
Towards Reliable Machine Learning Applications in Dynamic Manufacturing Environments

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The increasing deployment of advanced digital technologies such as Internet of
2 Things (IoT) devices and Cyber-Physical Systems (CPS) in industrial environments
3 is enabling the productive use of machine learning (ML) algorithms in the manu-
4 facturing domain. As ML applications transcend from research to productive use
5 in real-world industrial environments, the question of reliability arises. Since the
6 majority of ML models are trained and evaluated on static datasets, continuous
7 online monitoring of their performance is required to build reliable systems. Fur-
8 thermore, concept and sensor drift can lead to degrading accuracy of the algorithm
9 over time, thus compromising safety, acceptance and economics if undetected and
10 not properly addressed. In this work, we exemplarily highlight the severity of the
11 issue on a publicly available industrial dataset which was recorded over the course
12 of 36 months and explain possible sources of drift. We assess the robustness of ML
13 algorithms commonly used in manufacturing and show how uncertainty estimation
14 may be leveraged for online performance estimation as well as drift detection as a
15 first step towards continually learning applications.

16 1 Introduction and Motivation

17 Increasing digitization and the deployment of advanced technologies in the context of Internet of
18 Things (IoT) and Industry 4.0 are transforming manufacturing lines into Cyber-Physical Systems
19 (CPS) that generate large amounts of data. The availability of this data enables a multitude of
20 applications, including the development and deployment of ML algorithms for use cases such as
21 condition monitoring, predictive maintenance and quality prediction [1]. As ML applications are
22 deployed to productive usage, their continuous reliability has to be guaranteed to protect human
23 operators as well as the financial investments involved. Manufacturing environments are fast changing,
24 highly dynamic and inherently uncertain which poses the requirement for ML applications to be able
25 to adapt to changing environments with reasonable effort and cost [2]. While the ability of adapting
26 to a changing environment is often seen as a default property of machine learning [2], studies show,
27 that the generalization ability of a model mainly depends on the configuration and variety of the
28 available training data and is far from guaranteed [3, 4].

29 Long-term reliability and the handling of uncertainties caused by degrading equipment or faulty
30 sensors are seen as key factors and major hurdles when deploying ML systems in manufacturing
31 environments [5, 6]. With respect to safety certification and risk assessment, online performance
32 monitoring and uncertainty estimation are seen as critical [6]. In the context of quality management,
33 [7] showed, that the majority of the analyzed frameworks still lack any form of uncertainty estimation
34 or online monitoring.

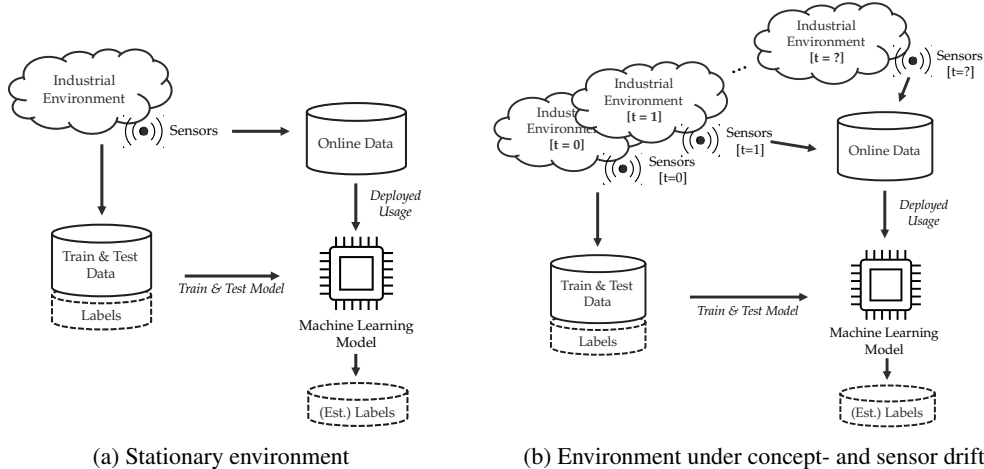


Figure 1: Supervised machine learning in stationary conditions with static sensors (a) in contrast to a dynamic system where the environment as well as the sensors used for perception are non-stationary, which applies to the majority of manufacturing usecases of ML (b). In the latter case, static training/testing sets do not provide continuous performance guarantees. Figure based on [9], adapted and extended with permission of the original authors.

35 In this work, we analyze the long-term reliability of ML applications in the manufacturing industry,
 36 highlighting the domain-specific issues and potential sources of drift. We benchmark a set of ML
 37 algorithms that are most relevant in this domain for robustness to time-dependent drift on an industrial
 38 dataset. Further, we assess uncertainty estimation techniques and highlight their potential utility for
 39 online performance estimation and drift detection in the context of continual learning to overcome
 40 the issue of silently failing ML applications.

41 Similar experiments have been conducted in [8] but did not consider non-deep learning algorithms,
 42 which are of high relevance in the manufacturing domain. Additionally, the introduced drifts were
 43 synthetic, while we evaluate on real-world time-dependent drifts.

44 2 Methodology and Experiments

45 *Concept drift* refers to a change in the underlying data distributions of machine learning applications.
 46 Especially in the context of pattern recognition, the terms *covariate shift* or *dataset shift* are used
 47 interchangeably [10]. Concept drift in the context of this publication, *cf.* Figure 1 (b), can be defined
 48 as $P_{train}(\mathbf{X}, Y) \neq P_{online,t}(\mathbf{X}, Y)$, where P_{train} and $P_{online,t}$ denote the joint distributions of
 49 input samples \mathbf{X} and target labels Y during training and deployed usage of the model at time t ,
 50 respectively.

51 Relevant short- and long-term sources of changes in manufacturing environments which may influence
 52 the reliability of ML models include: Tool- and machine wear [5, 6, 11, 12], changes in product
 53 configurations and material properties [11, 12], changes in upstream processes [12], changes in
 54 factory layout and machine placement [13], differences in operator preferences and -training [11, 12],
 55 seasons and time of day [6], environmental conditions such as temperature or humidity [6, 11, 12],
 56 sensor failure / -drift / -recalibration [5, 6] and data transmission problems [5].

57 Reliable machine learning applications in dynamic environments may be established in two ways:
 58 Either the model and data acquisition setup employed are robust against the relevant sources of drift,
 59 e.g. [14], or the model is continually assessed and, if required, adapted to the current environment in
 60 a continual learning setup [15].

61 Commonly, ML models with trained parameters θ in classification tasks produce probability estimates
 62 $p(\hat{y}_c | \mathbf{x}, \theta)$ for all classes $c \in \{1, \dots, C\}$, given a sample of data \mathbf{x} . The probability may be used to
 63 assess the models' confidence / uncertainty in its own prediction. The confidence is referred to as
 64 *well-calibrated*, when empirically, it is equal to the probability of the corresponding sample being

65 correctly classified [16]. Thus, confidence estimates that are well-calibrated even in the presence of
66 concept drift, may be used to reason about the reliability of a ML model and determine if it should be
67 adapted, i.e., retrained with new data. Additionally, well-calibrated confidence estimates may be used
68 to identify product configurations or situations that are difficult for the ML model to handle. In the
69 example of quality estimation, this could, e.g., trigger an additional human quality control for a given
70 part or the manual inspection of a machine.

71 2.1 ML Algorithms and Uncertainty Estimation Methods

72 To maximize the value for practical applications, we assess ML algorithms that have been identified
73 as most commonly used in the manufacturing domain by a recent review study [17]. We explicitly
74 include non-deep learning algorithms in the analysis, as these are highly relevant in the manufacturing
75 domain, where dataset sizes are often small. In the scope of this work, we focus on classification
76 tasks. Implementation is done using [18] and [19]. For each algorithm, we employ an uncertainty /
77 confidence estimation method as described below:

78 **Support Vector Machine (SVM)** confidence estimates are obtained using Platt scaling [20] of the
79 sample distances to the separating hyperplane. The parameters are fitted via 5-fold cross validation.

80 **Decision Tree (DT)** confidence estimates are computed as the fraction of training samples of the
81 same class in the leaf node [18].

82 **K-Nearest Neighbors (KNN)** confidence estimates are calculated similar to DTs. The probability of
83 a class is computed as the fraction of training samples of the same class in the set of nearest neighbors,
84 weighted by their distance.

85 **Random Forest (RF)** confidence estimates are computed as the mean predicted class probabilities of
86 the trees in the forest. The individual tree confidences are computed as described above. This method
87 of confidence estimation for RFs has been shown to be superior to more complicated extensions [21].

88 **Neural Network** For neural networks, we assess multiple recently proposed uncertainty estimation
89 methods: **Max. Softmax Probability (NN)** [22]; **Deep Ensembles (NN-Ens)** [23] with $M = 10$
90 ensemble members. Randomness is introduced by reshuffling of the training set as well as different
91 random initialization for each NN in the ensemble; **Monte-Carlo Dropout (NN-MCD)** [24] with
92 $M = 20$ forward passes for each sample. Dropout rate p is set to 0.2.

93 2.2 Dataset

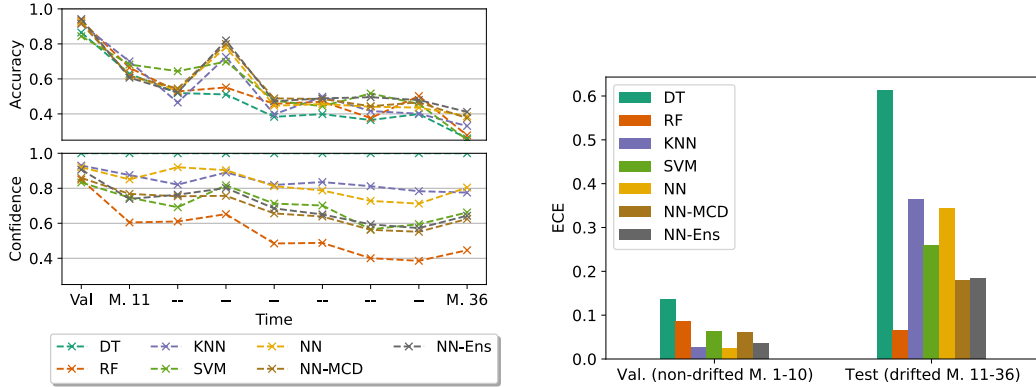
94 For our experiments, we use the Gas Sensor Array Drift dataset [25] that was recorded at the
95 University of California San Diego (UCSD). The dataset was recorded over 36 months at an industrial
96 test rig. Due to the long recording time, the dataset contains both sensor drift due to aging sensors
97 and concept drift due to external influences which resembles the expected environment conditions of
98 a real-world ML application in manufacturing, *cf.* Figure 1 (b). The dataset represents a classification
99 task, in which the target variable is the type of gas (one of six) that is currently present in the apparatus.
100 The experiments are perceived using 16 sensors and each row of the dataset contains 128 extracted
101 statistical features (8 per sensor) of the corresponding experiment run with a total of 13,910 runs.
102 The dataset is split into 10 consecutive batches, each capturing a varying amount of months.

103 2.3 Experiment Setting and Metrics

104 We train all the models on a random 50% split of the first 10 months of the available data and use the
105 remaining 50% as the validation set for performance evaluation on non-drifted data. To be able to
106 assess the robustness to drifts, we test on the remaining 26 months. For evaluation, we employ two
107 metrics capturing different aspects of interest:

108 **Accuracy** \uparrow is used to assess the performance of the model on the non-drifted test set as well as the
109 performance degradation under drift. The accuracy measures the percentage of correct classifications.

110 **Expected Calibration Error (ECE)** \downarrow [26] is used to evaluate the calibration of the confidences



(a) **Top:** Prediction accuracies over time. **Bottom:** Model confidences over time. *Val* indicates the validation set containing non-drifted data. (b) ECEs of the assessed models for the non-drifted validation set (months 1-10) as well as averaged over the drifted test set (months 11-36).

Figure 2: Results of the experiments on the UCSD Gas Sensor Array dataset. All experiments are repeated 10 times with different random seeds and the results consequently averaged.

111 produced by the model. The ECE corresponds to the average gap between model confidence and
 112 achieved accuracy. While the ECE has several shortcomings [8], we choose it over other calibration
 113 metrics for its simplicity and interpretability to strengthen the practical relevance.

114 2.4 Results

115 As visualized in the upper part of Figure 2 (a), the classification performance of all tested algorithms
 116 strongly degrades with an increasing time difference to the non-drifted validation set. This indicates,
 117 that none of the tested algorithms is robust against drifts in the environment. Thus, online monitoring
 118 and eventual model updates would be required to guarantee a reliable and safe application in real
 119 manufacturing environments. In parallel, the reduction in accuracy is well reflected by the lowering in
 120 confidence of a subset of the algorithms, most pronounced with RF, followed by NN-MCD, NN-Ens
 121 and SVM. Thus, the confidences may be used to identify drift in this scenario using frameworks such
 122 as [27]. The confidences are further analyzed in Figure 2 (b). Notably, most tested algorithms show
 123 well-calibrated confidences on the validation set, reflected by the low ECE, while the calibration
 124 strongly degrades for the drifted data. The calibration of the RF remains nearly unchanged, even
 125 for the drifted data. This supports the visualization (a) as the confidence degrades proportional to
 126 the accuracy under drift, indicating that the RF confidence may be used as a robust measure of the
 127 performance that can be expected of the algorithm as well as an indicator for drift.

128 3 Conclusion and Outlook

129 In this work, we highlighted the general relevance, implications and possible sources of drifts affecting
 130 the continuous reliability of ML applications in the manufacturing domain. Using an industrial dataset,
 131 we exemplarily show, that none of the most commonly used ML algorithms in manufacturing are
 132 robust against drifts in the data distribution inflicted by the environment or the sensors used for
 133 perception thereof. Positively, the experiments indicate that the confidence of algorithms such as
 134 random forests may be used to estimate the current performance and identify drifts. In a continual
 135 learning setup, the confidence could thus be used as a trigger signal for data collection and retraining
 136 of the respective model.

137 There are multiple opportunities for future work on continual learning and drift detection specific
 138 to ML usecases in manufacturing such as condition monitoring, predictive maintenance or quality
 139 prediction to enable further adaptation of ML applications in this domain. Especially the practical
 140 implementation of such systems on the shop floor level is still an open research issue.

141 References

- 142 [1] Nicolas Jourdan, Lukas Longard, Tobias Biegel, and Joachim Metternich. Machine Learning for
143 Intelligent Maintenance and Quality Control: A Review of Existing Datasets and Corresponding
144 Use Cases. volume 2, 2021.
- 145 [2] Thorsten Wuest, Daniel Weimer, Christopher Irgens, and Klaus-Dieter Thoben. Machine learn-
146 ing in manufacturing: advantages, challenges, and applications. *Production & Manufacturing*
147 *Research*, 4(1):23–45, 2016.
- 148 [3] Yeounoh Chung, Peter J Haas, Eli Upfal, and Tim Kraska. Unknown examples & machine
149 learning model generalization. *arXiv preprint arXiv:1808.08294*, 2018.
- 150 [4] Nicolas Jourdan, Eike Rehder, and Uwe Franke. Identification of uncertainty in artificial neural
151 networks. In *Proceedings of the 13th Uni-DAS eV Workshop Fahrerassistenz und automatisiertes*
152 *Fahren*, volume 2, 2020.
- 153 [5] Andrew Kusiak. Smart manufacturing must embrace big data. *Nature*, (544), 2017.
- 154 [6] Xinyang Wu, Mohamed El-Shamouty, and Philipp Wagner. White Paper: Dependable AI.
155 Using AI in Safety-Critical Industrial Applications. Technical report, Fraunhofer Institute For
156 Manufacturing Engineering and Automation (IPA), 2021.
- 157 [7] Beatriz Bretones Cassoli, Nicolas Jourdan, Phu H Nguyen, Sagar Sen, Enrique Garcia-Ceja,
158 and Joachim Metternich. Frameworks for data-driven quality management in cyber-physical
159 systems for manufacturing: A systematic review. *CIRP Conference on Intelligent Computation*
160 *in Manufacturing (ICME)*, 2021.
- 161 [8] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin,
162 Joshua V Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncer-
163 tainty? evaluating predictive uncertainty under dataset shift. *arXiv preprint arXiv:1906.02530*,
164 2019.
- 165 [9] Indrė Žliobaitė. Adaptive training set formation. 2010.
- 166 [10] Indrė Žliobaitė, Mykola Pechenizkiy, and Joao Gama. An overview of concept drift applications.
167 *Big data analysis: new algorithms for a new society*, pages 91–114, 2016.
- 168 [11] Shailesh Tripathi, David Muhr, Manuel Brunner, Herbert Jodlbauer, Matthias Dehmer, and Frank
169 Emmert-Streib. Ensuring the robustness and reliability of data-driven knowledge discovery
170 models in production and manufacturing. *Frontiers in Artificial Intelligence*, 4:22, 2021.
- 171 [12] Guangfan Gao, Heping Wu, Liangliang Sun, and Lixin He. Comprehensive quality evaluation
172 system for manufacturing enterprises of large piston compressors. *Procedia engineering*,
173 174:566–570, 2017.
- 174 [13] Fredy Sanz, Juan Ramírez, and Rosa Correa. Fuzzy inference systems applied to the analysis of
175 vibrations in electrical machines. *Fuzzy Inference Syst. Theory Appl*, 2012.
- 176 [14] Nicolas Jourdan, Tobias Biegel, Volker Knauth, Max von Buelow, Stefan Guthe, and Joachim
177 Metternich. A computer vision system for saw blade condition monitoring. *CIRP Conference*
178 *on Manufacturing Systems (CMS)*, 2021.
- 179 [15] Tom Diethe, Tom Borchert, Eno Thereska, Borja Balle, and Neil Lawrence. Continual learning
180 in practice. *arXiv preprint arXiv:1903.05202*, 2019.
- 181 [16] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
182 networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.

- 183 [17] Simon Fahle, Christopher Prinz, and Bernd Kuhlenkötter. Systematic review on machine learn-
184 ing (ml) methods for manufacturing processes—identifying artificial intelligence (ai) methods
185 for field application. *Procedia CIRP*, 93:413–418, 2020.
- 186 [18] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
187 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
188 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
189 *Learning Research*, 12:2825–2830, 2011.
- 190 [19] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
191 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow,
192 Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser,
193 Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek
194 Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal
195 Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete
196 Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-
197 scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- 198 [20] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized
199 likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- 200 [21] Henrik Bostrom. Estimating class probabilities in random forests. In *Sixth International*
201 *Conference on Machine Learning and Applications (ICMLA 2007)*, pages 211–216, 2007.
- 202 [22] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
203 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 204 [23] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable
205 predictive uncertainty estimation using deep ensembles. *arXiv preprint arXiv:1612.01474*,
206 2016.
- 207 [24] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
208 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.
209 PMLR, 2016.
- 210 [25] Alexander Vergara, Shankar Vembu, Tuba Ayhan, Margaret A Ryan, Margie L Homer, and
211 Ramón Huerta. Chemical gas sensor drift compensation using classifier ensembles. *Sensors*
212 *and Actuators B: Chemical*, 166:320–329, 2012.
- 213 [26] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated
214 probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*,
215 2015.
- 216 [27] Lucas Baier, Tim Schlör, Jakob Schöffner, and Niklas Kühl. Detecting concept drift with neural
217 network model uncertainty. *arXiv preprint arXiv:2107.01873*, 2021.