

48th CIRP Conference on MANUFACTURING SYSTEMS - CIRP CMS 2015

## Manufacturing lead time estimation with the combination of simulation and statistical learning methods

András Pfeiffer<sup>a,\*</sup>, Dávid Gyulai<sup>a,b</sup>, Botond Kádár<sup>a</sup>, László Monostori<sup>a,b</sup>

<sup>a</sup>Fraunhofer Project Center for Prod. Management and Informatics, Computer and Automation Res. Institute (SZTAKI), Kende str. 13-17, Budapest, HUNGARY

<sup>b</sup>Budapest University of Technology and Economics, Dept. of Manufacturing Science and Technology, Egy J. str. 1, Budapest, HUNGARY

\* Corresponding author. Tel.: +36 1 279 6176; E-mail address: [pfeiffer.andras@sztaki.mta.hu](mailto:pfeiffer.andras@sztaki.mta.hu)

### Abstract

In the paper, a novel method is introduced for selecting tuning parameters improving accuracy and robustness for multi-model based prediction of manufacturing lead times. Prediction is made by setting up models using statistical learning methods (multivariate regression); trained, validated and tested on log data gathered by manufacturing execution systems (MES). Relevant features, i.e., the predictors most contributing to the response, are selected from a wider range of system parameters.

The proposed method is tested on data provided by a discrete event simulation model (as a part of a simulation-based prediction framework) of a small-sized flow-shop system. Accordingly, log data are generated by simulation experiments, substituting the function of a MES system, while considering several different system settings (e.g., job arrival rate, test rejection rate).

By inserting the prediction models into a simulation-based decision support system, prospective simulations anticipating near-future deviations and/or disturbances, could be supported. Consequently, simulation could be applied for reactive, disturbance-handling purposes, and, moreover, for training the prediction models.

© 2015 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license

(<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Peer-review under responsibility of the scientific committee of 48th CIRP Conference on MANUFACTURING SYSTEMS - CIRP CMS 2015

**Keywords:** Lead time, Simulation, Robust prediction, Statistical learning

### 1. Introduction

Make to order production requires a proper estimation of manufacturing job lead times (LT) when dealing with production orders. Moreover, a reliable forecast on systems' load and output is also mandatory both for due-date quotation as well as for production control decisions. Since LT estimation is a difficult task, resulting in often unreliable output (e.g., when applying well-known shop-floor related characteristics and calculation methods as for instance combining total work content of the jobs and actual WIP), novel methods, considering a bigger set of system parameters influencing the LT are required. Though, discrete event simulation is well known and widely applied for predicting future systems' conditions, analytical interpretation of simulation outputs as prediction models would foster decision making on tactical level of production planning and control.

Simulation technologies are often used in supporting production control decisions and this is also particular for large-scale manufacturing systems. Several different applications of discrete-event simulation models in control of manufacturing systems were presented in [1] and [2].

A discrete-event simulator developed for the daily prediction of work in process (WIP) position in an operational wafer fabrication factory to support tactical decision-making is described in [3]. The model parameters are automatically updated by using statistical analyses performed on the historical event logs generated by the factory. A simulation study is presented in [4] which is applied to compare alternative WIP management policies, while [5] introduces a model, quantifying the effects of lot size changes, lot release controls and machine dispatching rules, on selected Key Performance Indicators (KPI-s) (throughput, process time and process time spread) for manufacturing steps. A simulation-based scheduling framework is presented in [6] for handling



[8]. In the on-line modes the simulation models represent the virtual mirror of the plant and run parallel to the real manufacturing environment, instantly simulating the future processes for a predefined short period.

In the paper, the off-line operation mode of the simulation and the prediction models are focused on. Interested readers might refer to [10], where on-line application of the simulation framework is introduced more in the details.

## 2.2. Prediction models for production control decisions

In order to extend the capabilities of the simulation towards prediction and estimation of future scenarios' results, the use of statistical learning models are proposed.

Basically, statistical learning refers to a set of methods for understanding and learning from data and providing solutions to understand the correlations among parameters and processes [11]. There are two main classes of these tools: the supervised and unsupervised learning techniques. Supervised learning aims at predicting some output parameters based on the input parameters and the priori known training set. The most fundamental supervised learning methods are the linear regression models capable of predicting a value of a quantitative output variable, assuming that there is approximately a linear relationship among the input/output variables. Other effective but simple techniques for practical applications are the tree-based methods that can be used both for regression and classification as well. The general idea behind these methods is the partition of the feature space into a set of disjoint rectangular regions, and fit a simple model in each one [11]. Building a regression tree over a given dataset is composed of two general steps. First, the feature space is divided into a set of disjoint regions, then for every observation which falls into a certain region the same prediction is made that is the mean of the region.

By building regression models over simulation data, one can estimate the production parameters even besides a dynamic environment. Tree-based models (e.g., random forests) are applied for estimating the capacity requirements of modular reconfigurable system by utilizing the results of several simulation runs in [12], while in [7], a tree construction approach is introduced for lead time estimation.

Regarding the simulation-related applications, regression and prediction models can be built over simulation-related data to support simulation-based optimization methods, in which some of the objective function or constraint(s) are represented by functions that are approximated by using the results of simulations [13]. The reason for applying simulation in these cases are usually the computational complexity or the lack of analytical expression of the objective function and/or constraints. These challenges are often faced when stochastic functions have to be represented in the optimization models [14]. A more general description of the input and output parameters of simulation models' is given by meta-models that are aimed at approximating the behaviour of system with mathematical functions [15]. In [16] regression models are introduced on simulation data to analyse the functional relationship among dispatching rules, due-date assignments and shop-load ration in job shops. Similar approach is applied

for a dynamic job-shop, however, simulation in this case was applied for evaluation only [17]. In [18], a multiple regression analysis platform was introduced, that enables prediction of different future scenarios considering the actual conditions of the production system.

It is assumed that by inserting the prediction models into the proposed simulation-based decision support system, prospective simulations anticipating near-future deviations and/or disturbances, could be more effectively supported. Consequently, simulation could be applied both for proactive and reactive, disturbance-handling purposes, and, moreover, for (off-line) training the prediction models.

## 3. Computational experiments

### 3.1. Description of the production system to be examined

The proposed method was tested on data provided by a discrete event simulation model of a small-sized parallel flow-shop system (Fig 2). Accordingly, log data are generated by simulation experiments, substituting the function of a MES system, while considering several different system settings (e.g., job arrival rate, test rejection rate). In each case a job is finished on a machine the actual status of the entire system is logged (timestamp, product type, location, WIP, buffer levels, etc.), considering information available in a real system.

The test production environment consists of five processing units (machines, grouped as stage 1 and stage 2) each with a buffer in front and a testing machine (stage 3) on which each product has to be tested (Fig 2). If this test is failed and the product is rejected, it is sent back to stage 2. Each job has a certain product type (A, B or C) assigned upon entering the system with equal probability and has to complete a process at each stage. There are two different levels of system load to be investigated, thus mean inter-arrival time ( $t_a$ ) of the jobs is 260 for a high and 295 seconds for a low system load, while no job collection or delay is applied. Routings are set dynamically, i.e., after entering the system or finishing processing on a machine the next machine is selected with the shortest input queue. Processing times are collected in Table 1. It is obvious that the processing times are product-type dependent at a selected stage, meanwhile the testing procedure may vary within a predefined time range (lower and upper bound), independently from the type of the job. Jobs are pulled to the machines from the input buffers by using the FCFS rule.

No setup times are considered in this model and the machines have a 90% availability during the production, but different mean times to repair (MTTR, Table 1).

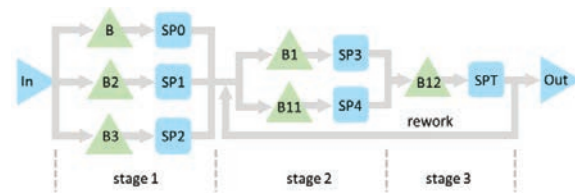


Fig 2. The layout and material flow of the test system.

Table 1. The main characteristics and parameters of the test production system.

	Stage 1 machines			Stage 2 machines		Stage 3 testing		
	SP0	SP1	SP2	SP3	SP4	SPT LB	SPT UB	
Processing time [min] (systematic fail ratio or reject rate)	Type A	12 (0.0)	12 (0.0)	12 (0.1)	6 (0.0)	6 (0.0)	2:50 (0.1)	3:30 (0.1)
	Type B	16 (0.1)	16 (0.0)	16 (0.1)	8 (0.0)	8 (0.0)	2:50 (0.1)	3:30 (0.1)
	Type C	10 (0.1)	10 (0.0)	10 (0.0)	3 (0.0)	3 (0.0)	2:50 (0.1)	3:30 (0.1)
Availability [%]	90	90	90	90	90	90		
MTTR [min]	10	10	10	20	20	10		

An important part of the experiments is to include some, so called, systematic failures during the processing of the jobs. It means that at a certain constellation (e.g. Type A on machine SP2) a certain probability is assigned for failing the process (e.g., the product type is difficult to assemble and the machine or operator usually makes errors). This does not influence the processing time, but the outcome of the testing process at stage 3. Consequently, in parallel to the normal reject rate (last two columns in Table 1), there are some jobs having a “systematic fail” built in, which will be recognized during the testing. After rework, jobs must be processed (reassembled) on stage 2 and tested again.

Thus, an expected total work content (TWK), indicating a minimal lead time, of each product types are the sum of the mean processing times along the manufacturing process (mean set-up time is zero while batch size equals 1): 1270, 1630 and 970 seconds, for type A, B and C, respectively.

The above detailed settings of the production systems resulted in a relatively high average machine utilization level on stage1 machines, 90 and for the testing machine at stage 3, 85. Stage 2 machine utilization was around 75 percent.

Moreover, by using the workload formula given in [19], the expected overall utilization level of the test system is 83% (without reject and rework). These levels could be considered as a proper level for lead time estimations [7], since, representing a real-world problem, there will be several jobs waiting in the system, making lead time estimation challenging, but, parallel keeping the expected WIP level manageable (e.g., by using Little’s law formula [19], the expected WIP is approx. 32 jobs). Note that as the utilization level (and system load) increases the estimation of the job lead times getting more difficult and unstable [7]. In Fig 3 the main characteristics of the job execution on the simulated test system are highlighted.

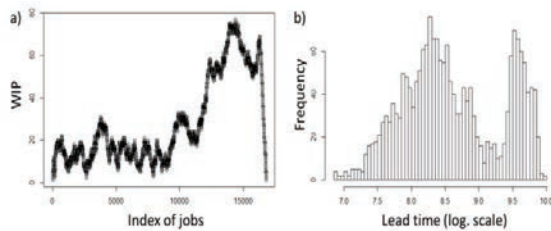


Fig 3. Main characteristics of a simulated execution of 2000 jobs in the test system for a lower system load level. a) Logged WIP level; b) Histogram of the job lead times on a logarithmic scale.

As the result of 2000 jobs going through the system the fluctuating WIP level (Fig 3a) and the histogram of the job lead times (Fig 3b, on a logarithmic scale) are given. It is obvious that the bimodal distribution of log-lead times are a consequence of the two distinct levels of WIP (before and after index 10000).

### 3.2. Selection of prediction parameters

The proposed statistical learning methods, introduced in Section 2.2 were applied on the historical log data provided by the simulation experiments, outlined previously. This section introduces first the feature selection applied in order to have the relevant parameters left in the models only. Then, two particular analysis are explained in details: 1) giving an explorative analysis, focusing on structure exploration (finding hidden failures); 2) applying prediction models for lead time estimation for different system conditions.

Once the test systems parameters are tuned in order to have a mostly stable and steady behaviour thorough the time horizon of the experiments in the simulation, the next step is to prepare the data for model formulation and estimation.

As it was stated before, collecting logs means that in each case a job is finished on a machine (or tested or event sent for rework) in the simulation system the actual status of the main parameter of the system is logged in a record of a log file (Table 2).

Here the main goal is to have particular entries in one line of the table containing all the relevant features (input variables) might describe the resulted outcome (output, or response variable) of the experiment (also referred to as observation). Therefore, log data must be preprocessed to be able to formulate the feature table (Table 3), necessary for model training and validation.

Table 2. Initial parameter set available from the log files

parameter	description
ID	Job ID
time_OUT	time stamp of when the process was finished
Location	machine ID, where process of the job was finished
rework	if the process was a rework
type	type of the job (product)
FailRate	1- the actual fail rate of the testing process (quality)
WIP	number of jobs in the system
Buffer levels	jobs waiting to be processed in the different queues

Table 3. Excerpt of a resulted feature table, aggregated to include relevant features describing jobs' behavior passing through the system.

ID	SP0	SP1	SP2	LT	type	FRate	WIP	B	...
1	1	0	0	1262	A	0.00	1	1	...
2	1	0	0	3335	C	1.00	7	1	...
3	1	0	0	4892	B	0.89	22	4	...
4	1	0	0	6769	A	0.92	20	6	...
5	0	1	0	4017	A	0.92	20	5	...

Routing of the jobs, as one of the major influencing factors, are assigned dynamically, but, however, can be tracked by the log entries (using *Location*, *ID* and *Type*). That means in the current representation, assigning a level one value to the certain machine (*SP0*, *SP1* or *SP2*) on which the job was processed at Stage 1. Other sections of the routing are irrelevant, thus are not mapped.

Similarly, job lead times are calculated from the logs (*LT* in Table 3) and assigned to the jobs by using the *time\_OUT* time stamps (Table 2). It was assumed that WIP level strongly influences *LT* in the system examined, thus, WIP and certain buffer levels (number of jobs waiting to be processed) in front of the machines are added to the feature table (*B*, *B2*, ... , *B12* in Table 3) as well.

Preliminary analysis of the features are given in Fig 4 in form of a correlation matrix of the selected variables. Note that the variable *type* is removed from this analysis, as it is a categorical variable [11]. As it was expected, the response variable *LT* (*lead* in Fig 4) is highly correlated with the features describing the system load (*WIP* and buffer levels).

3.3. Explorative analysis

After retrieving the feature table and all the log data from the simulation experiments, it is obvious to provide an explorative analysis, focusing on structure exploration, i.e., finding hidden relations resulting diverse response of certain inputs. In our case it means that (as introduced in Section 3.1) for some particular routing and job type combinations the probability to fail the test at Stage 3 will be above the "normal" rejection rate. It is presumed that the *LT*-s must reflect these constellations.

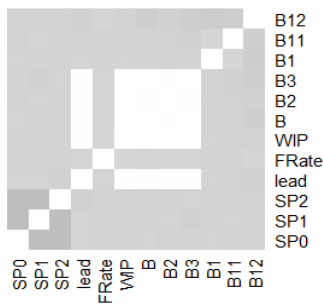


Fig 4. Graphical representation of the correlation matrix of the parameters in the feature table. White is a high positive, while dark grey is a negative correlation.

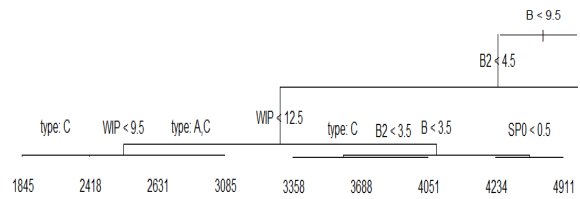


Fig 5. Prediction tree on the training set (excerpt).

However, in the current experiments introduced here these issues are known in advance, but in a real-world situation a structure exploration might be unique for identifying these underlying relationships.

In order to test the outcome of the predictions, first available data of the experiments were split into training and testing sets (one half each). A regression tree (

Fig 5, the left part is shown) was built by using the training set data only, but all the features available (general setting for the smallest allowed node size = 100). It can be seen that for low system load situations the shortest lead time (1845) is assigned to job type C (in accordance with the expected TWK). In contrast, e.g., if the WIP is above 12, the *LT* for a job type B is either 4234 or 4911 if in the routing *SP0* station was affected or not, respectively. Tree pruning was applied, and a new tree had been created in order to avoid over-fitting the model [11]. This is a systematic K-fold cross validation process to find the deviance or number of misclassifications as a function of the cost-complexity parameter [20]. In this new tree input variables used in the tree construction had been reduced, i.e., *B*, *B2* *WIP* and *SP2* were predictors with relevance.

3.4. Prediction on the test set

As the number of leaves (distinct and non-overlapping regions) in the tree are strongly limited, the prediction power of this model is expected to be low. The mean squared error (MSE) provided by the model when predicting from the test set was 952. Since it is apparent that the resulted tree based prediction model is useful for structure exploration as it was shown before, however, applying it for predicting *LT*, as the original goal, requires the application of other prediction models making comparison possible. Therefore, two other models, introduced in Section 2.2, were constructed on the training set. A multiple linear regression was formulated based on the strong correlation of system load related variables and the response variable *LT*. The general linear model resulted in the coefficients, where the higher intercept coefficient (3582) is compensated by the *FRate* (-3079). Regarding model quality, the *R*<sup>2</sup> value, showing the variance described by the model, is 0.975, while the MSE, calculated on the test set by the model, is 816.90, significantly lower compared to the previous tree based model.

Finally, an extension of the tree based method, a random forest model was constructed on the training data. This means that a number of decision trees are created on bootstrapped training samples [11] in order to try to reduce high



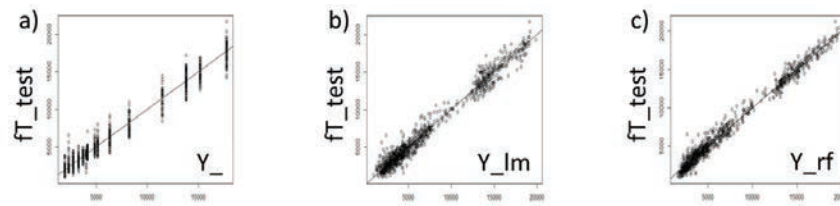


Fig 6. The function of model predictions ( $Y_{\_}$ ) and test set values ( $fT.test$ ) for the (a) prediction tree, (b) linear and (c) random forest models.

model variance. Building these decision trees, each time a split in a tree is considered, a random sample of  $m$  predictor variables is chosen as split candidates from the full set of  $p$  predictors. The split is allowed to use only one of those  $m$  predictors. This method is mostly helpful when dealing with a high number of correlated predictors. As expected, the quality of the prediction by applying this model on the test set, using 500 trees, is the best, compared to the other two solutions. MSE is 652.8 and  $R^2$  value is 0.98.

In Fig 6 the comparison of the prediction performance of the three models are highlighted. It can be seen that the prediction and test set values are the closest for the random forest model (values are close to the solid line representing a perfect prediction match) within the whole range of predicted LT values ( $Y_{\_}$ ). Contrary, the simple decision tree model provides a limited number of predicted values, thus in several cases the real output is away from the predicted one.

#### 4. Conclusions and future work

In the paper, a novel method was introduced for multi-model based prediction of manufacturing lead times. By inserting the prediction models into the proposed simulation-based decision support system, prospective simulations anticipating near-future deviations can be supported.

Future work covers testing the prediction models capabilities against changing, volatile system parameters. The degradation of the prediction power is expected, thus, however, a simulation supported re-training will be examined. This is important for the random forest model, which is highly sensitive for predictors taking values out of the expected boundaries (e.g., unexpected WIP level). However, the predictor variables are selected from a wider range of system parameters, applying systematic model selection methods would be a desirable way for improving model accuracy. Current results are to be compared to state-of-the-art analytical prediction solutions available in the literature.

#### Acknowledgement

Research has been partially supported by the European Union 7th Framework Programme Project No: NMP 2013-609087, Shock-robust Design of Plants and their Supply Chain Networks (RobustPlaNet). The authors express their thanks to the Hungarian Scientific Research Fund (OTKA) for its support (Pr. No.: 113038).

#### References

- [1] Banks, J., 1998, Handbook of Simulation, Principles, Methodology, Advances, Application and Practice, John Wiley & Sons Inc.
- [2] Law, A., Kelton, D., 2000, Simulation modeling and analysis, McGraw-Hill, New York.
- [3] Bagchi, S., Chen-Ritzo, C., Shikalgar, S.T., Toner, M., 2008, A full-factory simulator as a daily decision-support tool for 300mm wafer fabrication productivity, in Proceedings of the 2008 Winter Simulation Conference, S. J. Mason, R. R. Hill, L. Mönch, O. Rose, T. Jefferson, J. W. Fowler eds., p. 2021-2029.
- [4] Bureau, M., Dauzère-Pères, S., Yugma, C., Vermariën, L., Jean-Bernard M., 2007, Simulation results and formalism for global-local scheduling in semiconductor manufacturing facilities, Proc. of the 39th conference on Winter simulation, December, 2007, Washington D.C., pp. 1768-1773.
- [5] Sivakumar, A.I., Chong, C.S., 2001, A simulation based analysis of cycle time distribution, and throughput in semiconductor backend manufacturing, Computers in Industry, 45/1:59-78.
- [6] Lalas, C., D. Mourtzis, N. Papakostas, G. Chrysolouris, 2006, A Simulation-Based Hybrid Backwards Scheduling Framework for Manufacturing Systems, International Journal of Computer Integrated Manufacturing, 19/8:762-774.
- [7] Öztürk, A., Kayaligil, S., Özdemirel, N.E., 2006, Manufacturing lead time estimation using data mining, European J. of Op. Res., 173:683-700.
- [8] Monostori, L., Kádár, B., Pfeiffer, A., Karnok, D., 2007, Solution Approaches to Real-time Control of Customized Mass Production, CIRP Annals—Manufacturing Technology, 56/1:431-434.
- [9] Vánca, J., Monostori, L., Lutters, E., Kumara, S.R., Tseng, M., Valckenaers, P., Van Brussel, H., 2011, Cooperative and Responsive Manuf. Enterprises, CIRP Annals—Manuf. Technology, 60/2:797-820.
- [10] Pfeiffer, A., Kádár, B., Monostori, L., Vén, Z., Popovics, G., 2011, Co-analysing situations and production control rules in a large-scale manufacturing environment, The 44th CIRP Int. Conf. on Manufacturing Systems - New Worlds of Manufacturing, June 1 - 3, 2011, Wisconsin, Madison, USA, paper No. 101.
- [11] James, G., Witten, D., Hastie, T., Tibshirani, R., 2013, An Introduction to Statistical Learning, Springer, New York.
- [12] Gyulai, D., Kádár, B., Monostori, L., 2014, Capacity planning and resource allocation in assembly systems consisting of dedicated and reconfigurable lines, Procedia CIRP, 25:185-191.
- [13] Azadivar, F., 1999, Simulation optimization methodologies. In Proceedings of the 31st conference on Winter simulation: Simulation - a bridge to the future 1:93-100.
- [14] Gyulai, D., Kádár, B., Monostori, L., 2015, Robust production planning and capacity control for flexible assembly lines. In Proceedings of the 15th IFAC/IEEE/IFIP/IFORS Symposium, Information Control Problems in Manufacturing, IFAC, 2015, *In Press*.
- [15] Barton, R.R., 1992, Metamodels for simulation input-output relations. In Proceedings of the 24th conference on Winter simulation. ACM, 1992.
- [16] Cheng, T.C.E., 1988, Simulation study of job shop scheduling with due dates, International Journal of Systems Science, 19/3.
- [17] Sabuncuoglu, I., Comlekci, A., 2002, Operation-based flow time estimation in a dynamic job shop. Int J Management Sci. 30/6:423-442.
- [18] Han, S., Halpin, D.W., 2005, The use of simulation for productivity estimation based on multiple regression analysis. In Proceedings of the 37th conference on Winter Simulation, 2005.
- [19] Hopp, W.J., Spearman, M.L., 2001, Factory physics – Foundations of manufacturing management. Irwin McGraw-Hill.
- [20] R Core Team, R: A language and environment for statistical computing, 2013, R Foundation for Statistical Computing, <http://www.R-project.org/>