 representation thereof), α is a scalar bias and β is a weight vector. Historically,
322 LR dates back to the early 19th century when the logistic function was invented to
323 describe population growth and autocatalytic chemical reactions [36]. Recent ap-
324 plications of LR in PSE include methods that combine LR with dominant trend ex-
325 traction and dependent binary relevance classifiers to perform nonstationary fault
326 diagnosis and multi-label fault classification, respectively [37, 38].

327 **3.2.2. Support vector classifiers**

328 The basis for the L-SVC predictive classification technique used in this work
329 is the soft margin support vector machine (SVM) (or support vector network)
330 introduced in 1995 which is itself an extension of the hard margin SVM, concep-
331 tualized solved in 1965 [39, 40]. The difference between hard margin and soft
332 margin SVM is that hard margin SVM assumes the classes are linearly separable
333 and thus tries to find a hyperplane such that no point is misclassified whereas soft
334 margin SVM allows for some misclassification that is proportionally penalized
335 in the objective function. Binary SVC aims to construct a separating hyperplane
336 between the two classes of data such that the margin (i.e., distance) between the
337 hyper-plane and the nearest data points of each class is maximized [41].

338 Nonlinear formulations of SVMs utilize the kernel trick, i.e., kernel SVC (K-
339 SVC), such as the parametric polynomial kernel or the non-parametric radial basis
340 function kernel with important properties that allow for enhanced representation
341 capacity and efficient optimization [34, 42]. Recent applications of SVMs in PSE
342 include applying one-class SVM on finite impulse response (FIR) data to detect
343 model-plant mismatch (MPM) in a paper machine control system [43, 44] and
344 using nonlinear SVM-based feature selection for FDD [45]. This work studies
345 linear and kernel based SVCs with a variety of configurations (e.g., kernel and
346 regularizer choices) provided in Section 4.

347 **3.2.3. Artificial neural networks**

348 Conceptually, ANNs were inspired by the structure and function of neurons
349 in the human brain [46]. Neural networks have undergone at least three historical
350 waves of popularity beginning with cybernetics in the mid 20th century, connec-
351 tionism in the late 1900s and the current manifestation of deep learning that began
352 in 2006 [1]. The deep learning wave of popularity resulted from a breakthrough

353 in the efficiency of training deep networks by Geoff Hinton’s research group, re-
 354 ferred to as greedy layerwise unsupervised training [2]. Deep learning tackles the
 355 problem of representation learning by using complex neural architectures to gen-
 356 erate nested representations that are functions of simpler representations [1]. The
 357 versatility and non-linear representation capacity of ANNs has drawn immense
 358 interest from the scientific community as a classifier for modeling complex rela-
 359 tionships [47].

The perceptron, introduced by the psychologist Frank Rosenblatt, is the first and most simple example of a modern neural network that was explicitly used for binary classification of linearly separable functions [48][49]. The perceptron is a building block for complex multi-layered ANNs that involve input layers, hidden layers, and output layers consisting of neurons connected with learned weights [50]. Linear combinations of inputs and weights at each layer are followed by nonlinear activation functions that help with increasing the depth of the network and modeling nonlinear relationships [51]. It is important to select a suitable activation function as it can have a significant influence on the ANN performance [52]. Types of activation functions include sigmoid functions, rectified linear unit (ReLU) functions softmax functions and many more. The non-linear output after the application of the activation function is represented by:

$$Z = \left(\sum_{i=1}^n x_i w_i + b \right) \quad y = f(Z), \quad (3)$$

360 where x_1, x_2, \dots, x_n represent the n inputs of the perceptron, w_1, w_2, \dots, w_n are
 361 the weights given to the respective input, b is a bias term, and f represents the
 362 chosen activation function for this layer.

In the context of binary classification, the output layer of the ANN consists of a single output neuron that indicates the class of the segment by computing the weighted sum of hidden values from the last hidden layer, followed by a sigmoid function, i.e.,

$$Z = \left(\sum_{i=1}^n x_i^L w_i^L + b \right) \quad y = f(Z) = \frac{\exp(Z)}{1 + \exp(Z)}, \quad (4)$$

363 where the superscript L refers to values and weights from the neurons in the final
 364 hidden layer. If the output of the sigmoid neuron is greater than or equal to 0.5, it
 365 outputs 1 (i.e., arc loss). However, if the output is less than 0.5, it outputs 0 (i.e.,

366 stable operation). In this work, flattened input segments are fed to a multi-layered
 367 fully connected perceptron model to predict arc loss as illustrated in Figure 12.
 368 Back propagation is used to train the network by propagating the error from the
 output layer to the hidden layer to update the weight matrix [53].

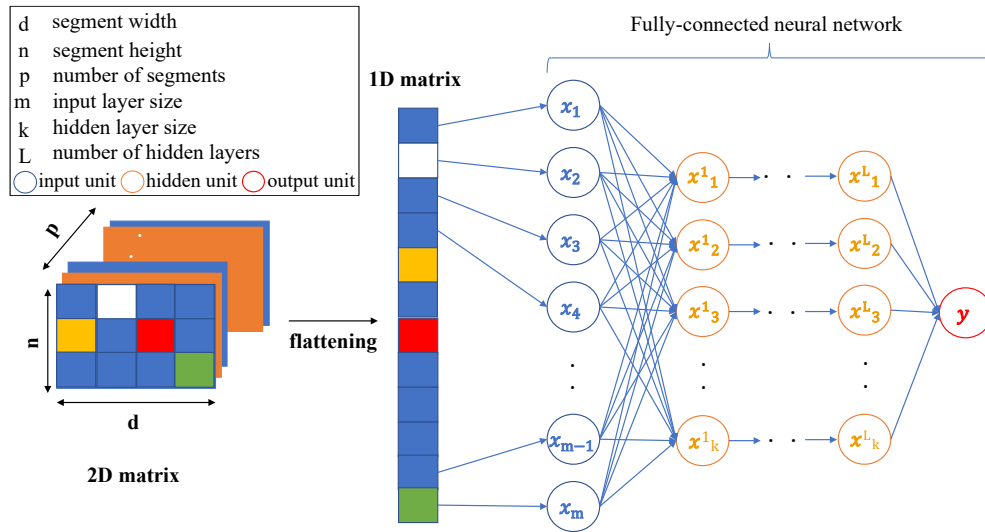


Figure 12: Illustrating the use of a fully-connected ANN to predict arc loss from input segments.

369

370 3.2.4. Convolutional neural networks

371 In the late 1980s CNNs were introduced to address visual pattern recognition
 372 problems such as handwritten digit recognition [54]. Instead of exclusively using
 373 fully-connected layers, CNNs use local connections (i.e., local receptive fields)
 374 to extract elementary features which are then combined by subsequent layers in a
 375 hierarchical feature extraction procedure [55]. Convolution with a kernel whose
 376 weights are learned through back-propagation creates the local receptive field for
 377 each feature map [56]. The receptive field and the dimensions of the resulting
 378 feature map are governed by the size of the kernel and the stride that the kernel
 379 takes over the input image (or input feature map). For a single input image, the
 380 number of output feature maps after the first convolution layer is equivalent to the
 381 number of learnable kernels specified for that layer.

For convolutional layer l , the output of the j th feature map is given by [57]:

$$\mathbf{x}_j^l = f\left(\sum_{i \in M_j} \mathbf{x}_i^{l-1} * \mathbf{k}_{i,j}^l + b_j^l\right) \quad (5)$$

where M_j is the set of input feature maps, $*$ is the discrete convolution operation, $k_{i,j}$ is the learnable kernel from input map i to output map j , and b_j is the additive bias for output map j . It is common to follow convolution layers in a CNN with a sub-sampling procedure known as a pooling layer. The output feature map from the convolutional layer is sub-sampled to create a lower dimensional feature map by applying a receptive field that converts the output at a certain location to a summary statistic of nearby outputs [1]. Two common types of sub-sampling operations are max pooling and average pooling. Pooling layers can help to improve computation and prevent overfitting [7]. As shown in Figure 13, the output of the last pooling layer is flattened before being passed to a fully-connected network.

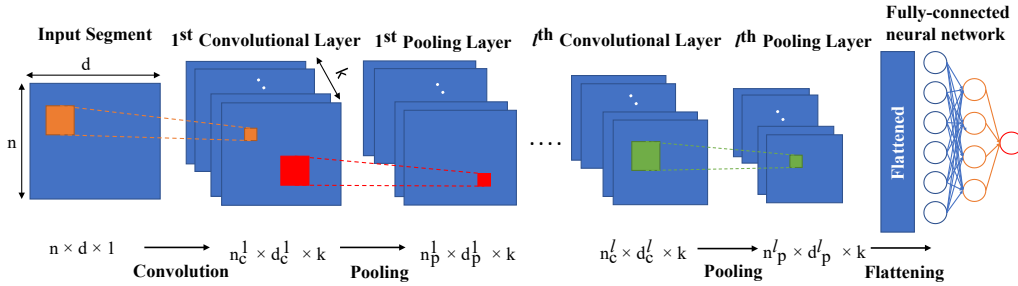


Figure 13: Illustration of the architecture for the CNN predictive classifier

State of the art performance has been achieved using CNNs on object recognition and natural language processing tasks [2]. Moreover, CNNs have been studied for PSE applications such as FDD on the TE process [7], the three phase flow facility at Cranfield University [9], and in a semiconductor manufacturing process [58]. To our knowledge this is the first time CNNs have been studied for fault prediction in an industrial manufacturing process with historical operating data. Many alternative algorithms could be chosen, but the purpose of this study is to compare traditional methods with contemporary methods in order to highlight areas for future investigation. We intend to explore more advanced representation learning methods (e.g., auto-encoders) and we also hope to release this dataset to the PSE community so that others can develop inferential sensors with better performance.

404 4. Experimental Setup

405 The experiments primarily consist of training, validating, and testing the fif-
406 teen experimental configurations shown in Figure 10. Our experimental setup has
407 the following two key factors that distinguish our work from previous FDD stud-
408 ies in PSE: i) simulating a production trial by preserving the temporal integrity
409 of our training data with respect to our testing data and ii) performing rigorous
410 cross-validation and hyperparameter optimization to compare models.

411 The preprocessed data segments are split into two groups; the training and
412 validation group consists of 2594 segments and the testing group consists of 458
413 segments. Prior experimental designs performed random selection of segments
414 for training and testing sets throughout the entire year of operation. Random sam-
415 pling is common in literature as well, but because we aim to develop an inferential
416 sensor for a real industrial process our experimental design mimics that of a pro-
417 duction trial. As shown in Figure 14, our simulated production trial trains and
418 validates models on the first ten months of operation while the last two months of
operation are strictly used for testing the final models.

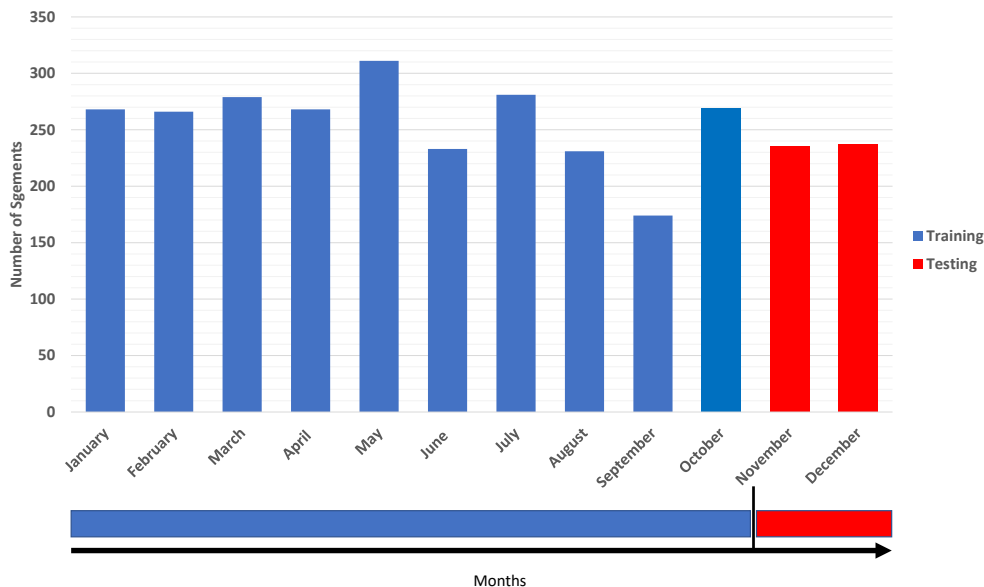


Figure 14: Dividing segments into training and testing sets based on their date.

419

4.1. Stratified k -fold cross-validation

A stratified k -fold cross-validation strategy is used to compare different hyperparameter configurations in our predictive classifiers. As shown in Figure 15, ten non-overlapping folds are created where each fold contains a balanced number of arc loss segments and stable operating segments. For each hyperparameter optimization trial (i.e., corresponding to a specific configuration) the model is trained on 90% of the training data and validated on the remaining 10%. This is repeated ten times, once for each fold, where the validation data changes as shown by the yellow highlight in 15. The result of the trial is the average accuracy of all ten validations which serves as a score to rank the hyperparameter configuration.

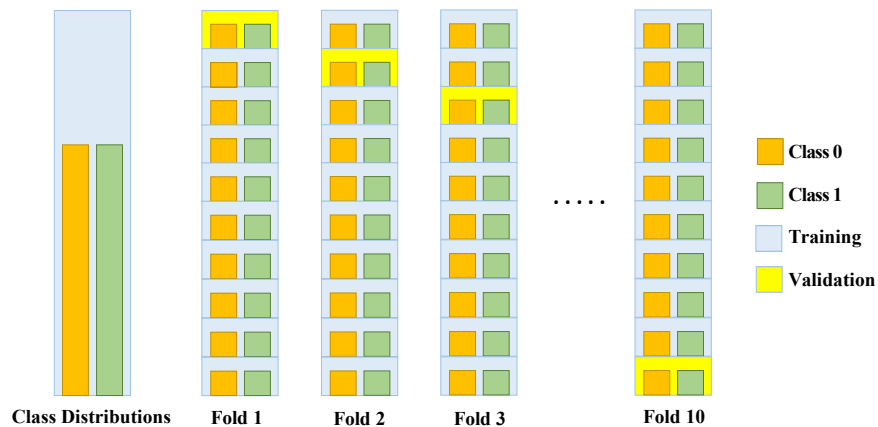


Figure 15: Each model is trained and validated with stratified k -fold cross-validation.

4.2. Hyperparameter optimization

A robust and transparent hyperparameter optimization strategy is critical for an impartial comparison of ML algorithms and the reliable development of a predictive inferential sensor. The efforts taken here aim to contribute a high level of rigor to hyperparameter optimization in the context of PSE. A broad space of possible hyperparameters is defined for each predictive classifier and then a Bayesian sequential model-based optimization (SMBO) algorithm searches this space using a tree-structured Parzen estimator (TPE) to suggest the best configurations by maximizing expected improvement (EI) [59][60]. Multiple trials are conducted for each hyperparameter configuration where the TPE specifies the configuration for the next trial based on the EI.

441 **4.2.1. Optimizing hyperparameters for learned representations**

442 Recall from Figure 10, the predictive classifiers are provided segments in the
443 raw feature space or a latent space of reduced dimensionality using either PCA
444 or PLS. The only hyperparameter that is considered while generating the PCA
445 and PLS representations is the number of components for each method, i.e., the
446 dimensionality of the latent space. Exhaustive search is performed by performing
447 ten-fold cross-validation with LR classification while iteratively increasing the
448 number of components for PCA and PLS. Figure 16 shows the resulting validation
449 accuracy as the number of components increases. The peak validation accuracy
450 occurs at 16 components for PLS (i.e., PLS-16) and 41 components for PCA (i.e.,
451 PCA-41). Each of the five predictive classifiers is optimized and tested with data
452 from three representations, i.e., raw features, PCA-41, and PLS-16.

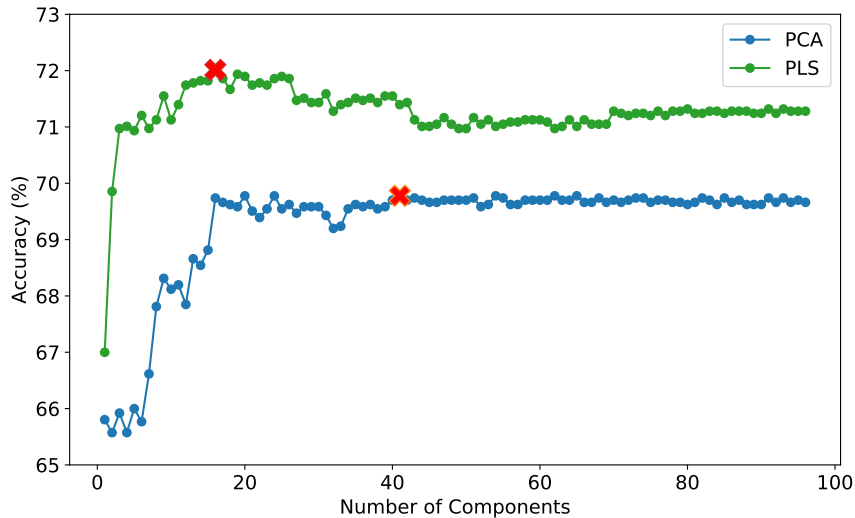


Figure 16: Selecting the number of components for PCA and PLS with cross-validation.

453 **4.2.2. Optimizing hyperparameters for predictive classifiers**

454 The Bayesian SMBO with a TPE is performed on the predictive classifiers,
455 i.e., the methods on the right hand side of Figure 10. The number of trials selected
456 for each method is manually selected based on the size of the hyperparameter

457 search space along with consideration for computational limitations. Optimizing
 458 the hyperparameters for these models, particularly for the deep learning models,
 459 is by far the most computationally demanding aspect of this work.

460 The classifiers are divided into traditional ML algorithms (i.e., LR, L-SVC,
 461 and K-SVC) and deep learning algorithms (i.e., ANNs and CNNs) for the purpose
 462 of describing the hyperparameter optimizations. The hyperparameter search space
 463 for the traditional ML algorithms and the deep learning algorithms are provided in
 464 Appendix A. Search options for the tolerance and the regularizer strength (λ) are
 465 the same for all traditional methods but some of the remaining hyperparameters
 466 only apply to one method (e.g., kernel type only applies to K-SVC). Notably,
 467 the penalty denoted elast. refers to elastic net, the squared-hinge (SH) loss is
 468 abbreviated, and the three kernels are abbreviated as poly. for polynomial basis,
 469 sig. for sigmoid basis, and the radial basis function (RBF). Thorough explanation
 470 of each hyperparameter is beyond our scope, inquisitive readers are referred to the
 471 literature for further information [61].

472 For each traditional ML classifier there are three explicit representations of the
 473 data that are separately optimized for hyperparameters. The result is nine exper-
 474 imental configurations of traditional ML algorithms with optimized hyperparam-
 475 eters specified in Table 1. Notably, an RBF kernel was selected for all K-SVCs,
 476 an SH loss was selected for all L-SVCs, and the optimized regularization strength
 varies significantly depending on the representation.

Table 1: Hyperparameter selection for traditional ML algorithms.

		λ	tolerance	penalty	loss	kernel
NA	LR	10	0.001	elast.	SH	RBF
	L-SVC	1000	0.001	ℓ_1		
	K-SVC	0.1	0.0001			
PCA-41	LR	10	0.001	elast.	SH	RBF
	L-SVC	100	0.00001	ℓ_1		
	K-SVC	0.01	0.0001			
PLS-16	LR	0.001	1e-5	ℓ_2	SH	RBF
	L-SVC	100	0.001	ℓ_1		
	K-SVC	0.1	1e-5			

477

478 For the deep learning methods the hyperparameter search space is significantly
 479 larger than the traditional ML methods. Although many more hyperparameter op-
 480 timization trials are conducted for the deep learning methods, the percentage of
 481 the search space covered by these trials is still significantly smaller. This is due
 482 to the fact that the search space for the deep learning methods is orders of magni-
 483 tude larger than the traditional ML methods and it is simply infeasible to conduct
 484 enough trials to search over an equivalent percentage of such a large space. Even
 485 with the use of cloud computation platforms to mitigate computational limita-
 486 tions, we are still constrained by the number of convolutional layers, batch size,
 487 and number of learned filters in our architecture.

488 After a series of challenging cross-validation trials the final choice of ANN
 489 and CNN hyperparameters are shown in Table 2 for each of the representations.
 490 The choice of optimizers, regularization strengths (λ), and fully connected layer
 491 (FCL) activation functions are the same for both ANN and CNN models. Some
 492 hyperparameters that are unique to the CNN include the number of convolutional
 layers (CLs) and the number of learnable filters.

Table 2: Hyperparameter selection for deep learning algorithms.

representation	NA		PCA-41		PLS-16	
classifier	ANN	CNN	ANN	CNN	ANN	CNN
optimizer	Adagrad	SGD	Adam	SGD	Adagrad	RMSprop
λ	0.1	0.05	0.01	0.001	0.1	0.0001
FCL activation	relu	elu	elu	tanh	elu	elu
no. of FCLs	5	2	9	1	10	1
FCL size	128	32	32	24	128	64
batch size	128	32	128	16	32	32
no. of CLs		1		1		1
CL activation		tanh		elu		tanh
filters		6		24		16
filter size		(5,5)		(20,20)		(3,3)
pool size		(1,1)		(2,2)		(2,2)

493 Ultimately, the fifteen experimental configurations in Figure 10 are outfitted
 494 with the parameters in Table 1 and Table 2. These models are tested on segments
 495 from two months of subsequent operation as shown in what follows.
 496

497 5. Experimental Results

498 Classification results from supervised learning studies can be represented as a
499 contingency table, known as a confusion matrix, with a dimension for the actual
500 class values (i.e., $y = 1$ or $y = 0$) and a dimension for predicted class values (i.e.,
501 $\hat{y} = 1$ or $\hat{y} = 0$). An arc loss event is considered a positive event with $y = 1$ and
502 no arc loss (i.e., stable operation) is considered a negative event with (i.e., $y = 0$).
503 The confusion matrix consists of four values; two of which correspond to correct
504 or truthful predictions, and two of which correspond to false predictions. False
505 predictions can be either false positive (FP) or false negative (FN), referring to
506 either type I error (false alarm) or type II error (missed alarm), respectively. True
507 predictions can be either true positive (TP) or true negative (TN), i.e., correctly
508 predicting arc loss or correctly predicting no arc loss, respectively.

509 For a particular experimental configuration (e.g., LR with PCA), the model
510 produces an output estimate for each segment in the testing set. Each output
511 estimate is compared to the true output label allowing the categorization of that
512 prediction as either FP, FN, TP, or TN. Therefore, the sum of these four values is
513 equivalent to the total number of segments in the testing dataset and the resulting
514 confusion matrix summarizes the prediction fidelity of the model with respect
515 to both the positive and the negative class. Various performance metrics (e.g.,
516 accuracy) can be derived for each model from the confusion matrix of that model.

517 The confusion matrix resulting from testing each of the fifteen experimental
518 configurations in Figure 10, with parameters shown in Table 1 and Table 2, is
519 provided in Table 3. In addition to the confusion matrix, two key performance
520 metrics are tabulated in Table 3 for each configuration, i.e., the accuracy (ACC)
521 and the recall, otherwise known as the true positive rate (TPR), with maximum
522 values emphasized in bold font. Accuracy is simply defined as the sum of true
523 predictions (i.e., TP and TN) divided by the sum of all predictions (i.e., the total
524 number of segments). The most accurate experimental configuration is with an
525 LR classifier on a 41 principal component representation followed very closely by
526 an LR classifier on the raw data itself.

527 The second critical performance metric provided in Table 3 is the recall which
528 focuses on the cases which precede an arc loss event. Specifically, recall is de-
529 fined as the number of times arc loss is correctly predicted divided by the number
530 of times arc loss occurs, i.e., $TPR = TP / (TP + FN)$. The experimental configura-

Table 3: Summary of the experimental results.

		TP	FP	TN	FN	ACC	TPR
NA	LR	182	49	151	76	0.727	0.705
	L-SVC	185	46	141	86	0.712	0.683
	K-SVC	167	64	156	71	0.705	0.702
	ANN	176	55	152	75	0.716	0.701
	CNN	181	50	141	86	0.703	0.678
PCA $d = 41$	LR	186	45	148	79	0.729	0.702
	L-SVC	176	55	151	76	0.714	0.698
	K-SVC	144	87	158	69	0.659	0.676
	ANN	176	55	142	85	0.694	0.674
	CNN	181	50	149	78	0.721	0.699
PLS $d = 16$	LR	184	47	130	97	0.686	0.655
	L-SVC	169	62	146	81	0.688	0.676
	K-SVC	166	65	149	78	0.688	0.680
	ANN	166	65	147	80	0.683	0.675
	CNN	147	84	161	66	0.672	0.690

531 tion with the best recall in this study is a logistic regression classifier on the raw
532 data. Interestingly, the runner-up for recall is a tie between a logistic regression
533 classifier on a 41 principal component representation and a kernel support vector
534 classifier on the raw data.

535 Aside from accuracy and recall, performance metrics for precision, also known
536 as positive predictive value (PPV), F_1 score, and F_β score are tabulated in Ap-
537 pendix A. The F_1 score represents the harmonic mean of precision and recall,
538 whereas the F_β score allows user specification of β which controls the weight-
539 ing of recall and precision, i.e., recall is β times more important than precision.
540 For this application recall is a much higher priority than precision because the
541 operating cost associated with false alarms is relatively negligible compared to
542 the operating cost associated with missed alarms (i.e., FN). Therefore, a choice
543 of $\beta = 0.25$ is selected for tabulating the F_β scores. The best configuration with
544 respect to precision and F_1 score is an LR classifier with a PCA representation,
545 whereas the configuration with the highest F_β score is, unsurprisingly, the same
546 as that with the highest recall, i.e., an LR classifier on the raw data.

547 Figure 17 provides a visually intuitive comparison of the classification accu-
548 racy results across the various experimental configurations. Comparing the differ-

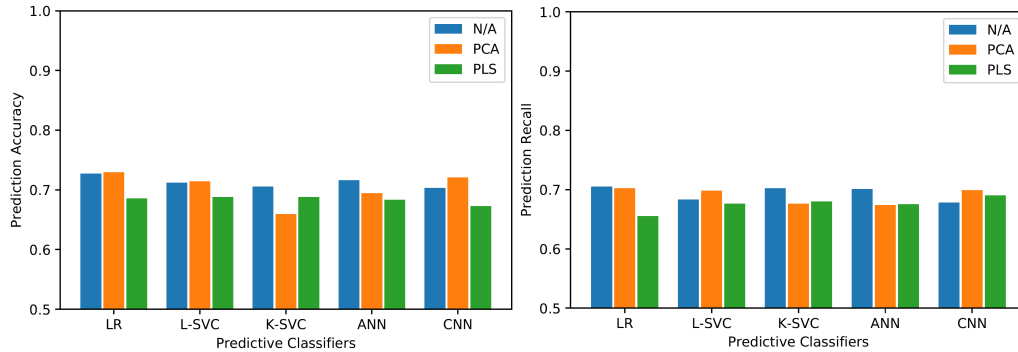


Figure 17: Comparing the testing accuracy (left) and recall (right) of each configuration.

548 ent representations in terms of accuracy, the raw data and the PCA representation
549 consistently outperform the PLS representation with the only exception being for
550 the K-SVC classifier. Comparing classifiers in terms of accuracy, the LR classifier
551 had the two highest accuracy scores with 72.9% and 72.7% on PCA components
552 and raw data, respectively. A similar situation arises with respect to the com-
553 parison of recall scores across the different representations, i.e., PLS is generally
554 outperformed by PCA and raw data representations with the raw data providing
555 the best recall score on average (across classifiers) as well as the highest recall of
556 70.5% with a LR classifier. Overall, logistic regression demonstrates better gener-
557 alization performance relative to the deep learning methods. Ultimately, given the
558 importance of recall in this application, a logistic regression classifier on the raw
559 data is the most promising configuration for development of an inferential sensor
560 to predict arc loss.
561

562 Deep learning methods contain a very large number of parameters which al-
563 lows them to model complex nonlinear functions if they have enough data to train
564 on. Although the logistic regression method performed slightly better in these ex-
565 periments, it is possible that the deep learning methods would perform better in
566 an experiment with multiple years of historical operating data. This is demon-
567 strated by Figure 18 which demonstrates the superior performance increases deep
568 learning methods relative to traditional methods when more data is provided. An-
569 other common issue with deep learning methods is over-fitting but special care
570 was taken to prevent over-fitting by introducing regularization and early stopping.

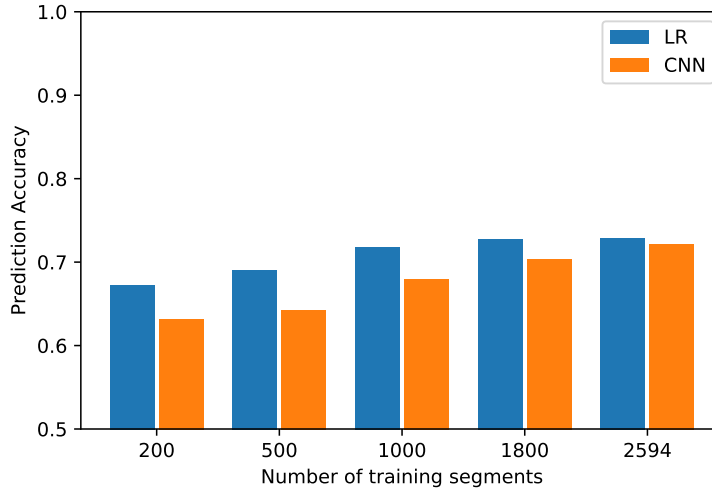


Figure 18: Comparing the classification accuracy of logistic regression with convolutional neural networks while varying the amount of training data.

571 Finally, sensitivity to hyper-parameters and network initialization is another po-
 572 tential concern for deep learning methods.

573 6. Conclusions and Future Directions

574 In this work we have introduced a novel industrial predictive classification
 575 problem, i.e., to predict arc loss in a DC EAF five minutes prior to the arc loss
 576 event. Moreover, we have proposed an end-to-end ML workflow that takes raw
 577 industrial data as input and yields an inferential sensor model that can predict
 578 arc loss events five minutes prior to occurrence with an accuracy of 72.7% and a
 579 true positive rate of 70.5% on unseen data from two subsequent months of oper-
 580 ation. Given that the unexpected loss of plasma arc in the DC EAF under study
 581 is an ongoing problem resulting in millions of dollars of lost production annu-
 582 ally, this work has the potential to contribute both significant economic savings
 583 and an improved environmental outcome as energy and material are consumed
 584 more efficiently. The final contribution of this work is the comprehensive em-
 585 pirical comparison between traditional and contemporary ML methods for rep-
 586 resentation learning and predictive classification during the development of the
 587 inferential sensor.

588 The use of representation learning algorithms for process data analytics is an
589 emerging research area that offers significant benefits to the process industry. This
590 work is a small part of a larger movement to migrate advanced data analytics tech-
591 niques from statistics and computing sciences to process industries. The explicit
592 representations considered in this work include traditional process analytics meth-
593 ods such as PCA and PLS, whereas implicit learned representations include the
594 multiple fully connected layers for the ANN and convolutional layers with multi-
595 ple learned filters for the CNN. Given all of the recent interest surrounding deep
596 learning methods we initially expected these methods to handily outperform the
597 traditional ML classifiers and representations. Having put forward a significantly
598 greater computational effort to optimize and train the deep learning methods re-
599 lative to the traditional ML classifiers it was indeed surprising to conclude that
600 applying the logistic regression classifier to the raw data was the best performing
601 experimental configuration.

602 Although the performance metrics in this study did not highlight the benefits
603 of representation learning, it is important to note that training and optimizing the
604 deep learning methods become much more computationally feasible when a lower
605 dimensional representation is used. For example, given a fixed computational
606 budget, using a CNN classifier with a lower dimensional representation enables
607 a much more comprehensive hyperparameter optimization, both in terms of trials
608 and search space, relative to using the raw data as the input to the CNN. Although
609 the PLS representations delivered slightly lower performance across various clas-
610 sifiers, the PLS representations also had the lowest dimensionality which led to
611 more convenient hyperparameter optimization. In future studies it is desirable to
612 implement more advanced variants of PCA and PLS as well as explicit nonlinear
613 deep representation learning methods such as deep autoencoders.

614 To further the development of the inferential sensor there are minor changes
615 to the experimental setup that can be made in future studies. For example, it is
616 desirable for future implementations to be tested on continuous raw data from the
617 process instead of relying on the procurement of a testing set composed of seg-
618 ments balanced by class. Such an experimental configuration would represent a
619 much more challenging evaluation of the inferential sensor, but it would also be a
620 significant step forward in terms of preparing for deployment. On the other hand,
621 one change that would likely improve performance while maintaining operational
622 integrity is to decrease the period over which the model is evaluated without being
623 updated. Two months is an unnecessarily large amount of time to have a model

624 be evaluated on new data without being updated. Instead, a more sensible config-
625 uration would be to evaluate the model over a shorter period (e.g., a week or less),
626 update the model weights with the new data from that period, and then continue
627 the evaluation with the updated model for the following week.

628 The arc loss predictor is designed to assist operators but ultimately it will be
629 up to operators to decide how it influence their actions. Operators may choose
630 to incorporate auxiliary information when deciding whether or not to act upon
631 the alarm. Ideally, the alarm will be accompanied by a confidence score which
632 could be designed using information about the quality of the model inputs and
633 prior knowledge of the process. Production trials will be critical for learning the
634 best practices for incorporating process knowledge into the model predictions.
635 Finally, ongoing work includes the procurement of a benchmark industrial arc
636 loss dataset that can be released to the process analytics community in order to
637 supplement the available simulation benchmarks such as the Tennessee Eastman
638 and penicillin fermentation datasets. Process analytics researchers could use this
639 dataset to evaluate their ML algorithms on a real industrial data with a process
640 fault that continues to have an unknown diagnosis.

641 **Acknowledgments**

642 The authors thank BBA Engineering Consultants, the National Science and
643 Engineering Research Council of Canada and the Izaak Walton Killam Memorial
644 Fund for funding this research through an Engage grant and a Killam Pre-Doctoral
645 Memorial Fellowship. This work is supported in part by the Institute for Comput-
646 ing, Information and Cognitive Systems (ICICS) at UBC.

647 **Appendix A. Supplemental Experiment Details**

648 The relatively small search space over which hyper-parameters are optimized
 for the traditional ML algorithms is presented in Table A.4.

Table A.4: Hyperparameter search space for traditional ML algorithms.

	LR	L-SVC	K-SVC
λ	(0.001, 0.01, 1, 10, 100)		
tolerance	(0.001, 0.0001, 0.00001)		
penalty	$(\ell_1, \ell_2, \text{elast.})$	(ℓ_1, ℓ_2)	
loss		(hinge, SH)	
kernel			(RBF, poly., sig.)
degree			(2, 3, 4, 5)

649

650 The vast search space over which hyper-parameters are optimized for the deep
 learning algorithms is presented in Table A.5.

Table A.5: Hyperparameter search space for deep learning algorithms.

	ANN	CNN
optimizer	(RMSprop, Adagrad, Adam, Adadelata, Adamax, SGD)	
λ	(0.1, 0.01, 0.001, 0.0001)	
FCL activation	(relu, tanh, selu, elu)	
no. of FCLs	(1, 2, ..., 10)	(1, 2, 3, 4)
FCL size	(32, 64, 128, 256, 512)	(32, 64, 128)
batch size	(32, 64, 128)	(16, 32, 64, 128)
epochs	(25, 35)	(20, 30, 40)
no. of CLs		(1, 2)
CL activation		(relu, tanh, selu, elu)
filters		(8, 16, 32, 64)
filter size		[(3,3), (5,5)]
pool size		[(2,2), (4,4)]

651

652 The precision (i.e., PPV), F_1 score and F_β ($\beta = 0.25$) score for each experi-
 653 mental configuration are provided as supplementary result metrics in Table A.6.

Table A.6: Supplemental experimental result metrics.

		TP	FP	TN	FN	PPV	F_1	F_β
NA	LR	182	49	151	76	0.788	0.744	0.708
	L-SVC	185	46	141	86	0.801	0.737	0.687
	K-SVC	167	64	156	71	0.723	0.712	0.702
	ANN	176	55	152	75	0.762	0.730	0.703
	CNN	181	50	141	86	0.784	0.727	0.681
PCA $d = 41$	LR	186	45	148	79	0.805	0.750	0.705
	L-SVC	176	55	151	76	0.762	0.729	0.701
	K-SVC	144	87	158	69	0.623	0.649	0.674
	ANN	176	55	142	85	0.762	0.715	0.677
	CNN	181	50	149	78	0.784	0.739	0.702
PLS $d = 16$	LR	184	47	130	97	0.797	0.719	0.659
	L-SVC	169	62	146	81	0.732	0.703	0.678
	K-SVC	166	65	149	78	0.719	0.699	0.682
	ANN	166	65	147	80	0.719	0.696	0.676
	CNN	147	84	161	66	0.636	0.662	0.688

654

References

- [1] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, MIT press, 2016.
- [2] Y. Bengio, A. Courville, P. Vincent, Representation learning: A review and new perspectives, IEEE transactions on pattern analysis and machine intelligence 35 (2013) 1798–1828.
- [3] N. Wagner, J. M. Rondinelli, Theory-guided machine learning in materials science, Frontiers in Materials 3 (2016) 28.
- [4] S. Heo, J. H. Lee, Fault detection and classification using artificial neural networks, IFAC-PapersOnLine 51 (2018) 470–475.
- [5] F. Lv, C. Wen, Z. Bao, M. Liu, Fault diagnosis based on deep learning, in: 2016 American Control Conference (ACC), IEEE, pp. 6851–6856.
- [6] D. Xie, L. Bai, A hierarchical deep neural network for fault diagnosis on tennessee-eastman process, in: 2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA), IEEE, pp. 745–748.
- [7] H. Wu, J. Zhao, Deep convolutional neural network model based chemical process fault diagnosis, Computers & Chemical Engineering 115 (2018) 185–197.
- [8] C. Ruiz-Cárcel, Y. Cao, D. Mba, L. Lao, R. Samuel, Statistical process monitoring of a multiphase flow facility, Control Engineering Practice 42 (2015) 74–88.
- [9] W. Yu, C. Zhao, Broad convolutional neural network based industrial process fault diagnosis with incremental learning capability, IEEE Transactions on Industrial Electronics 67 (2019) 5081–5091.
- [10] H. Legendijk, R. Jones, Production of ferronickel from nickel laterites in a dc arc furnace, Mintek.
- [11] D. Hurd, J. Kollar, Data for operating single electrode dc furnaces (1991).
- [12] I. Kotze, Pilot plant production of ferronickel from nickel oxide ores and dusts in a dc arc furnace, Minerals Engineering 15 (2002) 1017–1022.

- 683 [13] NRCan, Benchmarking energy intensity in the canadian steel
684 industry, [https://www.nrcan.gc.ca/sites/www.nrcan.gc.ca/files/oeefiles/pdf/industrial/](https://www.nrcan.gc.ca/sites/www.nrcan.gc.ca/files/oeefiles/pdf/industrial/SteelBenchmarkEnglish.pdf)
685 [SteelBenchmarkEnglish.pdf](https://www.nrcan.gc.ca/sites/www.nrcan.gc.ca/files/oeefiles/pdf/industrial/SteelBenchmarkEnglish.pdf), 2007.
686
- 687 [14] L. Rippon, I. Yousef, B. Hosseini, J.-F. Beaulieu, C. Prévost, S. Shah,
688 B. Gopaluni, Process analytics and machine learning to predict arc loss
689 in an electric arc furnace, in: 59th Conference of Metallurgists 2020 hosting
690 the 4th International Uranium Conference.
- 691 [15] S. Kotsiantis, D. Kanellopoulos, P. Pintelas, Data preprocessing for super-
692 vised leaning, *International Journal of Computer Science* 1 (2006) 111–117.
- 693 [16] L. Chiang, B. Lu, I. Castillo, Big data analytics in chemical engineering,
694 *Annual review of chemical and biomolecular engineering* 8 (2017) 63–85.
- 695 [17] E. Acuña, C. Rodriguez, On detection of outliers and their effect in super-
696 vised classification, *University of Puerto Rico at Mayaguez* 15 (2004).
- 697 [18] S. Visa, A. Ralescu, Issues in mining imbalanced data sets-a review paper,
698 in: *Proceedings of the sixteen midwest artificial intelligence and cognitive*
699 *science conference*, volume 2005, sn, pp. 67–73.
- 700 [19] B. Krawczyk, Learning from imbalanced data: open challenges and future
701 directions, *Progress in Artificial Intelligence* 5 (2016) 221–232.
- 702 [20] R. Longadge, S. Dongre, Class imbalance problem in data mining review,
703 *ArXiv abs/1305.1707* (2013).
- 704 [21] X. Guo, Y. Yin, C. Dong, G. Yang, G. Zhou, On the class imbalance problem,
705 in: *2008 Fourth international conference on natural computation*, volume 4,
706 *IEEE*, pp. 192–201.
- 707 [22] Z. Zheng, X. Wu, R. Srihari, Feature selection for text categorization on
708 imbalanced data, *ACM Sigkdd Explorations Newsletter* 6 (2004) 80–89.
- 709 [23] K. Pearson, Principal components analysis, *The London, Edinburgh, and*
710 *Dublin Philosophical Magazine and Journal of Science* 6 (1901) 559.
- 711 [24] H. Hotelling, *Multivariate quality control. techniques of statistical analysis*,
712 *McGraw-Hill, New York* (1947).

- 713 [25] V. Venkatasubramanian, R. Rengaswamy, S. N. Kavuri, A review of process
714 fault detection and diagnosis: Part ii: Qualitative models and search
715 strategies, *Computers & chemical engineering* 27 (2003) 313–326.
- 716 [26] N. F. Thornhill, S. L. Shah, B. Huang, A. Vishnubhotla, Spectral principal
717 component analysis of dynamic process data, *Control Engineering Practice*
718 10 (2002) 833–846.
- 719 [27] S. Narasimhan, S. L. Shah, Model identification and error covariance ma-
720 trix estimation from noisy data using pca, *Control Engineering Practice* 16
721 (2008) 146–155.
- 722 [28] H. Zhang, A. K. Tangirala, S. Shah, Dynamic process monitoring using
723 multiscale pca, in: *Engineering Solutions for the Next Millennium*. 1999
724 IEEE Canadian Conference on Electrical and Computer Engineering (Cat.
725 No. 99TH8411), volume 3, IEEE, pp. 1579–1584.
- 726 [29] E. L. Russell, L. H. Chiang, R. D. Braatz, Fault detection in industrial pro-
727 cesses using canonical variate analysis and dynamic principal component
728 analysis, *Chemometrics and intelligent laboratory systems* 51 (2000) 81–93.
- 729 [30] R. Sharmin, S. L. Shah, U. Sundararaj, A pca based fault detection scheme
730 for an industrial high pressure polyethylene reactor, *Macromolecular Reac-
731 tion Engineering* 2 (2008) 12–30.
- 732 [31] S. Skogestad, I. Postlethwaite, *Multivariable feedback control: analysis and
733 design*, volume 2, Wiley New York, 2007.
- 734 [32] S. Wold, M. Sjöström, L. Eriksson, Pls-regression: a basic tool of chemo-
735 metrics, *Chemometrics and intelligent laboratory systems* 58 (2001) 109–
736 130.
- 737 [33] G. Lee, T. Tosukhowong, J. H. Lee, C. Han, Fault diagnosis using the hybrid
738 method of signed digraph and partial least squares with time delay: The pulp
739 mill process, *Industrial & engineering chemistry research* 45 (2006) 9061–
740 9074.
- 741 [34] K. Severson, P. Chaiwatanodom, R. D. Braatz, Perspectives on process mon-
742 itoring of industrial systems, *Annual Reviews in Control* 42 (2016) 190–200.

- 743 [35] E. L. Russell, L. H. Chiang, R. D. Braatz, Data-driven methods for fault
744 detection and diagnosis in chemical processes, Springer Science & Business
745 Media, 2012.
- 746 [36] J. S. Cramer, The origins of logistic regression (2002).
- 747 [37] J. Shang, M. Chen, H. Ji, D. Zhou, H. Zhang, M. Li, Dominant trend based
748 logistic regression for fault diagnosis in nonstationary processes, *Control*
749 *Engineering Practice* 66 (2017) 156–168.
- 750 [38] T. W. Rauber, L. H. Mello, V. F. Rocha, F. M. Varejão, Multi-label fault clas-
751 sification experiments in a chemical process, in: 2014 Brazilian Conference
752 on Intelligent Systems, IEEE, pp. 265–270.
- 753 [39] C. Cortes, V. Vapnik, Support-vector networks, *Machine learning* 20 (1995)
754 273–297.
- 755 [40] V. N. Vapnik, A. Y. Chervonenkis, Necessary and sufficient conditions for
756 the uniform convergence of means to their expectations, *Theory of Probabi-*
757 *lity & Its Applications* 26 (1982) 532–553.
- 758 [41] J. Weston, C. Watkins, Multi-class support vector machines, Technical Re-
759 port, Citeseer, 1998.
- 760 [42] S. L. Brunton, J. N. Kutz, Data-driven science and engineering: Machine
761 learning, dynamical systems, and control, Cambridge University Press,
762 2019.
- 763 [43] Q. Lu, R. B. Gopaluni, M. G. Forbes, P. D. Loewen, J. U. Backström, G. A.
764 Dumont, Model-plant mismatch detection with support vector machines,
765 *IFAC-PapersOnLine* 50 (2017) 7993–7998.
- 766 [44] L. Rippon, Q. Lu, M. G. Forbes, B. Gopaluni, P. D. Loewen, J. Backstrom,
767 Machine direction adaptive control on a paper machine, *Industrial & Engi-*
768 *neering Chemistry Research* (2019).
- 769 [45] M. Onel, C. A. Kieslich, Y. A. Guzman, C. A. Floudas, E. N. Pistikopoulos,
770 Big data approach to batch process monitoring: Simultaneous fault detection
771 and diagnosis using nonlinear support vector machine-based feature selec-
772 tion, *Computers & chemical engineering* 115 (2018) 46–63.

- 773 [46] A. K. Jain, J. Mao, K. M. Mohiuddin, Artificial neural networks: A tutorial,
774 Computer 29 (1996) 31–44.
- 775 [47] J. V. Tu, Advantages and disadvantages of using artificial neural networks
776 versus logistic regression for predicting medical outcomes, Journal of Clin-
777 ical Epidemiology 49 (1996) 1225 – 1231.
- 778 [48] F. Rosenblatt, The perceptron: A probabilistic model for information storage
779 and organization in the brain, Psychological Review (1958) 65–386.
- 780 [49] L. Kanal, Perceptrons, in: N. J. Smelser, P. B. Baltes (Eds.), International
781 Encyclopedia of the Social & Behavioral Sciences, Pergamon, Oxford, 2001,
782 pp. 11218 – 11221.
- 783 [50] A. C. C. Coolen, A Beginner’s Guide to the Mathematics of Neural Net-
784 works, Springer London, London, pp. 13–70.
- 785 [51] C. Nwankpa, W. Ijomah, A. Gachagan, S. Marshall, Activation functions:
786 Comparison of trends in practice and research for deep learning, 2018.
- 787 [52] B. Ding, H. Qian, J. Zhou, Activation functions and their characteristics in
788 deep neural networks, in: 2018 Chinese Control And Decision Conference
789 (CCDC), IEEE, pp. 1836–1841.
- 790 [53] Y. Lecun, A theoretical framework for back-propagation (2001).
- 791 [54] Y. Le Cun, L. D. Jackel, B. Boser, J. S. Denker, H. P. Graf, I. Guyon, D. Hen-
792 derson, R. E. Howard, W. Hubbard, Handwritten digit recognition: Applica-
793 tions of neural network chips and automatic learning, IEEE Communications
794 Magazine 27 (1989) 41–46.
- 795 [55] Y. Le Cun, O. Matan, B. Boser, J. S. Denker, D. Henderson, R. E. Howard,
796 W. Hubbard, L. Jacket, H. S. Baird, Handwritten zip code recognition with
797 multilayer networks, in: [1990] Proceedings. 10th International Conference
798 on Pattern Recognition, volume 2, IEEE, pp. 35–40.
- 799 [56] Y. LeCun, Y. Bengio, et al., Convolutional networks for images, speech, and
800 time series, The handbook of brain theory and neural networks 3361 (1995)
801 1995.
- 802 [57] J. Bouvrie, Notes on convolutional neural networks (2006).

- 803 [58] K. B. Lee, S. Cheon, C. O. Kim, A convolutional neural network for
804 fault classification and diagnosis in semiconductor manufacturing processes,
805 IEEE Transactions on Semiconductor Manufacturing 30 (2017) 135–142.
- 806 [59] J. S. Bergstra, R. Bardenet, Y. Bengio, B. Kégl, Algorithms for hyper-
807 parameter optimization, in: Advances in neural information processing sys-
808 tems, pp. 2546–2554.
- 809 [60] J. Bergstra, D. Yamins, D. Cox, Making a science of model search: Hyper-
810 parameter optimization in hundreds of dimensions for vision architectures,
811 in: International conference on machine learning, pp. 115–123.
- 812 [61] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel,
813 M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,
814 D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine
815 learning in Python, Journal of Machine Learning Research 12 (2011) 2825–
816 2830.