

Analyze, Sense, Preprocess, Predict, Implement, and Deploy (ASPPID): An Incremental Methodology based on Data Analytics for Cost-Efficiently Monitoring the Industry 4.0

Jesus Para^a, Javier Del Ser^{b,c,d,*}, Antonio J. Nebro^e,
Urko Zurutuza^a, and Francisco Herrera^f

^a*Mondragon University, 20500 Arrasate-Mondragon, Spain.*

^b*TECNALIA, 48160 Derio, Spain.*

^c*University of the Basque Country (UPV/EHU), 48013 Bilbao, Spain.*

^d*Basque Center for Applied Mathematics (BCAM), 48009 Bilbao, Spain.*

^e*University of Málaga, 29071 Málaga, Spain.*

^f*University of Granada, 18071 Granada, Spain.*

Abstract

Industry 4.0 is revolutionizing decision making processes within the manufacturing industry. Among the technological portfolio enabling this revolution, the late literature has been rich in what regards to the potential of data analytics for improving the production cycle at different stages, from resource provisioning to planning, delivery and storage. However, such a promising role of data analytics has been so far explored without a proper, quantitative inspection of the cost-improvement trade-off, nor has the process of acquiring sensors and extracting valuable information from their captured data formalized in a chain of methodological steps. This paper introduces the Analyze, Sense, Preprocess, Predict, Implement and Deploy (ASPPID) methodology, an iterative decisional workflow that spans from the acquisition of sensing equipment to the quantitative assessment of their contribution to enhance the production step under focus. This methodology aims at helping improvement teams to make informed decisions about which parts of the process needs to be sensed and how to exploit this information towards a verifiable improvement of the production cycle. The implementation of this methodology is exemplified in a real use case focused on the automotive in-

*Corresponding author: javier.delser@tecnalia.com (Prof. Javier Del Ser). TECNALIA. P. Tecnológico Bizkaia, Ed. 700, 48160 Derio, Spain. Tl: +34 946 430 50. Fax: +34 901 760 009.

dustry, where the detection of defects in an annealing process can be modeled as a classification problem with a highly imbalanced dataset. The results obtained after several iterations of proposed ASPPID methodology show that the scrap ratio is reduced by sensing the correct part of the process at reduced investment costs, thus highlighting the crucial role of data science in the business chain of manufacturing plants.

Keywords: Industry 4.0, Methodological Data Analytics, Process Monitoring, Cost Efficiency, Imbalanced Learning

1. Introduction

Since its first use during the Hannover Fair in 2014 [1], the so-called Industry 4.0 paradigm has drifted the interest of the community towards technologies aimed at providing intelligence to industrial process. Industry 4.0 is widely conceived as the fourth industrial revolution following the introduction of steam engines, chain production and process automation. As such, this concept refers to the design and development of *intelligent* factories [2], where traditional manufacturing processes are evolved by collecting large amounts of operational data coming from sensors and other information sources. This massive deployment of sensing equipment throughout the plant empowers multi-criteria decision making based on data analytics, with a target on increasing the efficiency of the manufacturing plants at different levels, from stocking to production and delivery. Not in vain this challenging paradigm has been recently identified as one of the main catalysts of the worldwide economical recovery after the last recession [3].

In order to remain competitive and increase the value chain of their productive processes, manufacturing industries have included the adaptation to this new productive model within their technology roadmap, embracing the inclusion of new technologies and competences within their processes and staff that had been traditionally linked to other sectors. The change has been particularly disruptive in regards to the monitoring of the whole production and delivery chain, with new technological enablers such as Internet of Things (IoT [4]), cybersecurity, augmented reality and collaborative robotics, among many others. However, there is common consensus around the crucial role taken by data analytics (either by itself or as a constituent, core part of Big Data) when providing the intelligence sought for the process at hand. Examples abound in the literature underpinning and evincing the gains yielded by the application of predictive models and optimization algorithms in different industrial problems, such as predictive maintenance or quality assessment (see e.g. [5, 6, 7] and references therein).

29 Unfortunately, the ever-growing need for integrating sensors and computing
30 technologies for data collection, storage and analysis undergoes several issues
31 in practice. Although manufacturing industries allocate big investments for in-
32 stallng advanced sensing equipment in their processes and machinery, purchases
33 are often biased by the potential of the equipment itself to capture and store as
34 many diverse parameters of the monitored asset as possible, disregarding their in-
35 herent utility and contribution for subsequent data modeling stages. While new
36 positions appear within the staff of the majority industries (such as experts in
37 robotics, data architects and data scientists [8]), it is often the case that these roles
38 do not get involved within the purchase decision cycle, thus bounding their con-
39 tribution to a *best-effort* over data captured by such devices once they have been
40 purchased. This lack of involvement of data-based roles in managerial matters im-
41 plies significant over-expenditures and a high degree of uncertainty with regards
42 to the potential of data analytics for the industry [9, 10].

43 This noted fact becomes even more concerning when jointly assessed with
44 the traditional client-provider relational model dominating most of the purchasing
45 models within the small/mid-sized industries. Cost implications of a certain pur-
46 chase must be evaluated not only from the purely economical perspective (namely,
47 the rated balance between investment and depreciation of the purchased equip-
48 ment, along with extra expenditures for ancillary tools and maintenance), but also
49 with regards to the *value* of the captured data, the late referring to the contribution
50 of the information captured by the purchased sensor to the data modeling task for
51 which it was acquired. Such a contribution, when quantified in terms of predictive
52 performance, can be translated to economical terms so as to reflect the competitive
53 worth of the investment. However, studies dealing with the cost efficiency of data
54 analysis within the industrial ecosystem are surprisingly scarce, despite the intu-
55 itively central role that this technology should play in manufacturing companies
56 towards their convergence to the Industry 4.0 paradigm [11].

57 In this paper we postulate that all mechanisms and decision processes aimed
58 at the provision of data-based intelligence to the manufacturing cycle should grav-
59 itate on the value of the data, quantification that can be only reliably undertaken
60 by staff close to the inference of knowledge from the captured information. As
61 will be exposed in the next subsection, most methodologies used to date in in-
62 dustrial decisional flows aim at improving the production efficiency of the plant
63 by addressing technical aspects that regard data as a byproduct of the process that
64 could be exploited if available, rather than considering this asset as a central matter
65 for decision. In the next section we provide an overview of such methodologies,
66 which helps contextualizing the contribution of this work.

67 *1.1. Related Work*

68 In order to face an ever-growing competitiveness in their respective market, in-
69 dustries have been traditionally forced to be constantly improving their processes.
70 Several methodologies can be adopted for this purpose, among which SPC (Sta-
71 tistical Process Control), TQM (Total Quality Management), LEAN and 6SIGMA
72 can be regarded as the most widely adopted ones:

- 73 • SPC [12] was the first methodology to embrace statistical control techniques in
74 the early XX century, used years later by the United States in the World War II
75 to improve the weapons manufacturing process. SPC relies on the use of control
76 graphs where important variables are monitored, to evidence changes in the nor-
77 mal behavior of production processes, allowing operators to act and correct the
78 drift of the processes. More recently insightful reviews have been contributed
79 to the community, such as [13], where the application and experiences of SPC
80 in the food industry is surveyed over 30 years. Most of these reports conclude
81 that the general motivation beneath the adoption of SPC is to reduce defects
82 production, but that some barriers withstand, such as the resistance to change
83 of the personnel involved and the lack of sufficient statistical knowledge. The
84 study in [14] focuses on quality improvement and cost effectiveness. Benefits
85 of the application of SPC include troubleshooting and diagnosis in shorter time
86 frames at reduced production costs. However, SPC undergoes practical issues
87 in its implementation, as operators must be trained beforehand and new docu-
88 mentation must be produced, requiring an workload excess for its preparation.

- 89 • TQM [15] is a management strategy aimed at creating quality awareness in all
90 organizational processes. TQM has been widely used in manufacturing, edu-
91 cation, government and service industries. It is referred to as *total* because it
92 is concerned with the organization of the company as a whole, including its
93 staff. TQM comprises three general paradigms: 1) management, with specific
94 compounding steps and goals such as planning, organization, control, and lead-
95 ership; 2) total, standing for the coverage of the methodology all throughout
96 the organization; and 3) quality, denoting the quantitateness of the methodol-
97 ogy by means of usual indicator definitions and all derivatives thereafter. The
98 community has reported several success histories in regards to the particular-
99 ization of TQM to several industrial sectors; to give two examples, in [16] the
100 TQM methodology is applied to a hotel in Greece, identifying as a result influ-
101 ential TQM factors such as the quality practices of the top management board,
102 strategic quality planning, employee quality management, customer focus and
103 employee knowledge and education. Likewise, [17] studies, through TQM, the
104 impact of Organizational Culture in small and medium Indian automotion com-

105 panies. Principal factors were discovered to be openness, confrontation, trust,
106 authenticity, proaction, autonomy, collaboration and experimentation.

107 • LEAN [18] originates directly from the Toyota production system developed
108 between 1950 and 1980. After producing a large batch of cylinder heads and
109 assembling them in the corresponding engines, Toyota's production engineers
110 finally found out that the required power was not achieved because the produced
111 cylinder heads were defective. The enormous costs caused by this problem
112 taught Toyoda one of the lessons that have been maintained in the company's
113 culture: the acquisition and verification of the quality of the components at each
114 stage of the production process, before moving on to the next stage, conforming
115 the so-called Jidoka concept. This was the first motivation for creating a new
116 methodology to eliminate waste, improve quality and reduce production times
117 and costs, based on a set of target operational areas: overproduction, dwell time,
118 transport, excessive procedures, inventory, movements and defects.

119 • 6SIGMA [19] is defined as a business process that allows companies to improve
120 their final results. Originally created by Motorola in the 80's and consolidated
121 by General Electric, 6SIGMA is a methodology that uses rigorously measured
122 and analyzed data to identify causes of a problem and ways to eliminate them,
123 generating greater customer satisfaction and substantial financial savings. It is
124 supported by statistical and analytical tools and proposes the development of
125 dynamic work groups, working with data in its search for the root cause of the
126 problem studied. This is accomplished by a set of 5 steps: Define, Measure, An-
127alyze, Improve and Control. Recent experiences with 6SIGMA abound: Shokri
128 *et al* in [20] reported a good practice of the application of 6SIGMA to reduce
129 defects in an airbag manufacturing plant, obtaining a scrap reduction on the
130 first year of 0.63%. Similarly, [21] applied 6SIGMA to reduce the defect in
131 the manufacture of glass containers for a cosmetics company, using to this end
132 tools such as Pareto, Process Map and FMEA (Failure Mode Effects Analysis).

133 The methodologies examined above are nowadays applied to a vast number
134 of diverse industries; however, the advent and proliferation of industrial sensors
135 capable to record high volumes of data spurred the derivation of new methodolo-
136 gies based on the analysis of collected information. Some examples can be found
137 in recent works such as [22], where real-time data is exploited to increase the ef-
138 ficiency of a production process; [23], which focuses on the use of data mining
139 techniques to real industry cases; or [24], which proposes a manufacturing data
140 mining methodology for industrial environments, which spans from the analysis
141 of apparel industries manufacturing unit to the implementation and evaluation of
142 the mining model.

143 *1.2. Contribution*

144 The motivation to introduce a new methodology is that all previous method-
145 ological approaches bring into focus data that are registered and then analyzed in
146 a sequential manner, regarding the sensor providing such data as an legacy equip-
147 ment for the data modeling and exploitation phase. The acquisition of sensors is
148 usually performed under a *max-specs-min-cost* criterion, without taking into con-
149 sideration the valuable knowledge that data analysis can shed over the suitability
150 of one sensor or another for the predictive task in question. Sensor costs can be
151 dramatically reduced not only by including roles related to data science within
152 the industry personnel, but also by dilating their influence to early stages of the
153 decision workflow, specially those where the selection of sensing equipment is
154 discussed. There resides indeed the novelty of the Analyze, Sense, Preprocess,
155 Predict, Implement, and Deploy (ASPPID) methodology proposed in this paper,
156 which can be broken down in several novel aspects as summarized below:

- 157 1. A formal mathematical statement of the addressed cost-efficiency problem,
158 which evinces the interplay between cost and benefit in data-intensive indus-
159 trial processes, and motivates the iterative nature of the proposed methodology.
- 160 2. A data-centered decision workflow, which postulates the data scientist as a
161 member of the core decision team in all decisions made around data-sensitive
162 industrial processes, and a thorough description of the steps involved at each
163 stage of the workflow.
- 164 3. An empirical validation of the proposed methodology when adopted to reduce
165 the number of defects in a real aluminum injection plant from the Basque
166 Country (Northern Spain). The lack of previous references in the literature
167 dealing with this brand new injection technology makes the scenario specially
168 suitable for ASPPID, given the high uncertainty around which sensor(s) to use
169 towards decreasing the scrap ratio and the complexity of learning from the
170 highly imbalanced dataset registered by the sensor(s). Evidences will be given
171 on the valuable contribution of the data scientist to the ASPPID methodology.

172 [ASSPID methodogoly has been implemented on a real case study at Edertek,](#)
173 [the R&D center of Fagor Ederlan, a tier1 automotive manufacturer, to predict and](#)
174 [reduce defective parts in their aluminum injection experimental line, reducing the](#)
175 [costs of sensor the process and getting value from data analysis of it, with the](#)
176 [particularity that the data response, quality part, is imbalance, with an imbalance](#)
177 [rate of 7.69.](#)

178 So, this paper is structured as follows: Section 2 formulates the optimization
179 problem that models the decision problem tackled by ASPPID, whereas Section

180 3 and subsections therein details the fundamentals of the proposed methodology.
181 The aforementioned real use case is presented in Section 4 along with a discussion
182 on the obtained results, and finally Section 5 concludes the paper and outlines
183 future research lines springing from this work.

184 2. Problem Formulation and Rationale

185 For a better understanding of the motivation for our proposed methodology we
186 herein formulate the decision workflow in which it is framed as an optimization
187 problem. This problem is driven by the balance between the costs derived from the
188 acquisition and installation of the sensing equipment, the collection and analysis
189 of the captured data and the benefits derived from the results of the analysis.

190 To formally pose this problem, we first mathematically model the cost incurred
191 by acquiring and installing a sensor equipment s from a portfolio of S possible
192 choices $\mathcal{S} \doteq \{s_1, \dots, s_S\}$ aimed at capturing data from a machine towards a
193 certain predictive modeling task (e.g. energy consumption characterization and
194 minimization, predictive maintenance or defect detection). Let the net price of
195 the sensor(s) be denoted as $C_{sens}(s)$ (in monetary units), to which installation,
196 calibration (both fixed quantities, $C_{cal}(s) + C_{inst}(s)$, including in-machine adjust-
197 ments and workforce) and maintenance (recurrently imputed cost $C_{mnt}(s)$ with
198 periodicity $T_{mnt}(s)$ [time units]) should be added for a fully operational sensing
199 deployment. Once installed, data must be retrieved from the machinery and stored
200 for its further analysis (either in the Manufacturing Execution Server – MES – or
201 the Cloud). Such a collection should span a minimum time frame $T_{col}^{min}(s)$ so as to
202 provide predictive models with enough data substrate for characterizing the pat-
203 tern between the collected information and the variable to be predicted (following
204 the previous examples, defect or consumed energy). This minimum time dura-
205 tion depends not only on the sampling specifications of the deployed sensors, but
206 also on the stationarity of the data and the timing requirements of the predictive
207 modeling task under consideration.

208 In any case, additional costs arise during the data collection period, such as the
209 storage of data $C_{stor}(s, t)$ (which depends on both the sensor – through their data
210 sampling specifications – and the time t for which it is utilized). Depending on the
211 requirements of the use case at hand it might be the case that the machine at hand
212 would be kept aside from the production line during the data capturing period for
213 experimentation purposes, causing further costs due to e.g. cease of production
214 and expenditures in feedstock and raw materials used to build samples for experi-
215 mentation. All such possible costs will be denoted as $C_{other}(s, t)$, which not only
216 depends on the sensor and time span of the experimentation themselves, but also
217 on the net contribution $R(s, t)$ (in monetary units per unit of time) of the halted

218 machine to the production of the plant. For instance, should a manufacturing
 219 chain in a parallel production chain be completely stopped and left inoperative
 220 during the trial, $R(s, t)$ would account for the aggregate difference between the
 221 net income of all products that would have been processed by the machine during
 222 $T_{col}^{min}(s)$ and the cost of raw material not acquired/consumed during the trial¹.

223 Once enough data have been retrieved from the sensed machinery, the work-
 224 ing flow proceeds by performing different data mining steps aimed at building
 225 a model based on the collected data. As has been previously argued this model
 226 construction phase aims at achieving a productive gain that can be translated to
 227 economical terms. For instance, a predictive model aimed at detecting defective
 228 products in a manufacturing machine could be evaluated in terms of its predic-
 229 tive score (e.g. ratio of true positives or recall). In this case, the recall of the
 230 developed predictive model could be combined with the scrap ratio of the ma-
 231 chine under analysis and the cost of scrap samples to yield a economical benefit
 232 $B(s, t)$ [monetary units] associated to the model ergo the acquisition of sensor s
 233 and measured over a certain time horizon t after the trial has been completed. The
 234 development of the model itself plus the prior preprocessing of the captured data
 235 would add an additional time gap $T_{model}(s)$ to the experimentation phase.

236 Before delving into further derivations, it is of paramount relevance for the
 237 scope of this study to note that the estimated economical benefit $B(s, t)$ cannot be
 238 quantified beforehand, as the predictive value for the information collected by sen-
 239 sor s is uncertain until captured and fed to the model for its development. This un-
 240 certainty should have profound practical implications when deciding which sensor
 241 to acquire, as well as the sales and purchase scheme. At this point we mathemati-
 242 cally pose two different purchase options:

- 243 A. Traditional sales and purchase scheme, by which the sensor supplier receives
 244 a monetary quantity $C_{sens}(s)$ for the selected sensor, plus additional payments
 245 $C_{cal}(s)$, $C_{inst}(s)$ and $C_{mnt}(s)$ should it also offer maintenance, installation and
 246 calibration services.
- 247 B. Sensor rental scheme, in which a time-dependent cost $C_{rent}(s, t) \ll C_{sens}(s)$ is
 248 negotiated and agreed between manufacturer and supplier, which is billed and
 249 paid during the total time of the trial, namely $T_{xp}(s) \doteq T_{col}^{min}(s) + T_{model}(s)$.

Intuitively, a rental scheme would incur a penalty cost $C_{rent}(s, T_{xp}(s))$, which
 adds to the final cost of the trial if the predictive value of the data captured by

¹This computation implicitly assumes that the sensed machine is not productive at all during the trial. Otherwise a mixed cost model should be instead adopted so as to reflect in $R(s, t)$ the interplay between material savings and productivity losses.

the sensor is good enough to make a difference in regards to the balance between overall costs and the benefit $B(s, t)$ of the developed model. Based on the above notation, the overall cost of the trial with sensor s when a traditional sales and purchase scheme (A in the above list) is adopted can be expressed as

$$C_A(s) = C_{sens}(s) + C_{cal}(s) + C_{inst}(s) + C_{mnt}(s) \left[\frac{T_{xp}(s)}{T_{mnt}(s)} \right] \quad (1)$$

$$+ C_{stor}(s, T_{xp}(s)) + C_{other}(s, T_{xp}(s)) + R(s, T_{xp}(s)), \quad (2)$$

where (1) refers all costs derived to the deployment of the sensor equipment, and (2) collects all expenditures and production losses (if any) during the data collection and model construction and assessment. In case the sales and purchase model is implemented on a rental basis (correspondingly, B in the list) during the time frame $T_{xp}(s)$ of the trial, the overall cost is given by

$$\begin{aligned} C_B(s, B_{th}(t)) = & \\ & C_{rent}(s, T_{xp}(s)) + C_{stor}(s, T_{xp}(s)) + C_{other}(s, T_{xp}(s)) + R(s, T_{xp}(s)) \\ & + \mathbb{I}(B(s, t) \geq B_{th}(t)) [C_{sens}(s) + C_{cal}(s) + C_{inst}(s) + C_{mnt}(s)], \quad (3) \end{aligned}$$

250 where $\mathbb{I}(\cdot)$ is an auxiliary binary function taking value 1 if its argument is true and
 251 0 otherwise. It can be noted that $C_A(s) \leq C_B(s)$ whenever the final decision is to
 252 acquire the sensor and exploit its information within the production line, decision
 253 that is made if the obtained benefit $B(s, t)$ from the data collection and modeling
 254 phases using sensor s exceeds an expected minimum threshold $B_{th}(t)$ [monetary
 255 units per unit of time] agreed by the team before starting the trial. We compile the
 256 above formulae as $C(s, m, B_{th}(t)) \in \{C_A(s), C_B(s)\}$, which denotes the cost of
 257 the trial with sensor s depending on its purchase mode $m \in \{A, B\}$.

258 As argued before, the unpredictability of the value of $B(s, t)$ for any t before
 259 acquiring sensor s imposes that it cannot be adopted as a metric for selecting a
 260 sensor from \mathcal{S} before its acquisition. Therefore, a prior analysis of the available
 261 choices for sensing is crucial, yet not always possible, specially in avant-garde
 262 manufacturing processes as the one later tackled in this manuscript. The lack
 263 of previously reported good sensing practices poses an unavoidable risk when
 264 selecting the sensor to be tested.

265 Based on the above rationale, methodologies as the one proposed in this pa-
 266 per should aim at maximizing the return of investment when a series of trials –
 267 each testing a sensor $s \in \mathcal{S}$ – is performed without any a priori knowledge nor
 268 estimation of the value of the captured data by each of the choices. Such trials
 269 are assumed to occur sequentially in time, such that $s_k \in \mathcal{S}$ denotes the sensor
 270 selected in trial k . At every trial k a two decisions must be made by the team in

271 regards to 1) the acquisition model m_k of selected sensor s_k (A or B, if available
272 upon agreement with the supplier); and 2) a target value $B_{th}(t)$ of the estimated
273 economical benefit of the developed model based on the retrieved data by the
274 purchased sensors over the series of trials. It should be clear that if sensor s_k is
275 selected at trial k , purchased ($m_k = A$) and deployed on the machine/process at
276 hand, it might occur that the captured data and the modeling phase do not achieve
277 the aforementioned minimum value, i.e. $B(s_k, t) < B_{th}(t) \forall t$. However, since
278 sensor s_k was purchased, sensors selected at subsequent trials may benefit from
279 the information captured by s_k , potentially – yet uncertainly – furnishing a better
280 model performance than if sensor s_{k+1} was utilized in isolation. On the contrary,
281 benefit yielded by the model constructed from s_k could undergo a collinearity
282 between variables collected at different trials, leading to a lack of predictive rele-
283 vance of the information provided by the new sensor. This implicitly unveils that
284 the choice of sensors is order-sensitive.

By introducing a slight abuse in notation, the accumulated benefits after K successive trials can be computed as $B(\mathbf{s}_K, \mathbf{m}_K, t)$, where $\mathbf{s}_K = [s_1, s_2, \dots, s_K]$ and $\mathbf{m}_K = [m_1, m_2, \dots, m_K]$ (with $m_k \in \{A, B\}$) denote the selected sensors at trials $\{1, \dots, K\}$ and their purchase mode, respectively. A similar redefinition of the overall costs of the trials can be made as

$$C(\mathbf{s}_K, \mathbf{m}_K, B_{th}(t)) = \sum_{k=1}^K C(s_k, m_k, B_{th}(t)), \quad (4)$$

which must be kept below an upper bound C_{max} that represents the budget allocated by the manufacturing plant to perform trials. A purchase and trial methodology should pursue the optimal number K^* and sequence \mathbf{s}_K^* of selected sensors and their purchasing policy \mathbf{m}_K^* so that the sought benefit from the developed model $B_{th}(t)$ is achieved at a minimum overall cost. Mathematically,

$$[K^*, \mathbf{s}_K^*, \mathbf{m}_K^*] = \arg \min_{K, \mathbf{s}_K, \mathbf{m}_K} C(\mathbf{s}_K, \mathbf{m}_K, B_{th}(t)), \quad (5)$$

$$\text{subject to } C(\mathbf{s}_k, \mathbf{m}_k, B_{th}(t)) \leq C_{max}, \forall k \in \{1, \dots, K\}, \quad (6)$$

$$B(\mathbf{s}_K, \mathbf{m}_K, t) \geq B_{th}(t), \quad (7)$$

$$m_k \in \{A, B\} \cap \mathcal{M}_{sup}(s), \quad (8)$$

$$s_k \neq s_{k'} \text{ if } k \neq k' \forall k, k' \in \{1, \dots, K\}, \quad (9)$$

285 where (6) denotes that the overall cost is constrained, (8) indicates that the choice
286 of purchase mode should also take into account the selling policy $\mathcal{M}_{sup}(s) \in$
287 $\{A, \{A, B\}, B\}$ of the sensor supplier; and (7) imposes that the target benefit
288 should be achieved after the K trials. The process to implement a hypothetical case

289 comprising $K = 3$ trials is exemplified in Figure 1, where it is important to notice
 290 the impacts of the decisions made after every trial in the cost and the accumulated
 benefit along the process.

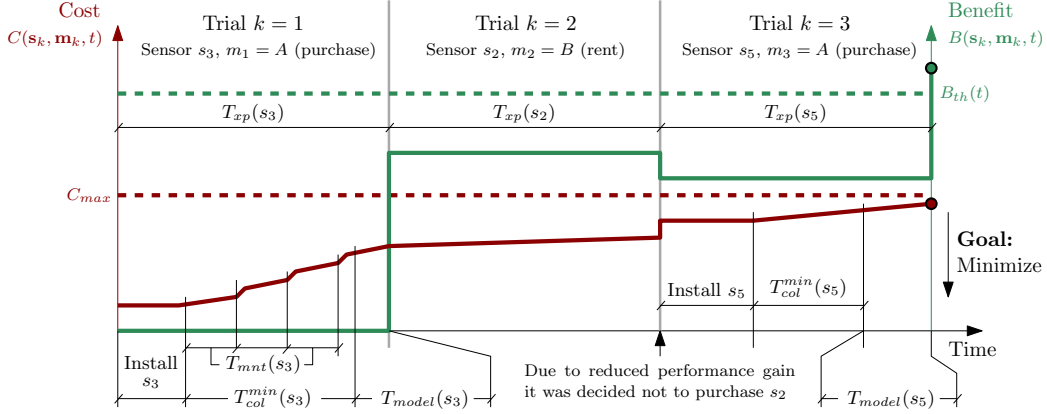


Figure 1: Example of $K = 3$ sensor trials driven by a methodology proposing $\mathbf{s}_K = [s_3, s_2, s_5]$ and $\mathbf{m}_K = [A, B, A]$ for a given predictive modeling task (e.g. defect detection). A double vertical axis has been used to reflect benefits and costs. Due to the uncertainty of the team over the utility of s_2 for the task at hand it was decided to rent it for the duration of the trial, saving costs after its lack of predictive relevance was confirmed.

291
 292 It should be clear that decisions made on m_k along the trials should be made
 293 not only depending on the purchasing options allowed by the supplier, but also
 294 on the prior knowledge (if any) about the potential of the sensor to provide value
 295 for achieving the goal $B_{th}(t)$. Should this a priori estimation be reliable enough,
 296 purchase mode A would incur less costs than B , but in practice this entails a risk
 297 of very difficult clearance. The ASPPID methodology proposed in this work is
 298 aimed at lowering this risk by including the data scientist all through the decision
 299 loop, not only in the data modeling stage. The adoption of our proposed method-
 300 ology not only reduces the aforementioned risk, but also establishes realistically
 301 the expected benefits from the data model by matching a business-related goal
 302 $B_{th}(t)$ to data-based insights within an iterative process.

303 3. Proposed ASPPID Methodology

304 In order to face the problem posed in the previous section, a methodology
 305 is proposed to make the entire decision process gravitate on cost and data-based
 306 performance of the considered sensing equipment. The methodology is detailed
 307 in what follows.

308 **3.1. General Description**

309 ASPPID is an iterative methodology based on data analysis that pursues a cost-
 310 efficient improvement of a certain productive process within manufacturing plants
 311 by involving a data scientist from the very beginning of the methodology, i.e. from
 312 the moment when economical decisions are made on the acquisition model of sensor
 313 equipment to the moment at which the developed model is deployed on the
 314 production line. Traditionally, data scientists are part of the improvement team
 315 when sensors are already installed on the machines, being in charge of exploit-
 316 ing the information collected by them. In this way, data scientists are restricted to
 317 mine the information that the installed sensors are recording. In practice this strat-
 318 egy leads to very frequent errors: already installed sensor(s) may not be recording
 319 relevant variables and parameters for the targeted mission of the model, which
 320 ultimately jeopardizes the discovery of sufficiently explanatory models to achieve
 321 the imposed benefit goal $B_{th}(t)$. As anticipated before, the reason for this mis-
 322 match between the information and the performance of the sensor is that in data
 323 science, the predictive value of data not only has to do with its cardinality, but
 324 also with its relevance: it is often the case when staff responsible for acquiring in-
 325 dustrial equipment hinge their decision on the power of sensors and price, sinking
 326 in excess of information, incorporating redundant sensors or sensors that register
 variables with no influence on the predictive model, leading to additional costs.



Figure 2: Schematic diagram illustrating the iterative workflow of the ASPPID methodology.

327 In ASPPID the data scientist plays a central role within the organization, partic-
 328 ipating not only in the development of advanced models for data analytics, but
 329

330 also in the choice and validation of acquired sensors, after in-depth analysis of
331 previous practices reported in the literature and taking into account the cost of the
332 sensors $C(s, m, B_{th}(s))$.

333 3.2. ASPPID Steps

334 Specifically ASPPID can be split into 6 fundamental steps: 5 comprising an
335 iterative cycle aimed at fulfilling the project goal(s); and a final step performed
336 only once after the target has been met and the developed model validated. Such
337 steps are depicted in Figures 2 and described in the following subsections:

338 3.2.1. Analyze

339 Fueled by the momentum gained by the Industry 4.0 paradigm, it is usual
340 to find organizations machinery with a myriad of sensors capable of recording
341 huge amounts of data captured from different stages of the productive process.
342 In practice this information is almost never exploited and, if used anytime, much
343 of the captured information is highly redundant or does not add any value to the
344 production efficiency of the plant. This is a real problem for organizations: they
345 have large collections of data, but they do not know how to exploit them nor have
346 they analyzed the cost overrun as the result of purchasing decisions.

347 ASPPID faces this problem by defining a methodology particularly suitable
348 for monitoring industrial machines and processes with sensors. An in-depth anal-
349 ysis of the target problem or task that the organization aims to solve should be
350 the first step before sensing. If this analysis is not properly done or should it
351 steer towards mistaken goals, the entire methodology could give rise to consid-
352 erably higher costs, hence this step is crucial in our proposed methodology. In
353 this regard, the entire set of decisions to be made has to be driven by a concern,
354 by something to improve quantitatively. This target must be established by the
355 plant/business management board, but has to be achievable, measurable and veri-
356 fiable. Once the goal has been set, a so-called *improvement* team must be arranged
357 to carry out all the trials needed to meet the target.

358 3.2.1.1. *Contract*. ASPPID dictates that trials cannot be started unless the im-
359 provement team agrees with the terms and goals imposed by the management
360 board. This agreement is formalized in a *contract*, which reflects and certifies
361 the confidence and support of the management of the plant or business to which
362 the improvement project is destined, feeling it as its own. This contract compels
363 both parts to fulfill the responsibilities posed in its clauses: on one hand, the man-
364 agement board will use it as tool to trace and verify the results obtained by the
365 improvement team within the timing plan; on the other hand, the improvement
366 team will be allocated to an internal project having agreed in their terms, which

367 is a guarantee of involvement, commitment and identification with its goals, far
368 beyond the usual implications of a project (e.g. work labor, company support for
369 investment). The evolution of all the improvement teams on the plant with respect
370 to their signed contracts will be analyzed at the end of the working year to detect
371 common causes of deviation, unveil management problems in the organization
372 and potential needs of prevailing improvement teams. The contract should have
373 *at least* the following elements:

- 374 • Title: it must be descriptive, concise and easy to understand for both client (in
375 ASPPID this role is taken by the management board) and improvement team.
- 376 • Description of the problem for which the improvement team is hired. This part
377 must be filled in by the customer and reviewed by the improvement team.
- 378 • Objective of the project $B_{th}(t)$, which must be described in detail and quantified
379 in monetary terms as an expected target benefit of the model to be developed
380 within the project. It is of utmost importance to specify the units in which
381 the benefit is measured, either as a rate over time (as in the problem statement
382 in Section 2) or as an alternative measure, e.g. benefit per product (such as
383 monetary units per produced unit). In all cases this goal will be used to decide
384 the success/failure of the project and its closure. It is important to remark that
385 this objective can be reviewed along the iterative ASPPID workflow, yet any
386 change must be reviewed and agreed by both parties.
- 387 • Maximum costs of the project C_{max} , which establishes an upper bound on the
388 economic costs associated to the implementation of ASPPID including work-
389 force, sensor purchases, ancillary equipment and external services. Since costs
390 can be revised between trials and for transparency and traceability it is advisable
391 to break down the maximum cost of the trials into the aforementioned concepts.
- 392 • Team: the definition of the improvement team should not only determine their
393 members, but also their hours of dedication. Team leadership will be shared by
394 the project manager and the leading data scientist. In constant communication
395 to each other, the project manager is responsible for managing the infrastruc-
396 ture part, purchasing equipment and contacting sensor suppliers, whereas the
397 leading data scientist will be responsible for data analysis. The rest of the im-
398 provement team must be composed by personnel aware of the process from a
399 theoretical/technical point of view, as well as by maintenance engineering per-
400 sonnel and/or data architects. The role of them is crucial to determined which
401 parts of the process they suspect to be, under their experience, most promis-
402 ingly originating the problem to be solved. Data architects and/or maintenance
403 engineers will ensure that sensors are installed and recording data properly.

- 404 • Timing plan, including estimated duration of the planned trials $T_{xp}(s) \forall s \in \mathcal{S}$,
 405 intermediate milestones and reserved work hours for each member of the im-
 406 provement team (see Figure 3 for an exemplifying schedule of the ASPPID
 407 phases for the shown in Figure 1). Creating realistic timing plans helps assess-
 408 ing the progress of the project, and keeps the team motivated. It is important
 409 that the customer knows the real status of the project, so presentations will
 410 be performed periodically (at an agreed period) and once every trial k is over.
 411 When detected, any punctual deviation from the original planning must be thor-
 412 oughly documented and recorded, and countermeasures and corrective actions
 413 must be proposed to the client in duly time, i.e. in the presentation held closest
 414 to the detection of the deviation. The client and the project leader will evaluate
 415 the cause of the deviation, resolving it (if possible).

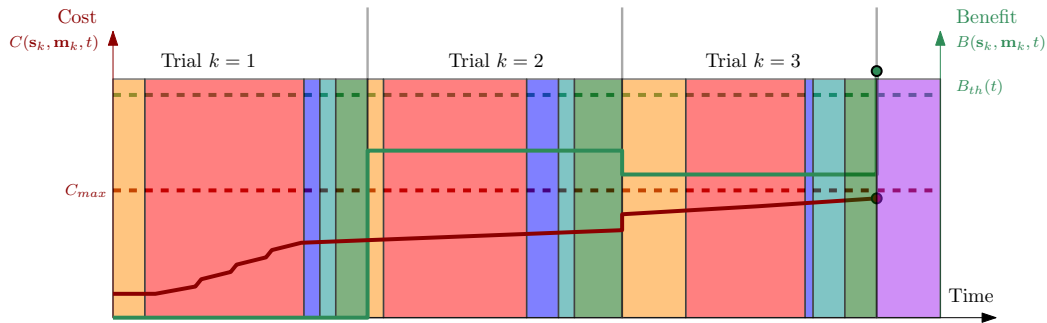


Figure 3: Indicative duration of each phase of ASPPID in the example of Figure 1.

416 After the contract has been completed, it must be signed by the parts: the
 417 customer (plant/business directors) and all team members. A copy of the signed
 418 contract must be sent to all managers of the team members, so that they have a
 419 record of the hours that their collaborators will devote to the ASPPID project.
 420 Once the project has been signed, a first team meeting must be held, where team
 421 leaders explain in detail to the rest of the team why the team has been formed, and
 422 the role that all the members will play on the team. The objectives, timing and
 423 dedication must be clear to all attendants.

424 *3.2.1.2. Brainstorming.* Once this step has been taken, the team can begin to ap-
 425 proach the project. The first step is to do a brainstorming exercise to determine
 426 which causes may be affecting the problem. A brainstorming exercise is a group
 427 creativity technique, where all participants must freely contribute ideas about the
 428 causes that affect a given problem. It is very important for this dynamic to be
 429 effective, that all participants feel free from expressing the ideas that they think
 430 can have an influence. This is why team leaders must communicate the rules of

431 participation to the rest of the team. No opinion can be underestimated nor any
432 person held back. The ranks of the organization must be forgotten. If the presence
433 of a responsible can hinder the correct development of the dynamics, then this
434 exercise should be carried out at different times, isolating the conflicting parts.

435 To efficiently implement this brainstorming session, ASPPID embraces the
436 so-called Ishikawa diagram [25] – also known as the fish-bone diagram because
437 of its shape – which consists of an axis on which the problem to be tackled re-
438 sides. From this axis comes the *thorns*, grouping together the main aspects related
439 to the problem: equipment, process, people, materials, environment and manage-
440 ment, as depicted in Figure 4. The objective of the Ishikawa procedure is that the
441 team considers all possible causes that may affect the problem at hand, not only
442 those in which each individual member probably focuses biased by his/her role in
443 the plant. When all ideas are contributed, those that are believed to possess the
444 highest influence in the problem are selected. Hypotheses will be established with
445 them. Other variables deemed irrelevant / less relevant are initially discarded, but
446 can be taken again into account in forthcoming trials: due to the uncertainty of
447 the selection of sensors and variables, what seemed to be the root cause and a
448 relevant variable for the task/problem at hand does not often meet the expecta-
449 tions. Therefore, new hypotheses must be again generated by bringing back the
Ishikawa methodology along with all previously discarded alternatives.

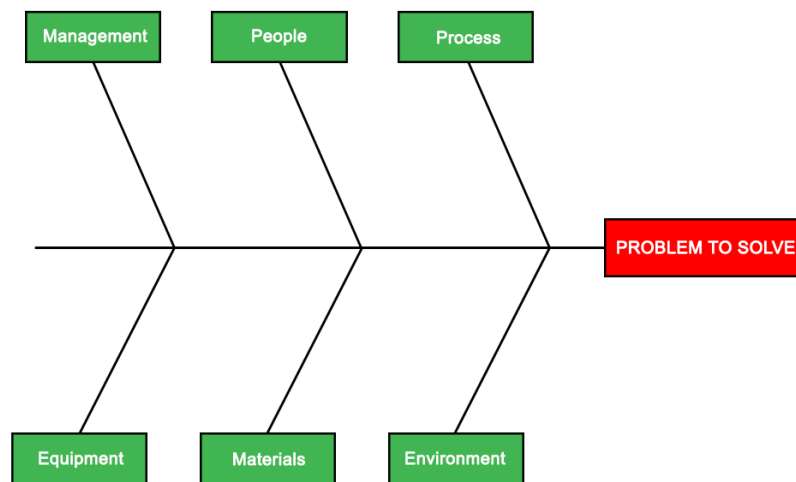


Figure 4: Ishikawa diagram utilized within the *analyze* phase of ASPPID.

450 At this point we highlight that the participation of the data scientist from the
451 very beginning of the process (including the establishment and agreement of the
452 project's goals) is unprecedented, and sheds certainty on the reachability of the
453

454 benefits expected by the rest of the team on the model to be developed. To this
455 end the data scientist can contribute to the brainstorming session with literature on
456 the data-based problem to be solved, drafty estimations based on reported good
457 practices of similar setups, and an explanation of the practical implications of
458 the scores by which the performance of the model to be developed. As a result
459 the brainstorming session should produce a calendar of K trials where to test
460 different monitored variables, and a qualitative yet informed indication of the level
461 of confidence with respect to their utility for the addressed problem.

462 3.2.2. *Sense*

463 Once the brainstorming exercise has been completed, it is time to test the set of
464 hypothesized variables at every trial. To do this, the project manager and the lead-
465 ing data scientist must inspect whether the variables are already being recorded
466 within the plant, and if so, whether the sensor capturing such information is reli-
467 able, properly calibrated, in optimum operating conditions [26, 27] and capturing
468 enough data for the needs of the subsequent modeling phase. To verify the latter,
469 the sensor supplier must carry out tests and a calibration plan to avoiding opera-
470 tional errors, repeatability and bias problems.

471 If the information is not being recorded, the team must gauge the acquisition of
472 sensing equipment aimed at capturing the required variables. At this point ASP-
473 PID underpins the importance of performing this acquisition in a cost-efficient
474 way, since purchasing sensors based only on their specifications can be counter-
475 productive due to the unpredictability of the contribution of its captured data to the
476 problem at hand. When the leading data scientist declares the utility of s_k at trial
477 k as extremely uncertain (due to the lack of previous history with the sensor or
478 the problem under analysis), the cost derivation in Section 2 suggests that renting
479 the sensing equipment should be preferred over a traditional purchase, if possible
480 (i.e. if the sensor supplier allows for this model in the acquisition of s_k , i.e. if
481 $B \in \mathcal{M}_{supp}(s_k)$).

482 Sensors can be a major part of the cost of the project, so acquiring them only
483 after validating if they are estimated to be truly effective is critical to reduce the
484 overall cost $C(s_K, \mathbf{m}_K, B_{th}(t))$ of the project. As described in Section 2, when-
485 ever there is a risk of lack of predictive power of the captured data, sensor renting
486 is more cost-efficient than purchasing it from the very beginning of the trial. This
487 is the reason why ASPPID imposes that before purchasing a sensor it must be
488 proven reliably that its captured data can be expected to yield a quantitative gain
489 towards achieving the project's goal. This a priori estimation must be done by the
490 leading data scientist (and his/her team, if any) based on a prior literature study,
491 interviews with plant personnel, and the study of reports with similar cases.

492 Once the process/machine has been sensed, it is time to collect and analyze

493 information of the process. To do this, the team must establish the minimum time
494 period $T_{col}^{min}(s_k)$ for the sensor s_k to collect representative data of all the variability
495 and casuistry of the monitored parameters of the current trial. This estimation can
496 be performed by exploring several aspects of the scenario under analysis, such as
497 the elements that cause the variability of the response variable to be predicted by
498 the model, the variability pattern (e.g. seasonality) of the magnitudes from which
499 data is retrieved, and other statistical features of the data alike.

500 3.2.3. Preprocess

501 Once data have been captured and stored, ASPPID proceeds by preprocess-
502 ing it towards the construction of a model aimed at achieving the project's goal.
503 This is an essential step in the data analysis process, as the extraction of valuable
504 information from data not only depends on the models themselves, but also on a
505 complete data preprocessing process including the following steps [28]:

- 506 • Cleansing consists of the removal of data observations that significantly differ
507 from the rest of observations, prone to misfit the subsequent model.
- 508 • Data scaling, by which each recorded variable is scaled prior to the use of mod-
509 els that are sensitive to the range of values taken by their input variables (e.g.
510 neural networks). A survey of scaling techniques can be found in [29].
- 511 • Imputation, by which data samples with missing values are removed or filled
512 depending on their degree of incompleteness [30, 31].
- 513 • Noise identification, by which all instances that result in a loss of predictive
514 model accuracy are removed and not used for model building [32].
- 515 • Data transformation, which combines, projects and/or summarizes the recorded
516 data in order to make samples fulfill a priori assumptions of the consequent
517 predictive model (e.g. linearity, normality).

518 Although the relevance of data preprocessing has been widely acknowledged
519 by the research community, ASPPID reflects explicitly the need for undertaking
520 this step within its cycle due to its high share of the modeling costs $T_{model}(s_k)$.
521 When agreeing on the schedule and costs of the project, the management team
522 must be made aware of the specific data preprocessing needs of the considered
523 sensor portfolio and be warned of eventual over-costs due to this step. To mini-
524 mize the potential budgetary impact of data preprocessing an informed data qual-
525 ity estimation report should be prepared by the team, if possible before contract
526 negotiation and agreement.

527 *3.2.4. Predict*

528 After data preprocessing, the next step is to generate models that permit to
529 achieve the benefit goal $B_{th}(t)$ established in the contract agreement and, ulti-
530 mately, to enhance the performance of the monitored machine/process within the
531 manufacturing plant. To this end the retrieved, already preprocessed dataset is
532 used as the starting point of a data analysis phase. As per the majority of prac-
533 tical use cases foreseen for ASPPID, this phase will aim at the development of
534 a *predictive* model capable of explaining the pattern between a set of monitored
535 parameters and a target variable, the latter impacting on $B_{th}(t)$ either directly or
536 indirectly (e.g. consumed power or defective output).

537 All in all, an exploratory data analysis must be first done [33], comprising
538 techniques such as descriptive statistics, dimensionality reduction and clustering,
539 among others. Once the exploratory analysis is done, the part of the improvement
540 team in charge for data analysis and modeling will be in a better position to con-
541 struct the planned data model(s); in this regard, the nature of the variable(s) to
542 be predicted and the problem at hand will require different procedures for con-
543 structing the model. Intuitively, in industrial environments productive scenarios
544 embracing ASPPID for improving their performance are based on sensors, whose
545 collected information is used to train and build a predictive model for classifi-
546 cation or regression, depending on the kind of target variables to be predicted.
547 In industrial environments this family of scenarios often requires the adoption of
548 techniques for class imbalance, feature selection/construction and model valida-
549 tion, since it is likely that the collected data lacks of enough positive examples for
550 properly modeling the pattern between the monitored parameters and the target
551 variable(s). This is the reason why filter, wrapper or embedded feature selection
552 methods become essential in the *predict* phase of ASPPID, along with oversam-
553 pling, under-sampling or cost-sensitive techniques for class balancing.

554 Methodologically speaking, during this phase the improvement team must
555 bear in time the *interpretability* of the model towards easing its deployment should
556 it perform satisfactorily with respect to the imposed target benefit $B_{th}(t)$. To this
557 end, ASPPID fosters the consideration of rule mining techniques so as to infer in-
558 terpretable rule sets from the developed models that can be implemented straight-
559 forward by the personnel of the plant upon its approval by the management team.

560 *3.2.5. Implement*

561 The *implement* phase is aimed at verifying that models generated on the previ-
562 ous phase are valid on production and therefore, can be deployed if they meet the
563 imposed performance goal and eventually approved. Depending on the scenario
564 and problem in hands, the implementation can imply 1) that the developed model
565 is tested with new unseen data to assess its estimated generalization capabilities

566 (as in e.g. a predictive model forecasting the power consumption of a certain ma-
567 chine); 2) and/or that the parameters configuring the monitored machine/process
568 are varied to those optimized based on the insights provided by the data modeling
569 phase. In both cases this phase must prove that models generated in the previous
570 phase perform as expected during real production, so that an empirical guarantee
571 of the improvement yielded by the developed model can be granted to the man-
572 agement board.

573 To this end, an implementation period must be determined. It should be long
574 enough to allow checking the impact of all process components with a potential in-
575 fluence on the target variable (e.g. changes of tools, molds, etc), but short enough
576 for the equipment tension to be maintained and costs reduced. It is necessary
577 to record the start and end timestamps of this phase so as to compute the benefit
578 $B(\mathbf{s}_k, \mathbf{m}_k, t)$ from the savings obtained during this time frame. Once this period is
579 over, conclusions from the implementation phase must be fully documented, along
580 with an estimation of the economical impact of the tested model, its comparison to
581 the established goal $B_{th}(t)$ and, if deemed necessary, informed recommendations
582 from the improvement team (specially, from the data scientist(s)) on an update
583 of the goals and targets of the ASPPID project. This reported information will
584 be gathered within a deliverable for the client's perusal, who will be called for a
585 review meeting so as to close the k -th loop of the ASPPID project.

586 3.2.6. *Deploy*

587 In the aforementioned meeting a decision on the fulfillment of the project's ob-
588 jective is to be made in light of the reported information and conclusions from the
589 current ASPPID iteration and the evolution of costs incurred within the project.
590 Should the target $B_{th}(t)$ has been met by the developed model, a cost study of
591 the extra cost required to deploy the model definitely within the production line
592 is elaborated, including the purchase of sensors if they have been rented during
593 previous trials. If the estimated overall costs for deployment are contained be-
594 low C_{max} , the project is declared as successful. Otherwise, both the team and
595 the management team must agree on either a relaxation of the cost constraint or a
596 reversal of the ASPPID cycle to a previous trial, which is decided depending on
597 the suitability of each alternative for the plant. For instance, if the gap between
598 $B(\mathbf{s}_k, \mathbf{m}_k, t)$ and $B_{th}(t)$ is high, the management board could decide to accom-
599 modate extra funding for deploying the model; conversely, expenditures beyond
600 C_{max} should be leveraged by assessing small benefit over $B_{th}(t)$ with the invest-
601 ment required to deploy the model.

602 A platform to keep documents compiling the knowledge acquired during the
603 project is advisable, since many of the obtained data-based insights can be extrap-
604 olated to other scenarios dealing with products with equal or similar character-

605 istics. In these documents the best possible detail should be given as to identify
606 the key conditions and aspects that make the obtained conclusions generalizable
607 to other use cases. The responsible for recording this information is the project
608 manager. Once the information is saved in the platform, a project closing meeting
609 should be organized to crosscheck that the produced documentation is complete.
610 It is also important that the rest of the organization visualizes its success, so a
611 good practice is to install visible panels inside the plant with the most synoptic
612 description possible of the steps taken along the ASPPID cycle.

613 **4. Real Case Study of the Proposed Methodology**

614 ASPPID was originally conceived and put to practice in the manufacturing
615 industrial case of the Basque Country (northern Spain). This section elaborates
616 on the details of its implementation, placing an emphasis on how the process
617 evolved towards getting a cost-efficient operational gain in the production of the
618 manufacturing line under study.

619 *4.1. The Scenario: Edertek (Fagor Ederlan)*

620 In particular the project was run over Edertek, a corporate research and devel-
621 opment center of Fagor Ederlan. Fagor Ederlan is a supplier for the automotive in-
622 dustry specialized in chassis and power train applications in different technologies
623 and materials. Nowadays Fagor Ederlan has production plants in Spain, Mexico,
624 China and Slovakia, a worldwide presence for which the company is particularly
625 in pursuing technological advances for production efficiency – through its center
626 Edertek – that can be extrapolated and exploited in several albeit related scenar-
627 ios. For confidentiality reasons, sensitive information will not be disclosed in the
628 conclusions extracted from the application of ASPPID to this real use case, yet we
629 consider this information of minimal relevance to illustrate how ASPPID helped
630 this company.

631 The case of use aims at the prediction of defective parts in an aluminum injec-
632 tion line reproduced experimentally in Edertek. From a data analysis perspective,
633 this predictive task can be modeled as a highly-imbalanced binary classification
634 problem where the model to be developed classifies products as ACCEPTED (i.e.
635 compliant with quality requirements) or REJECTED. Although the literature is
636 rich in tackling this particular problem instance by means of predictive models,
637 the aluminum injection is a brand new manufacturing technology for which no
638 previous references are known. This particularity motivated further the participa-
639 tion of the leading data scientist of the company all over the decision workflow,
640 proposing deep changes in the *modus operandi* that prevailed in the company be-
641 fore this use case was started. Those changes gave rise to the definition of an

642 early iterative procedure to relocate data analysis at the core of the project, which
643 matured towards the herein proposed ASPPID methodology.

644 4.2. *Data mining Methodology and Metrics used on implementation*

645 The objective of the ASPPID project was agreed to be the design and deploy-
646 ment of a process robust enough to manufacture parts with a sufficiently high
647 quality rate for this first experimental phase.

648 As described on previous subsection, this project has been run on the experi-
649 mental line of Edertek, a R&D center, to get a predictive analysis on the quality
650 parts, and reduce the costs of sensor. This kind of process has usually a high
651 imbalance rate. In this case study, the IR is 7.68.

652 To deal with imbalance domain it is necessary to use imbalance data-mining
653 techniques to make imbalance class distribution has no influence on model gen-
654 eration. To do this, over and under-sampling techniques has been tested: minor-
655 ity class oversampling, Synthetic Minority Over-sampling Technique (SMOTE,
656 [36]), SMOTE with Tomek links removal [37], the borderline variant of SMOTE
657 [38], Adaptive Synthetic Sampling (ADASYN, [39]), Edited Nearest Neighbours
658 (ENN, [40]) and One-Sided Selection (OSS, [41]).

659 The most common metrics to validate classification problems is Acc (10). This
660 metric is not appropriate on imbalanced domains, since it is not capable to distin-
661 guish between the number of correctly classified instances on the majority and
662 minority class [34].

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} \quad (10)$$

663 On Edertek all the manufactured parts pass an internal quality control, so met-
664 rics used were F1 score (13) and principally, Recall (12). Precision (11) is by-
665 passed because of the fact that all parts pass a quality control, so the focus of
666 the model is to be capable to create good models for majority class (ACCEPTED
667 parts) and minority one (REJECTED parts).

$$Precision = \frac{TP}{TP + FP} \quad (11)$$

$$Recall = \frac{TP}{TP + FN} \quad (12)$$

$$F1 = \frac{(1 + \beta)Recall * Preccision}{\beta^2Recall + Precission} \quad (13)$$

668 AUC 14 was also calculated to demonstrate is not a good metric, as is influ-
 669 enced by the majority class.

$$AUC = \frac{1 + TPrate - FPrate}{2} \quad (14)$$

670 Following the ASPPID guidelines for data analysis in Subsection (3.2.4) a
 671 Random Forest (RF, [43]) classifier with 500 trees was used as the inner predic-
 672 tive model for choosing among class imbalance techniques, and to discover the
 673 features that had more influence on the response. Once the most influent vari-
 674 ables are selected, a C4.5 Decision Tree model was utilized as the interpretable
 675 predictive model due to its simplicity and straightforward readability in terms of
 676 deployable rule sets.

677 Both, the over & undersampling thecniques + Random Forest and C4.5 tree,
 678 were runned in a 5×5 -fold cross validation to avoid bias on modelling.

679 For costs controlling on sensors, the mathematical approach detailed on Sec-
 680 tion 2 is used.

681 The way on that the process was sensed and data analyzed are described on
 682 the next subsection, deploying ASPPID methodology.

683 4.3. Description of ASPPID Applied to the Use Case

684 In this case $K = 2$ iterations of ASPPID sufficed to achieve the imposed goal:

685 4.3.1. First Iteration ($k = 1$)

686 (1.a) Analyze: Once the target was defined quantitatively (translated from mone-
 687 tary units to a maximum Non-Quality Rate, NQR, equal to 8% that should
 688 be achieved by the project to meet the target) and after a budget was provi-
 689 sioned for the project, a multidisciplinary improvement team was arranged,
 690 with a project manager and a leading data scientist as its coordinators. The
 691 rest of the team was selected based on the technical knowledge of the pro-
 692 cess to be monitored, and a timing was established to develop the project.
 693 Once the contract was created and signed, a first team meeting was launched
 694 aimed at developing an Ishikawa exercise to conjecture which parameters
 695 may have an influence on the scrap ration of this new injection technology.
 696 To do this, a blackboard was used and the data scientist, leading the session,

697 explained the rules and coordinated the brainstorming through an Ishikawa
698 diagram. As a result, a selection and prioritization of hypothesized factors
699 were agreed by the team. For this first iteration, the following factors were
700 selected:

- 701 • Furnace number that processes each injected part.
- 702 • Cycle time, i.e. time during which the part is processed in the furnace.
- 703 • Cavity number of the part within the mould.
- 704 • Mould number within the stock of moulds in the plant.
- 705 • X_t^1 [kept confidential], a time series representing a magnitude recorded
706 during the furnace processing.

707 (1.b) Sense: once identified, the above factors must be sensed from the moni-
708 tored process. To do this and given the high uncertainty with this injection
709 technology, one of the principles of the ASPPID methodology, to spend the
710 minimum amount of money on sensors, was adopted, avoiding to purchase
711 equipment that would increase the cost of the project without any guaranteed
712 reliability of their contribution to the project's goal, as described in Section
713 2. In this first iteration the team got advantage of sensing equipment that
714 had previously been installed in another process within the plan and that
715 was out of use by the time when the trial was made. Should this equipment
716 be eventually relevant for the project, new sensing units would be purchased
717 and installed in the monitored process. A period $T_{col}^{min}(s_1)$ equal to 2 months
718 was decided in order to balance the trade-off between losses due to the non-
719 productivity of the monitored process and the collection of enough data to
720 capture the underlying variability of the process. A detailed cost breakdown
721 was done to jointly evaluate all decisions impacting on the budget of the
722 project (use of provisional sensors, losses during the *sense* period, etc.).

723 (1.c) Preprocess: once collected (at a rate of 1 sample per second in the case
724 of X_t^1), data were preprocessed using one-hot encoding for the categorical
725 variables (furnace, cavity and mould numbers), whereas the cycle time was
726 converted to its Z-score. As for the time-dependent variable X_t^1 , a series
727 of interviews with empirical findings from the personnel allowed extracting
728 features of potential importance for defect detection from the time series
729 recorded for each part, as shown² in Figure 5. Once extracted, the fea-
730 ture values were each normalized to their Z-score. The preprocessing stage

²Plotted curves must be considered just for illustrative purposes, since for the sake of confi-
dentiality the originally processed time series cannot be explicitly shown in the Figure

yielded a total of 19 binary and real-valued predictors.

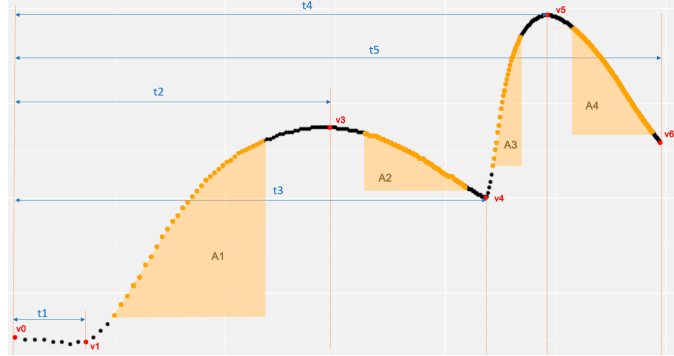


Figure 5: Features extracted from the time series X_t^1 captured by the selected sensor at trial $k = 1$: extreme points $\{v_0, \dots, v_6\}$, areas $\{A_1, \dots, A_4\}$ and delays $\{t_1, \dots, t_5\}$, resulting in 16 features.

731

732

733

734

735

736

737

No imputation of lost data was needed in this first cycle, whereas outliers were only found in X_t^1 due to reading errors undergone by the installed sensor. A statistical correction was performed over identified atypical samples, substituting by the average of the median of the 3 previous and subsequent values those points whose difference with respect to the previous and posterior point value is greater than a certain percentage. (see Figures 6.a and 6.b).



Figure 6: (a) Computation of medians after an outlier \bullet has been detected; (b) Replacement of the atypical sample using the average value of the computed medians.

738

739 (1.d) Predict:

740

741

742

743

744

Results in terms of precision, recall, F1 score and Area Under the Curve (AUC, [42]) obtained by each scheme in a 5×5 -fold cross-validation for under and over-sampling + RF are shown in Table 1. Given the poor results obtained in this benchmark, a final check was performed by the team to determine whether the number of examples in the dataset was enough for

Table 1: Average model scores of the first iteration of the ASPPID methodology.

| Model | Recall | Precision | F1 | AUC |
|-----------------------|--------|-----------|--------|--------|
| OVERSAMPLING + RF | 0.2325 | 0.6427 | 0.3404 | 0.8745 |
| SMOTE + RF | 0.2988 | 0.4527 | 0.3480 | 0.8214 |
| BORDERLINESMOTE + RF | 0.2325 | 0.4510 | 0.3051 | 0.8633 |
| SMOTE-TOMEKLINKS + RF | 0.313 | 0.4346 | 0.3504 | 0.8105 |
| ADASYN + RF | 0.2331 | 0.5020 | 0.3160 | 0.8666 |
| ENN + RF | 0.3199 | 0.4201 | 0.3463 | 0.8032 |
| OSS + RF | 0.2164 | 0.6584 | 0.3249 | 0.8837 |

745 predicting the target variable. This study was done by observing the evolu-
 746 tion of the models’ performance as the size of the dataset is progressively
 747 reduced. Only shown for ADASYN + RF (Figure 7) for brevity, results
 748 obtained for all methods lead to the same conclusion: the amount of data
 749 retrieved in this first trial was sufficient, since additional data produced by
 750 more samples (by e.g. extending $T_{col}^{min}(s_1)$) would not improve significantly
 751 the performance of the model.

752 As mentioned on the previous subsection 4.2, we can confirm that AUC is
 753 not a good measure when imbalance class is present. It can be seen that
 754 despite the model is not well classifying the positive class, high values of
 755 AUC are given. This is because of the weight that majority class has on
 756 model generation.

757 (1.e) Implement: the unsatisfactory results obtained with this first sensor were not
 758 convincing to invest on a *implement* phase to validate them in production.

759 Based on the data-based insights drawn from this first study, the improvement
 760 team proposed a new iteration of the ASPPID methodology to the management
 761 board of the plant, where to test a new set of sensors within the monitored pro-
 762 duction process. The first conclusion drawn and agreed by both parties was that
 763 sensors used in this first trial were not enough to meet the established target value
 764 $B_{th}(t)$, yet the use of sensors already installed in another machine instead of pur-
 765 chasing them had reduced dramatically – more than 90% – the planned costs for
 766 this first iteration. This fact buttressed the consideration of cost efficiency as a key
 767 aspect in subsequent iterations of ASPPID, and steered thereafter the negotiations
 768 with sensor suppliers.

769 4.3.2. Second Iteration ($k = 2$)

770 This second iteration, and in general, all iterations within an ASPPID imple-
 771 mentation, starts always from the Ishikawa exercise done in previous iterations:

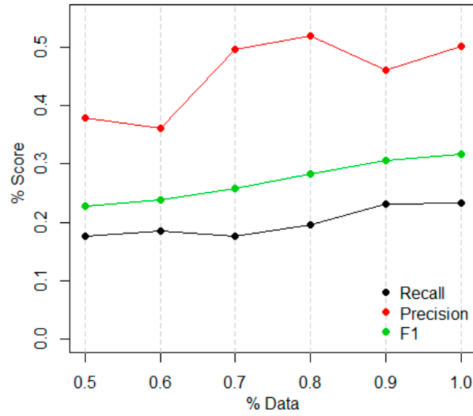


Figure 7: Average model scores versus amount of data for the ADASYN + RF approach.

- 772 (2.a) Analyze: the improvement team gathered again and generated new hypothe-
773 ses with the variables that were initially discarded. Consequently it was
774 decided and agreed to sense new magnitudes and elements along the pro-
775 duction line, withstanding the uncertainty of their predictive relevance with
776 the incorporation of new members (upon approval of the management team)
777 with knowledge in Physics and Thermodynamics. Since sensors in the first
778 ASPPID iteration yielded no costs for the project, they were again used in
779 this second trial. New sensors needed to capture the newly decided set of
780 parameters were rented instead of purchased after negotiation with the sup-
781 plier, who finally admitted a change in their relationship model with the
782 company (i.e. the supplier accepted to include mode B in its sales portfolio
783 $\mathcal{M}_{sup}(s_2)$). At this point it is important to prescribe that in general, vari-
784 ables and sensors capturing them used in early trials should not be discarded
785 from the ASPPID cycle unless the associated sensing equipment is expen-
786 sive enough to potentially cause a budgetary problem along the project, even
787 if by discarding them the team takes the risk of losing valuable predictive
788 interactions between variables.
- 789 (2.b) Sense: Once rented sensors were installed and fully operational, 3 new time
790 series were monitored: X_t^2 , X_t^3 and X_t^4 (masked for confidentiality), col-
791 lected over a period of $T_{col}^{min}(s_2) = 2$ months at 1 sample per second.
- 792 (2.c) Preprocess: repeating the feature extraction procedure made on X_t^1 with the
793 3 additional time series monitored by the new sensors, a dataset composed of
794 13522 parts, each characterized by 35 features, was produced, with a non-
795 quality sample rate – namely, imbalance ratio – of 13.44%. Standardization
796 and cleansing were implemented by using the same techniques as in X_t^1 on

797

 X_t^2 , X_t^3 and X_t^4 .

798

(2.d) Predict: an identical data analysis to that of the first iteration was done with this newly produced dataset: techniques for handling class imbalance and an inner RF classifier using 5×5 -fold cross-validation were used when building the model, rendering the scores listed in Table 2.

799

800

801

802

803

804

805

806

807

808

809

810

811

812

From all the generated models, the combination with the highest average F1 and recall values was selected (SMOTE + RF), since all parts undergo a quality control after being processed, thus precision is less relevant in the choice of the model. The study of the predictive importance for each variable (evinced by the RF model) suggested the discrimination of two feature groups comprising more than 50% of the predictive importance upon which to construct 2 different decision tree classifiers: one only with the first variable importance group (Group 1), and the second one with the first and the second group (Group 1+2). The reason for proceeding in this way hinges on the eventual implementation and deployment of the model: operating and fine tuning fewer variables within the production line is more cost-efficient and easily adoptable by the plant personnel.

Table 2: Average model scores of the second iteration of the ASPPID methodology.

| Model | Recall | Precision | F1 | AUC |
|-----------------------|---------------|------------------|-----------|------------|
| OVERSAMPLING + RF | 0.4604 | 0.7604 | 0.5733 | 0.7585 |
| SMOTE + RF | 0.5989 | 0.5728 | 0.5854 | 0.6658 |
| BORDERLINESMOTE + RF | 0.5160 | 0.6757 | 0.5850 | 0.7227 |
| SMOTE-TOMEKLINKS + RF | 0.6070 | 0.5595 | 0.5821 | 0.6593 |
| ADASYN + RF | 0.5324 | 0.6604 | 0.5893 | 0.7125 |
| ENN + RF | 0.6243 | 0.5345 | 0.5756 | 0.6455 |
| OSS + RF | 0.4971 | 0.7155 | 0.5864 | 0.7361 |

813

814

815

816

817

818

819

820

821

822

823

Unlike Random Forest, classification trees can be sensitive to collinearity [43], so the first step is to verify the multi-linearity of the variables involved in Group 1, namely, V09, V10, V15 and V16. In this regard, a pairwise correlation matrix revealed a high linear dependency between the selected variables. Based on this insight, tree models were constructed over different data subsets based only on variable v16, whose decision threshold was proven to be quite stable with respect to different training samples (by comparing among trees trained over different subsets of the available data). If both groups of features were used, variables v14, v1, v8, v4 and v3 composing Group 1+2 would be included in the dataset. A similar correlation

824 analysis to the one performed previously evinced that collinearity was only
 825 present between the variables of the first group. Therefore, all features from
 826 the second group were fed to the decision tree model, along with feature
 827 v16 from the first group.

828 Figures 8.a and 8.b show the most frequent classification tree resulting from
 829 creating the explanatory model over 5 folds of the balanced dataset with
 830 the variables of groups 1 and 1+2, respectively. Each node shows the pre-
 831 dominant predicted class (ACCEPTED/REJECTED), the percentage of obser-
 832 vations and the occurrence of the predominant predicted within such obser-
 833 vations at every node. It can be noted that v16 is capable, by itself, of
 834 separating the two classes. As for tree built using group 1+2, the number of
 835 variables to be controlled is greater, yielding a too complex set of rules to
 836 be implemented in production. Again, variable v16 is evinced to split both
 classes perfectly.

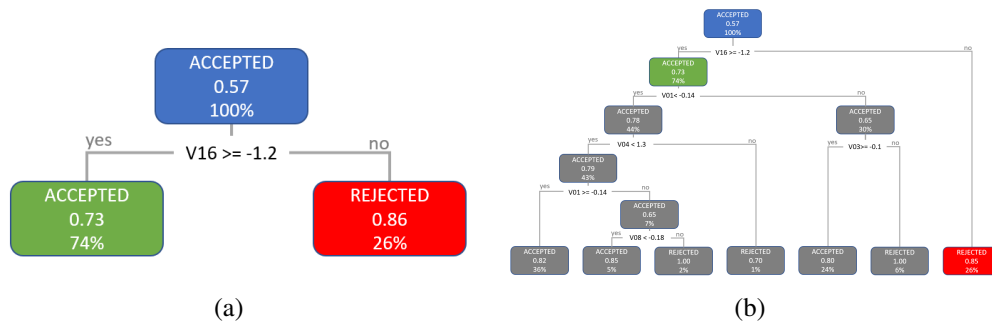


Figure 8: Most frequent decision trees inferred from the collected information with different feature subsets (a) Group 1; (b) Group 1+2.

837
 838 Results corresponding to the application of the above trees are shown in
 839 Table 3. Despite the NQR is lower in the tree trained over groups 1 and
 840 2, the rule set extracted from the tree grown by the model was definitely
 841 considered to be overly complex to be implemented in the production line.
 842 Even though both models enhanced their NQR scores with respect to the
 843 benchmark in trial $k = 1$ (supporting their potential to produce a maximum
 844 non-quality rate equal to 8% in the implement phase), the Pareto trade-off
 845 between implementability of the rules extracted from the model and their
 846 performance gain in terms of recall was reported by the improvement team
 847 to the engineering manager. Driven by its simplicity, the manager opted for
 848 the simplest model based on thresholding variable v16.

849 (2.e) Implement: in this phase the rule extracted from the explanatory decision
 850 tree model selected after the *predict* phase was implemented and run in pro-
 851 duction for $T_{imp}(s_2) = 3$ months. Specifically the plant personnel ensured
 852 that the conditions dictated by the selected tree model to yield less defective
 853 parts (since the simplest model was selected, it reduced to $V16 < -1.2$)
 854 were fulfilled by the machinery. This new production directive produced
 855 23966 parts, with 1672 REJECTED and 22294 ACCEPTED parts, i. e.
 856 a non-quality rate of 6.98%, achieving the target for the ASPPID project
 within $K = 2$ iterations.

Table 3: NQR (mean± confidence interval) of rules extracted from feature groups.

| Features | Original data | SMOTE + Decision Tree | |
|------------------------------|---------------|-----------------------|----------------|
| | | Min NQR branch | Max NQR branch |
| Group 1+2 (<i>predict</i>) | 13.34 ± 0.47% | 7.05 ± 0.49% | 59.67 ± 2.04% |
| Group 1 (<i>predict</i>) | 13.34 ± 0.47% | 7.37 ± 0.52% | 57.76 ± 2.21% |
| Group 1 (<i>implement</i>) | – | 6.98% | – |

857

858 Once verified by the improvement team and presented to the management
 859 team, the rented second sensor was purchased, increasing only 8.62% the
 860 overall cost of the trial with respect to the case where sensor s_2 was pur-
 861 chase from its beginning. However, as shown on Figure 9 the overall cost of
 862 sensing equipment spent during the ASPPID project, including the first and
 863 second iteration, was reduced down to 43.11% with respect to a traditional
 864 sales and purchase scheme. These cost savings were acknowledged by the
 865 management team to a right, informed choice of the sensing equipment, and
 866 a close match between the performance metrics of the developed predictive
 867 model and the requirements of the use case, which blended together ease of
 868 deployment and compliance with the imposed performance target.

869 4.3.3. Deploy

870 Once the results of the *implement* phase were verified, the project manager
 871 updated the FMEA information in the knowledge base of the company [44]. Al-
 872 though the case study was done on an experimental process within Edertek, the
 873 deployment of the model on a real production line will be straightforward thanks
 874 to the insertion of this simple yet efficient rule in the FMEA database: a rule
 875 extracted by virtue of the proposed ASPPID methodology.

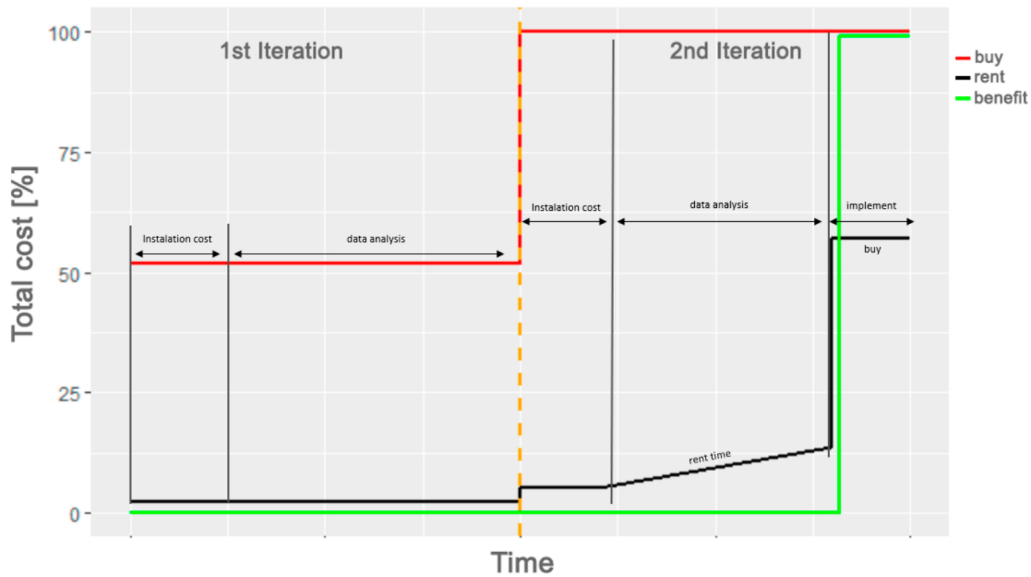


Figure 9: Graphical comparison between costs derived from renting (—) and buying (—) strategies in the acquisition of sensors within the ASPPID methodology. The plot also includes the progression of the benefit (—) along the SPPID cycle.

876 5. Concluding Remarks and Open Research Lines

877 This paper has presented ASPPID, a novel iterative methodology for effi-
 878 ciently selecting sensing equipment in industrial processes and machinery, es-
 879 sentially capitalizing on a predominant role of the data scientist along the entire
 880 decision workflow, not only within the stages where the captured data are ex-
 881 ploited. The proposed methodology comprises steps where all decisions are made
 882 towards minimizing the costs and fulfilling the target benefit of the overall project.

883 For this to be realized, ASPPID proposes to arrange a team where the lead-
 884 ing data scientist bridges the traditional gap between the purchase of sensors and
 885 its expected utility with respect to the objectives of the project. Furthermore, if
 886 allowed by the supplier ASPPID recommends exploring new forms of temporal
 887 sensor provisioning and installation, such that the extra renting costs are compen-
 888 sated by a lower risk taken by the team, specially when the project is held with
 889 data generated in application scenarios with scarce or null good practices reported
 890 by the community. The adoption of this methodology is specially suitable when
 891 there is no a priori information about the potential of the data captured by the
 892 sensor to bias for the project at hand; in this case a deep involvement of the data
 893 scientist in the selection of the sensing equipment to be tested along the trials is
 894 deemed crucial not to oversize the sensing requirements of the plant.

895 A use case held in a real manufacturing industry was used to validate the
896 proposed methodology. The application of SPPID gave rise to a reduction of the
897 sensor costs over 43% with respect to the case where all sensors were acquired
898 from the beginning, while meeting a scrap ratio (NQR) less than the target value
899 imposed at the beginning of ASPPID.

900 Future research will be devoted towards the inclusion and hybridization of ex-
901 isting production models for data mining (e.g. CRISP-DM) and agile software
902 development (corr. SCRUM or Lean Software, among others) within the ASPPID
903 methodology, placing an special emphasis on how to incorporate multiple conflict-
904 ing decision criteria to the decision making process. Also new real applications of
905 ASPPID to other industrial scenarios will be performed, such as predictive main-
906 tenance of remote machinery and the detection of non-technical losses in Smart
907 Grids, among others.

908 **Acknowledgments**

909 The real case could not be possible without the participation of Fagor Ederlan
910 S.Coop. and its R&D center, Edertek. This work has been supported in part
911 by the ELKARTEK (ref. KK-2016/00096, BID3ABI project) and EMAITEK
912 programs of the Basque Government. Antonio J. Nebro is supported by Grants
913 TIN2014-58304 (Ministerio de Ciencia e Innovación, Spain), and P11-TIC-7529
914 and P12-TIC-1519 (Plan Andaluz I+D+I, Spain).

915 **Bibliography**

- 916 [1] Bosch, Industry 4.0 – Germany takes first steps toward the next industrial
917 revolution, [http://blog.bosch-si.com/industry-4-0-germa-
918 ny-takes-first-steps-toward-the-next-industrial-r
919 evolution/](http://blog.bosch-si.com/industry-4-0-germany-takes-first-steps-toward-the-next-industrial-revolution/) (accessed 28 September 2017)
- 920 [2] Wang, S., Wan, J., Li, D., Zhang, C. 2016. Implementing smart factory of in-
921 dustrie 4.0: an outlook. *International Journal of Distributed Sensor Networks*
922 12(1): 3159805.
- 923 [3] Latham, S. F., Braun, M. R. 2008. The performance implications of financial
924 slack during economic recession and recovery: observations from the soft-
925 ware industry (2001-2003). *Journal of Managerial Issues* 20(1): 30-50.
- 926 [4] Hwang, K., Chen, M. 2017. *Big-Data Analytics for Cloud, IoT and Cognitive*
927 *Computing*. John Wiley & Sons.

- 928 [5] Browne, W., Yao, L., Postlethwaite, I., Lowes, S., Mar, M. 2006. Knowledge-
929 elicitation and data-mining: Fusing human and industrial plant information.
930 *Engineering Applications of Artificial Intelligence* 19(3): 345-359.
- 931 [6] Camacho, E. F., Bordons, C. 2012. *Model predictive control in the process*
932 *industry*. Springer Science & Business Media.
- 933 [7] Zhang, Z., He, X., Kusiak, A. 2015. Data-driven minimization of pump op-
934 erating and maintenance cost. *Engineering Applications of Artificial Intelli-*
935 *gence* 40: 37-46.
- 936 [8] Pfohl, H. C., Yahsi, B., Kurnaz, T. 2017. Concept and Diffusion-Factors of
937 *Industry 4.0 in the Supply Chain*. *Dynamics in Logistics*, 381-390.
- 938 [9] Gunasekaran, A., Papadopoulos, T., Dubey, R., Wamba, S. F., Childe, S. J.,
939 Hazen, B., Akter, S. 2017. Big data and predictive analytics for supply chain
940 and organizational performance. *Journal of Business Research*, 70: 308-317.
- 941 [10] Chehbi-Gamoura, S., Derrouiche, R. 2017. Big Valuable Data in Supply
942 Chain: Deep Analysis of Current Trends and Coming Potential. In *Working*
943 *Conference on Virtual Enterprises* 230-241. Springer, Cham.
- 944 [11] Lasi, H., Fettke, P., Kemper, H. G., Feld, T., Hoffmann, M. 2014. *Industry*
945 *4.0*. *Business & Information Systems Engineering* 6(4): 239-242.
- 946 [12] Best, M., Neuhauser, D. 2006. Walter A Shewhart, 1924, and the Hawthorne
947 factory. *Quality and Safety in Health Care* 15.2 (2006): 142-143.
- 948 [13] Lim, S. A. H., et al. 2017. A systematic review of statistical process control
949 implementation in the food manufacturing industry. *Total Quality Manage-*
950 *ment & Business Excellence* 28(1-2): 176-189.
- 951 [14] Madanhire, I., Mbohwa, C. 2016. Application of Statistical Process Control
952 (SPC) in Manufacturing Industry in a Developing Country. *Procedia CIRP* 40:
953 580-583.
- 954 [15] Deming, W. E. 1981. Improvement of quality and productivity through ac-
955 tion by management. *Global Business and Organizational Excellence* 1(1):
956 12-22.
- 957 [16] Bouranta, N., et al. 2017. Identifying the critical determinants of TQM and
958 their impact on company performance: evidence from the hotel industry of
959 Greece. *The TQM Journal* 29(1): 147-166.

- 960 [17] Sinha, N., et al. 2016. Mapping the linkage between organizational culture
961 and TQM: the case of Indian auto component industry. *Benchmarking: An*
962 *International Journal* 23(1): 208-235.
- 963 [18] Becker, R. M. 1998. Lean manufacturing and the Toyota production system.
964 *Encyclopedia of World Biography*.
- 965 [19] Harry, M. J. 1998. Six Sigma: a breakthrough strategy for profitability. *Qual-*
966 *ity Progress* 31(5): 60.
- 967 [20] Shokri, A., Bradley, G., Nabhani, F. 2016. Reducing the scrap rate in an
968 electronic manufacturing SME through Lean Six Sigma methodology.
- 969 [21] Delgado López, E. 2016. Propuesta de un plan para la reducción de la merma
970 utilizando la metodología six sigma en una planta de productos plásticos. Doc-
971 toral dissertation (in Spanish), Pontificia Universidad Católica del Perú, Es-
972 cuela de Posgrado.
- 973 [22] Winter, D., et al. 2017. Using real-time data for increasing the efficiency
974 of the automated fibre placement process. *International Journal of Vehicle*
975 *Structures & Systems* 9.1: 11.
- 976 [23] Braha, D., ed. 2013. *Data mining for design and manufacturing: methods*
977 *and applications*. Vol. 3. Springer Science & Business Media.
- 978 [24] Rahim, M. S., Rahman, M. Chowdhury, A. E. 2017. Mining Industrial En-
979 gineered Data of Apparel Industry: A Proposed Methodology. *International*
980 *Journal of Computer Applications* 161(7): 0975-8887.
- 981 [25] Enarsson, L. 1998. Evaluation of suppliers: how to consider the environ-
982 ment. *International Journal of Physical Distribution & Logistics Management*
983 28(1): 5-17.
- 984 [26] Burdick, R. K., Borrer, C. M., Montgomery, D. C. 2003. A review of meth-
985 ods for measurement systems capability analysis. *Journal of Quality Technol-*
986 *ogy* 35(4): 342.
- 987 [27] Dietrich, E. 2012. Capability of measurement processes based on ISO/FDIS
988 22514-7 and VDA 5. XX IMEKO World Congress Metrology for Green
989 Growth, Republic of Korea, 9-14.
- 990 [28] Chmielewski, M. R., Grzymala-Busse, J. W. 1996. Global discretization
991 of continuous attributes as preprocessing for machine learning. *International*
992 *journal of approximate reasoning* 15(4): 319-331.

- 993 [29] Shalabi, L. A., Shaaban, Z. 2006. Normalization as a preprocessing engine
994 for data mining and the approach of preference matrix. International Confer-
995 ence on Dependability of Computer Systems, Poland, 207-214.
- 996 [30] Santos, M. S., et al. 2017. Influence of data distribution in missing data impu-
997 tation. Conference on Artificial Intelligence in Medicine in Europe. Springer,
998 Cham, 285-294.
- 999 [31] Little, R. J. A., Rubin, D. B. 2014. Statistical analysis with missing data.
1000 John Wiley & Sons.
- 1001 [32] Gamberger, D., Lavrac, N., Dzeroski, S. 2000. Noise detection and elimina-
1002 tion in data preprocessing: experiments in medical domains. Applied Artifi-
1003 cial Intelligence 14(2): 205-223.
- 1004 [33] Velleman, P. F., Hoaglin, D. C. 1981. Applications, basics, and computing
1005 of exploratory data analysis. Duxbury Press.
- 1006 [34] Kotsiantis, Sotiris, Dimitris Kanellopoulos, and Panayiotis Pintelas. "Han-
1007 dling imbalanced datasets: A review." GESTS International Transactions on
1008 Computer Science and Engineering 30.1 (2006): 25-36.
- 1009 [35] He, H., Garcia, E. A. 2009. Learning from imbalanced data. IEEE Transac-
1010 tions on knowledge and data engineering 21(9): 1263-1284.
- 1011 [36] Chawla, N. V., Bowyer, K. W., Hall, L. O., Kegelmeyer, W. P. 2002.
1012 SMOTE: synthetic minority over-sampling technique. Journal of Artificial In-
1013 telligence Research 16: 321-357.
- 1014 [37] Tomek, I. 1976. Two modifications of CNN. IEEE Transactions on Systems,
1015 Man and Cybernetics 6: 769-772.
- 1016 [38] Han, H., Wang, W. Y., Mao, B. H. 2005. Borderline-SMOTE: a new over-
1017 sampling method in imbalanced data sets learning. Advances in intelligent
1018 computing, 878-887.
- 1019 [39] He, H., Garcia, E. A., Li, S. 2008. ADASYN: Adaptive synthetic sampling
1020 approach for imbalanced learning. IEEE International Joint Conference on
1021 Neural Networks, 1322-1328.
- 1022 [40] Wilson, D. L. 1972. Asymptotic properties of nearest neighbor rules using
1023 edited data. IEEE Transactions on Systems, Man, and Cybernetics 2(3): 408-
1024 421.

- 1025 [41] Kubat, M., Matwin, S. 1997. Addressing the curse of imbalanced training
1026 sets: one-sided selection. International Conference on Machine Learning, 97.
- 1027 [42] López, V., et al. 2013. An insight into classification with imbalanced data:
1028 Empirical results and current trends on using data intrinsic characteristics.
1029 Information Sciences 250: 113-141.
- 1030 [43] Breiman, L. 2001. Random forests. Machine learning 45(1): 5-32.
- 1031 [44] McDermott, R., Mikulak, R. J., Beauregard, M. 1996. The basics of FMEA.
1032 SteinerBooks.