*mper*

# APPLICATION OF INSTANCE-BASED LEARNING FOR CAST IRON CASTING DEFECTS PREDICTION

Robert Sika, Damian Szajewski, Jakub Hajkowski, Paweł Popielarski

*Poznan University of Technology, Faculty of Mechanical Engineering and Management, Poland*

*Corresponding author:*
*Robert Sika*
*Poznan University of Technology*
*Faculty of Mechanical Engineering and Management*
*Piotrowo 3, 60-965 Poznań, Poland*
*phone: +48 61 6652459*
*e-mail: robert.sika@put.poznan.pl*

ABSTRACT
The paper presents an example of Instance-Based Learning using a supervised classification method of predicting selected ductile cast iron castings defects. The test used the algorithm of $k$-nearest neighbours, which was implemented in the authors' computer application. To ensure its proper work it is necessary to have historical data of casting parameter values registered during casting processes in a foundry (mould sand, pouring process, chemical composition) as well as the percentage share of defective castings (unrepairable casting defects). The result of an algorithm is a report with five most possible scenarios in terms of occurrence of a cast iron casting defects and their quantity and occurrence percentage in the casts series. During the algorithm testing, weights were adjusted for independent variables involved in the dependent variables learning process. The algorithms used to process numerous data sets should be characterized by high efficiency, which should be a priority when designing applications to be implemented in industry. As it turns out in the presented mathematical instance-based learning, the best quality of fit occurs for specific values of accepted weights (set #5) for number $k = 5$ nearest neighbours and taking into account the search criterion according to "product index".

KEYWORDS
Soft modelling, instance-based learning, $k$-nearest neighbours algorithm, cast iron casting defects, computer application.

## Introduction

Industry, as a field of material production, must face growing globalization these days. The market confronts enterprises with the necessity of constant adjusting to increase customer demands [1]. Production companies are purchasing increasingly newer machines and devices [2]. Multiproduct manufacturing [3], advanced methods of process evaluation [4] using specific Data Mining methods (soft modelling) for multidimensional evaluation of the manufacturing process [5] and the use of hard modeling methods, which significantly reduce production costs [6] are increasingly used. Classical approaches, based on the principle of 'paper and pencil' are in use for many years are no longer sufficient or effective. It is known in many companies, quality management systems (QMS) are implemented, that is associated with specific activities, which can be less or more efficient. The authors assume [7] that level of efficiency is an indicator of the maturity of a QMS, without which company is uncompetitive. In the constantly developing Industry 4.0 concept integration of all areas of an enterprise activities is assumed, starting with business or manufacturing processes relating to machines, employees, customers and orders [8]. It is possible thanks to Acquisition and Data Mining (A&DM) approach, i.e. collection and processing of data generated by the basic activity of the company, obviously with no data redundancy. Tools for collaborative production based on the cloud are increasingly used [9]. Mostly, all these tasks are taken over by management support IT systems which are used on various levels, i.e. ERP systems (Enterprise Re-

sources Planning) and MES systems (Manufacturing Execution Systems). They turn raw data into useful information, which can be used by relevant people in various positions. The Industry 4.0 Concept presented in Hanover in 2011 encompassed four key assumptions:

- Internet and Internet of Things availability and use,
- technical and business processes integration within the enterprise,
- digital mapping and virtualisation of the real world,
- use of SMART factory including SMART means of production and SMART products.

Benefits resulting from following Industry 4.0 are relatively easy to measure in many enterprises. They mainly relate to reduce production costs by 10–30%, reducing logistic costs by 10–30% and reducing quality control costs by 10–20% [10]. Today's enterprises face the huge challenge of processing large amounts of daily data generated. The amount is bound to grow each year. This trend, called the Big Data, has 5 features [11]:

- large amount of data,
- large variety of data,
- high rate of data generation,
- significant data reliability
- data value (the degree of data valuable).

The above are a great challenge for companies. It is thus necessary to develop tailored IT solutions. Without the right tools it is not possible to draw significant conclusions from the large amount of industrial data. This aspect has recently been particularly visible in foundries [12, 13].

It is important that the technological and organisational complexity of casting production results from material diversification and the course and relations of partial processes susceptible to random and/or systematic factors. Casting production processes can be a large extend reproducible and controllable under proper monitoring with the use of a systematic data acquisition [14]. Appropriate Data Mining tools use for support of process management may be of great assistance.

## State of art

Soft modelling is data modelling (drawing the right conclusions) with the use of specialist knowledge (expert knowledge) on the analysed process through appropriate measures on empirical data (mathematical soft modelling) [15]. Information on such a process may be residual or limited and models may be built on the basis of assumptions and dependencies obtained from the provided data analysis. Such a modelling method has flaws as many processes in enterprises are not permanent. Deviations from processes are varied and cannot be compared. Consequently, the resulting relationships have low reliability [12–14, 16].

However, the increasing amount of data generated results in measurable compliance with the reality. Thus created models do not ensure an accurate course mapping, but help in estimating the expected result. Nevertheless, advanced methods allow to determine the impact of individual parameters on the final result as well as to evaluate the weight of specific parameters [17, 18].

Modelling quality mainly depends on the amount of input data, i.e. the fewer the data the more coarse the determination of possible outcomes. At the same time the more data the greater the probability of errors resulting from the data analysis. This is why monitoring and adapting the most important production processes on the basis of the actual process are necessary. High flexibility is soft modeling advantage, as it allows quick and easy analysis of process behaviour [19].

The $k$-nearest neighbours ($k$-nn) algorithm belongs to the group of instance-learning algorithms, where the learning set is selected, and new objects (test, reference) classification is made by comparing with the most similar objects from the learning set. The algorithm works as follows [20, 21]:

- data scaling by applying standardization or normalization,
- setting distance between vectors from the test and the training sets,
- sorting distances and selecting the most frequently used $k$ label.

The key step in the $k$-nn algorithm is to determine the distance between the values from the reference data set and the values stored in the database. This allows to specify $k$-nearest neighbours, including the input values in regards to the output values most often using the Euclidian metric. There are also other ways of finding $k$-nearest neighbours, e.g. matrices multiplication composed of input reference values and values from the learning and testing sets. Although the most effective way of finding the optimal result taking into account the accuracy (matching) criterion is: Manhattan, Chebyshev, Minkowski, Hamming, GEMM (matrices multiplication), but it is used more rarely than Euclidian metric [22]. This is caused by the necessity of calculating all the distances with the use of matrices (even for a strictly defined number of input parameters), which is not necessary for $k$-nearest neighbours. In view of the

above, the $k$-nn algorithm has the advantage over matrices multiplication in view of significant time reduction and a very small error of the model relation to real phenomena [20–23].

## Research methodology

The primary objective of the research was to predict selected casting defects for ductile cast iron, the causes of which should be sought in the group of parameters relating to the moulding sand formation and the pouring process (surface defects i.e.: surface roughness, metal penetration, burn-on, veining, sand holes). In case of input parameters (independent variables) one database record was built as follows:

- for a specific casting series produced during the day,
- for nominal values of selected green sand parameters,
- for average values of selected pouring parameters,
- with about 30-minutes frequency.

The above assumptions were dictated by the time interval between green sand's properties measurements in the moulding sand laboratory. Additionally, it was important when the series was being poured for more than half an hour – a more extensive database for learning could be obtained in such a case. The output parameter (dependent variable) in turn referred to the percentage share of casting defects estimated for a specific one day series. Therefore, the same defect percentage share for the dependent variable was assigned for each record built by independent variables. The database (several XLS files chronologically related according to the cast assortment) from selected production processes was built in real production conditions with the use of an original database system, according to assumptions adequate to the foundry conditions [24]. Input and output parameters selected for modelling are presented in Table 1.

The data modelled were relating to production of technologically similar small-sized castings (similar shape and casting weight, one cast in mould, one core located in the casting made by the cold-box method) on the Künkel Wagner automatic molding line with a horizontal mould division. The paper presents the results of comparison of the rate and quality of action of the $k$-nn algorithm using weighing to predict casting defects. The aim was to select the efficient $k$-nn algorithm with a good quality of casting defects predict based on similar work in this area [23, 25–28].

Table 1
Parameters influencing for selected defect occurrence.

| Symbol | Parameter | Tolerance limit | Unit |
|--------|-----------|-----------------|------|
| \multicolumn{4}{c}{Input parameters – GREEN SAND PARAMETERS} | | | |
| $M$ | Moisture | 2.4–4.5 | % |
| $Pm$ | Permeability | 100–220 | m$^2$/Pa*s*$10^{-8}$ |
| $Cms$ | Compression | 0.133–0.299 | MPa |
| $Cs$ | Compaction | 26–62 | % |
| $T$ | Temperature | 25–52 | °C |
| \multicolumn{4}{c}{Input parameters – POURING PARAMETERS} | | | |
| $T_{\text{POUR}}$ | Pouring temperature | 1285–1502 | °C |
| $T_{\text{TIME}}$ | Pouring time | 8–14 | s |
| \multicolumn{4}{c}{Output parameters – SURFACE DEFECTS} | | | |
| $SR$ | Surface roughness | | |
| $MP$ | Metal penetration | | |
| $BO$ | Burn-on | zero defects | % |
| $VE$ | Veining | | |
| $SH$ | Sand holes | | |

## Results

Having the selected research problem in mind, it was assumed that the $k$-nn algorithm would be tested by multiplying the square distance of the $m$-attribute vector by the weight of individual attribute sets, including their scaling. Based on a specific metric (Euclidean distance was used in the study), the algorithm searches for $k$ matched cases nearest to the new reference object (1) [28].

$$d_E(x,y) = \sqrt{\sum_i (x_i - y_i)^2 \times w_i}, \qquad (1)$$

$x_{i-1}^m$, $y_{i-1}^m$ – $m$ attributes vectors (input parameters) of two records ($x$ – reference ones and $y$ – learning and testing ones), $w_i$ – weight for individual attribute sets including their scaling.

The searched cases are the $k$-nearest records with the estimated shortest Euclidean metric for the selected input parameters for which the actual number of defective castings (with unrepairable casting defects) is shown (the output parameter, not included in the Euclidean metric) [29–31].

A mobile application was designed for the prepared numerical $k$-nn code. It constitutes a database management system along with a decision support algorithm within prediction of surface cast defects occurrence (raw surface and after tooling). Seven parameters (expert knowledge) with the greatest influence on specific casting defects occurrence were selected for the analysis. These parameters can be divided into two groups:

- relating to green sand properties (moisture, permeability, compression, compaction, green sand temperature),
- directly relating to casting process (pouring temperature and pouring time).

Considering the database structure, three SQL algorithms were used to sort estimated Euclidean distances. These algorithms use various mechanisms called: MergeSort, QuickSort and HeapSort [32–35]. They were automatically selected according to the database filtering method (optionally according to the product ERP index and/or production line number) taking into account different weight values [28, 32]. The percentage share of cast defects for 5 weight sets was predicted with the use of over 10.000 training records built in real industrial conditions in an cast iron foundry in the period between 03 January and 30 June.

For the tested algorithm, the recommended (although not necessary) data normalisation was not performed. Individual independent variables' weights were selected so as to ultimately achieve a balanced load not causing any parameter's rank increase (Weight 1, Table 2). Such values were initially set in the designed algorithm. Unfortunately, no satisfactory algorithm results were obtained – the REAL number of defective castings sometimes differed several times from the PREDICTED number of defective castings. The quality of the prediction for the first weight set fluctuated around the value of 0.5–0.6. At a later stage, weight selection was based on specific process expert knowledge (co-operation with moulding sand production and smelting departments technologists) and further 4 test sets differentiating selected production parameters were selected in accordance with the significance of their impact on surface defect formation (Weights 2–4, Table 2). Particularly important parameters affecting surface defects were marked as "sig-aff" (significantly affecting). After another iterative application of the $k$-nn algorithm with the possibility of weight modification, the weight value for the $M$ and $T_{POUR}$ parameters was finally increased, for the $Cms$ and $T$ parameters it was reduced and for the $Cs$ parameter the weight value was slightly corrected. The $T_{TIME}$ parameter value weight was left unchanged (Weight 5, Table 2).

Figure 1 presents the example of pre-processing stage, i.e. the screen with the parameters assigned to the new reference object. New values of specific production parameters registered during the production process are entered into the specially designed program form. This stage called pre-processing enables:

- selection of a specific product index (selecting the 'Without ERP index' option results in selecting all technologically similar series of cast),
- selecting a specific production line (selecting the 'Without ERP index' as above),
- entering new values specific production parameters from the GREEN SAND PARAMETERS and POURING PARAMETERS groups.

Table 2
List of weight sets used in instance-based algorithm for selected independent variables.

| Symbol | Weight set #1 original rescaling | Weight set #2 sig-aff change | Weight set #3 sig-aff change | Weight set #4 sig-aff change | Weight set #5 sig-aff final correction |
|---|---|---|---|---|---|
| GREEN SAND PARAMETERS | | | | | |
| $M$ (sig-aff) | 7 | 9 | 11 | 10 | 10 |
| $Pm$ | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |
| $Cms$ (sig-aff) | 125 | 115 | 105 | 95 | 75 |
| $Cs$ | 0.45 | 0.45 | 0.45 | 0.45 | 0.5 |
| $T$ | 0.6 | 0.6 | 0.6 | 0.6 | 0.35 |
| POURING PARAMETERS | | | | | |
| $T_{POUR}$ (sig-aff) | 0.015 | 0.017 | 0.019 | 0.021 | 0.02 |
| $T_{TIME}$ | 1 | 1 | 1 | 1 | 1 |



Fig. 1. Pre-processing stage. New independent variables involved in the process of a dependent variables predicting.
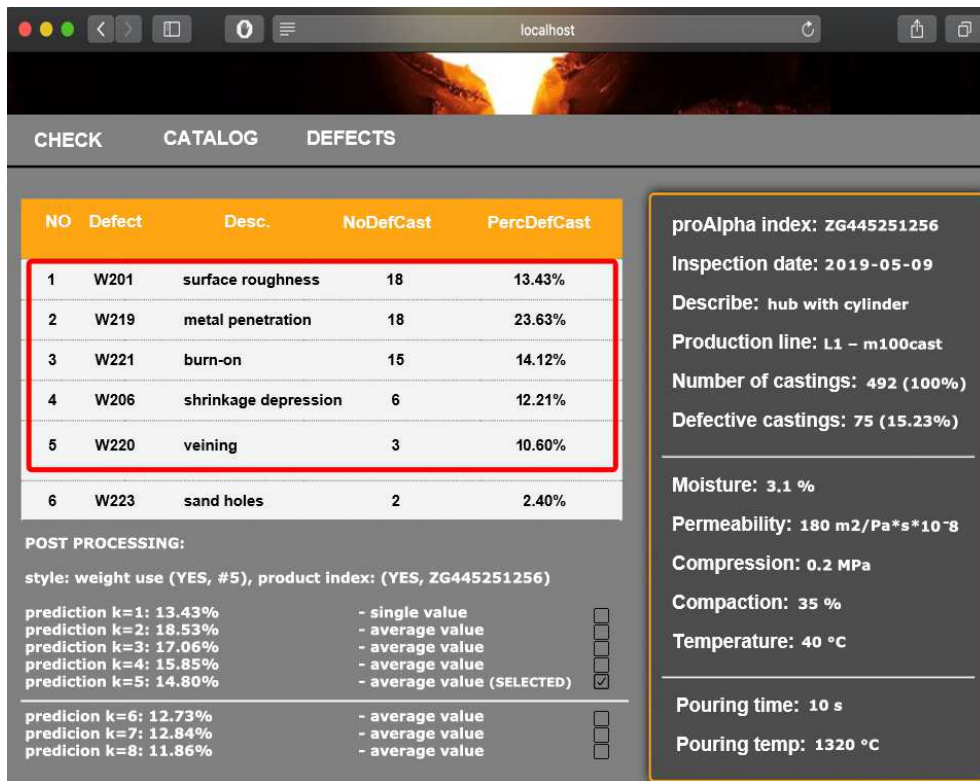
Fig. 2. Post-processing stage. $k$-nn algorithm result for finally selected weights set (set #5) and $k = 5$.

Learning records are put through learning by calculating the Euclidean distance and then sorted in ascending order with $k$ nearest neighbours (main-processing, solver). Finally the algorithm takes into account the possibility of manual adjusting particular parameter weight so that they better reflect the real results (expertise). At the post-processing stage the example of results are presented in a graphic form (Fig. 2).

Averaged results of the most advantageous weight (according to the criteria presented in the first column) set are presented in Table 3. Extended study of instance-based learning with the use of the $k$-nn algorithm taking into account the weighing system was used for the first weight set (original scaling without the use of selected parameters' gain) as well as for the fifth finally selected weight set.

Research using the search index (one assortment type) or not using it (many technologically similar assortments) was separately conducted. Sorting according to ProdLine was discontinued (the option from the algorithm level was turned off) due to the fact that the selected assortment type (technologically similar) was cast only on the L1 line production.

Table 3
Comparison of casting defect predicting results using the $k$-nn algorithm for original rescaling weight set (set #1) and finally selected weight set (set #5).

| PREDICTED CRITERIA | Weight set #1 | | Weight set #5 | |
|---|---|---|---|---|
| Product index | No | Yes | No | Yes |
| Predicting time | 11.03 sec | 0.16 sec | 10.93 sec | 0.16 sec |
| The PREDICTED number of defective castings (5-nn/best match of $k$-nn) | 8%/2-nn = 11.2% | 11.4%/1-nn = 12.5% | 12.3%/3-nn = 13.5% | 5-nn = 14.8% |
| The REAL number of defective castings | 15.65% | 15.24% | 15.65% | 15.24% |
| The quality of the prediction (5-nn/best match of $k$-nn) | 0.51/0.73 | 0.75/0.82 | 0.78/0.86 | 0.97 |

Obviously the quality of the prediction without the use of modified weights is correspondingly 0.73 and 0.82 depending on the use of search index sorting, whereas in both cases it was only obtained for the average value of two neighbours ($k = 2$) and one neighbour ($k = 1$). This means that a larger record number obtained at the post-processing stage ($k > 2$ and $k > 1$) strongly worsens the quality. On the other hand, too large $k$ number strongly generalises the result, which can be seen in independent variables' various values. And so, if take it assume a fairly representative value of $k$ neighbours [22], then the prediction quality is much lower – at the level of 0.51 and 0.75. In the specific case the value of $k$ was assumed to be $k = 5$ due to the fact that for larger $k$ values responses of independent parameters varied considerably (e.g. in an extreme case green sand moisture differed by up to 0.9% for $k = 7$).

After applying the finally selected weights (set #5) the prediction quality has improved and amounts to 0.78 for the case without using the search index. The best prediction quality and execution time of the $k$-nn algorithm takes place for indicating a specific index (index-search, as the specific name of the assortment) for the fifth weight set (set #5) as well as for fifth neighbour ($k = 5$) and gives a result of 0.97.

## Conclusions

The paper presents an example of mathematical modelling (soft modelling) using a supervised classification method for selected ductile cast iron castings defects predicting. The research uses the $k$-nn algorithm taking into account different weight sets for the examined attribute sets. The algorithm described was implemented in the authors' computer application. To ensure its correct operation it is necessary to have real historical data such as moulding sand, pouring parameters, chemical composition and percentage share of defective castings (not suitable for further repairs).

The use of a computer application using the $k$-nn algorithm in a particular foundry requires a tailored approach. The scope of the application use is wide. An example of this is the design of new processes and selection of significant production parameters within the pilot production management with the application presented.

The modified parameter values introduced to the reference object will allow to predict a variable share of specific defects. The application will also be useful in case of a sudden process unstable increasing production of defective castings. Then, on the basis of parameter values currently measured the level of defective castings can be predicted.

## References

[1] Ignaszak Z., Popielarski P., Hajkowski J., *Problem of Acceptability of Internal Porosity in Semi-Finished Cast Product as New Trend "Tolerance of Damage"*, Present in Modern Design Office, Defect and Diffusion Forum, 326–328, 612–619, 2012.

[2] Hamrol A., Zerbst S., Bozek M., Grabowska M., Weber M., *Analysis of the conditions for effective use of numerically controlled machine tools*, Lec. Not. in Mech. Eng., Springer, 201519, 3–12, 2018.

[3] Ivanov V., Dehtiarov I., Pavlenko I., Liaposhchenko O., Zaloga V., *Parametric optimization of fixtures for multiaxis machining of parts*, [in:] Hamrol A., Kujawińska A., Barraza M. [Eds], Advances in Manufacturing II, Lec. Not. in Mech. Eng., Springer, pp. 335–347, 2019.

[4] Kujawińska A., Vogt K., Diering M., Rogalewicz M., Waigaonkar M.D., *Organization of visual inspection and its impact on the effectiveness of inspection*, Advances in Manufacturing, Lec. Not. in Mech. Eng., Springer, 201519, 899–909, 2018.

[5] Rogalewicz M., Kujawińska A., Piłacińska M., *Selection of data mining method for multidimensional evaluation of the manufacturing process state*, Management and Production Engineering Review, Production Engineering Committee PAN, 3, 2, 27–35, 2012.

[6] Ignaszak Z., Hajkowski J., *Contribution to the Identification of Porosity Type in AlSiCu High-Pressure-Die-Castings by Experimental and Virtual Way*, Arch. of Found. Eng., 15, 1, 143–151, 2015.

[7] Grabowska M., Takala J., *Assessment of quality management system maturity*, Lec. Not. in Mech. Eng., Springer, 201519, 889–898, 2018.

[8] Presentation at the French Embassy in the Germany, "Industry of the future", http://www.amba-france-de.org/Vorstellung-des-neuen-franzosischen-Plans-Industrie-du-Futur-in-der-Botschaft, accessed: 05.06.2019, (2015).

[9] Varela M.L.R., Putnik G.D., Manupati V.K., Rajyalakshmi G., Trojanowska J., Machado J., *Collaborative manufacturing based on cloud, and on other I4.0 oriented principles and technologies: a systematic literature review and reflections*, Management and Production Engineering Review, 9, 3, 90–99, 2018.

[10] Rojko A., *Industry 4.0 Concept: Background and Overview*, Internat. J. of Interact. Mob. Technol. (iJTM), 111, 5, 77–90, 2017.

[11] Reis S.M., Kenett R., *Assessing the value of information of data-centric activities in the chemical processing industry 4.0*, AIChE – Proc. Sys. Eng., 64, 11, 3868–3881, 2018.

[12] Perzyk M., Kozłowski J., *Methodology of Fault Diagnosis in Ductile Iron Melting Process*, Arch. of Found. Eng., 16, 4, 101–108, 2016.

[13] Vijayaram T.R., Sulajman S., Hamouda A.M.S., Ahmad M.H.M., *Foundry quality control aspects and prospects to reduce scrap rework and rejection in metal casting manufacturing industries*, J. of Mater. Process. Technol., Elsevier, 178, 1, 39–43, 2006.

[14] Ignaszak Z., Sika R., Rogalewicz M., *Contribution to the Assessment of the Data Acquisition Effectiveness in the Aspect of Gas Porosity Defects Prediction in Ductile Cast Iron Castings*, Arch. of Found. Eng., 18, 1, 35–40, 2018.

[15] Khan S., Finkelstein L., *Mathematical modelling in the analysis and design of hard and soft measurement systems*, Measurement, Elsevier, 46, 8, 2936–2941, 2013.

[16] Stekh Y., Lobur M., Shvarts M., *Some methods for improving the accuracy of prediction recommendations*, Bulletin of Lviv Polytechnic National University, Series: Computer Systems Design Theory and Practice, 882, 46–49, 2017.

[17] Mahanta B.K., Chakraborti N., *Evolutionary Data Driven Modeling and Multi Objective Optimization of Noisy Data Set in Blast Furnace Iron Making Process*, Steel Research Int., Wiley, 89, 8, 1–11 (1800121), 2018, doi: 10.1002/srin.201800121.

[18] Zdobytskyi A., Lobur M., Iwaniec M., Breznitskyi V., *Optimization of the Structural Characteristics of the Robotic System Holder*, 15th International Conference on the Experience of Designing and Application of CAD Systems (CADSM), Polyana (Svalyava), Ukraine, IEEE, 2019.

[19] Grzegorzewski P., Kochański A., *From Data to Reasoning*, Part I – Theory in Soft Modeling in Industrial Manufacturing, Grzegorzewski P., Kochański A., Kacprzyk J. [Eds], Soft Modeling in Industrial Manufacturing, 2018.

[20] Gweon H., Schonlau M., Steiner S.H., *The k conditional nearest neighbor algorithm for classification and class probability estimation*, Neural Computation, 17, 3, 731–40, 2019.

[21] Sanodiya R., Saha S., Mathew K., *A kernel semi-supervised distance metric learning with relative distance: Integration with a MOO approach*, Expert Systems with Applications, 125, 233–248, 2019.

[22] Ryoo J., Arunachalam M., Khanna R., Kandemir M.T., *Efficient K nearest Neighbor Algorithm Implementations for Throughput-Oriented Architectures*, 19th Int'l Symposium on Quality Electronic Design, IEE, pp. 144–150, 2018.

[23] Juan L., *An Improved K-Nearest Neighbor Algorithm Using Tree Structure and Pruning Technology*, Intell. Autom. and Soft Comp., 25, 1, 35–48, 2019.

[24] Sika R., Ignaszak Z., *Data Acquisition procedures for A&DM systems dedicated for foundry industry*, Advances in Design, Simulation and Manufacturing II, pp. 692–701, 2019.

[25] Comac E., Arslan A., *A new training method for support vector machines: Clustering k-NN support vector machines*, Expert Systems with Applications, Elsevier, 35, 564–568, 2018.

[26] Balan K.P., *Metallurgical Failures Analysis*, Elsevier 2018.

[27] Thacker K.B., *Analysis of parameters for casting ductile iron pipe*, International Journal of Engineering Research and General Science, 3, 1, 496–503, 2015.

[28] Rogalewicz M., Sika R., Popielarski P., Wytyk G., *Forecasting of steel consumption with use of nearest neighbors method*, Modern Technol. Manufact. (MTeM 2017 – AMaTUC), 137, 01010-1–01010-6, 2017.

[29] Henon G., Mascre C., Blanc G., *Investigation of the cast products quality*, International Committee of Foundry Technical Associations, France, 1986.

[30] PN-85 H-83105, Casting Division and terminology of defects, Poland, 1985.

[31] Baler J., Koppen M., *Manual casting defects*, IKO-Erbsloh, 2004.

[32] All you need to know about sorting in Postgres, https://madusudanan.com/blog/all-you-need-to-know-about-sorting-in-postgres, accessed: 09.05.2019.

[33] Hills M., Klint P., Vinju, J. J., *Enabling PHP software engineering research in Rascal*, Science of Comp. Programm, Elsevier, 134, 37–46, 2016.

[34] Advanced features of PHP7, 2017, https://etrix-tech.com/advanced-features-of-php-7, accessed: 09.05.2019.

[35] Horwat W., *JavaScript 2.0: Evolving a Language for Evolving Systems*, 2005.