# A machine learning model predicting supplier delivery delays under partial shipments conditions: a case study in the automotive sector

**Lorenzo Civolani\*, Matteo Gabellini\*, Alberto Regattieri\*,
Francesca Calabrese\*, Michele Ronchi \***

*\*Department of Industrial Engineering, University of Bologna, Viale del Risorgimento 4, 40136 – Bologna – Italy
(matteo.gabellini5@unibo.it, lorenzo.civolani@unibo.it, francesca.calabrese9@unibo.it, alberto.regattieri@unibo.it, michele.ronchi8@unibo.it)*

**Abstract**: In today's increasingly vulnerable supply chain landscape, the ability to anticipate risks is paramount for business survival. Of particular importance is the estimation of supplier delivery delays, especially for companies heavily reliant on outsourcing and just-in-time practices, where late deliveries can disrupt production flow and result in significant revenue loss. Recognizing this critical need, researchers have developed machine learning models to forecast supplier delivery delays. However, existing models often overlook the possibility of a single order being delivered in multiple shipments by the supplier. To address this limitation, this study thus proposes a novel multioutput regression model to deal with delivery delay predictions in presence of partial shipments conditions. The proposed model is thus built to be able to estimate four key variables for each order: the days between the planned delivery date and the date of the first partial shipment, the days between the planned delivery date and the date of the second partial shipment and the amount of quantity delivered respectively in the first e second partial shipments. An empirical investigation of the predictive accuracy reachable by the proposed approach, based on real-world data from an automotive case study, is conducted to evaluate the proposed approach's effectiveness. Moreover, the capability of the proposed approach to properly estimate the real cost impact generated by the non punctual delivery of purchased components is compared with the capability to estimate the same effect using a model not specifically designed to consider situations involving partial shipments.

**Keywords**:  Artificial intelligence, machine learning, supply chain, supply chain risk management

## 1.Introduction

In recent years, Artificial Intelligence (AI) (Cannas et al. 2023; Helo & Hao 2021) and more specifically Machine Learning (ML) (Akbari & Do 2021; Ni, Xiao & Lim 2020) have witnessed significant advancements leading to a deep transformation of established supply chain management (SCM) practices.

Simultaneously, the increased number of disruptions faced by supply chains in different sectors caused by natural disasters, geopolitical tensions, or global pandemics, has highlighted the need for more robust and adaptive risk management strategies and has led industries to rethink their risk management practices by implementing these new technologies to identify, assess and mitigate future risks (Baryannis et al. 2019; Gabellini et al. 2022, 2023; Ganesh & Kalpana 2022; Regattieri et al. 2024).

The problem of designing efficient ML algorithms in the field of Supply Chain Risk Management (SCRM) has thus assumed increasing relevance. In particular, among various risks affecting supply chains, the prediction of supplier delivery delay risks has received remarkable attention from scholars. Indeed, the widespread adoption of lean management practices and the consequent need for on-time delivery to avoid inventory stockout has motivated several researchers to propose models to anticipate possible nonpunctual delivery of purchased components.

One of the first studies addressing the problem is the one by Baryannis, Dani, et al.(2019). In their study, authors proposed a ML approach to estimate if future order will be in late or not. Moreover, in the study authors also investigated the advantages and disadvantages of adopting more or less interpretable ML models for the predictions. Later on, Cavalcante et al., (2019) proposed to adopt a similar model to estimate in a binary form the on-time delivery of suppliers and integrate these predictions in a rule base system to support the allocation of orders to suppliers based on their forecasted reliablity. Afterwards, in Brintrup et al., (2020) the advantages of introducing engineered features to enhance predictions for the problem of classifying deliveries as punctual or not has been investigated. The study of Steinberg et al., (2023) represents instead the first attempt to explicitly design a ML model not only able to identify if an order would have been delivered in late or not, but also able to estimate the exact amount of days of delay. Another approach, implementing regression models can be found in Gabellini, et al. (2024). Here, in addition, authors have considered the value of macroeconomics indicators to increase prediction accuracy. Lastly, several other recent contributions to solve the problem can be found in Bodendorf et al., (2023) and Zheng et al., (2023). In the former, authors solved the

problem of classifying future orders as on time or not leveraging not only internal to company data but also relying on macroeconomic indicators. Conversely, Zheng et al., (2023), proposed for the first time a federated learning approach to allow companies having the same suppliers to share their delivery punctuality data to increase prediction accuracy.

Although these studies provide a reference point in supplier delivery delay predictions, some gaps can however be noticed. In particular, an often-overlooked aspect when dealing with real-world delivery is the phenomenon of partial shipment. Specifically, partial shipment refers to the delivery of an order in more than one shipment (Banerjee et al. 2001; Gabellini, et al. 2024). The presence of partial shipments can thus complicate delivery delay prediction efforts, as traditional models may not adequately account for the variability introduced by partial shipments. Indeed, in the proposed literature, a single delivery punctuality value is estimated for each order, while the presence of partial shipment requires associating a number of predictions equal to the number of partial shipments to each order.

To cover this gap, this study thus proposes a novel ML-based approach for delivery delay prediction that explicitly considers partial shipments. Furthermore, an empirical investigation is conducted based on data coming from a real case study in the automotive sector to address the following research questions:

1. Which predictive accuracy can the proposed approach reach in estimating the days of delay and the delivered quantity in each partial shipment?

2. Which advantages can the proposed model lead in the estimation of the cost impact of non-punctual delivery compared to a model not specifically designed to consider partial shipments?

3. Does the proposed model require more computational time to be trained than a model not specifically designed to consider partial shipments?

The remainder of this paper is organized as follows. Section 2 outlines the proposed model and the methodology adopted to investigate its potential. Section 3 presents the empirical results. Finally, Section 4 discusses the result and concludes the paper.

## 2. Materials and Methods

The proposed model and the experimental design adopted to investigate its performance are presented in this Section.

### 2.1 Proposed model

The design of a ML model involves several steps. Among these steps the definition of the problem to solve and identification of the variables to predict, the selection of the variables adopted as input for the predictions and the specification of the ML architecture to adopt, represent some of the most essential aspects to take into consideration. For this reason, each of these aspects is detailed in the following sections.

### 2.1.1 Problem definition and selected target variables

According to the contributions stated in Section 1, the aim of the proposed model is to punctually estimate the delivery punctuality of purchased components also when partial shipments occur. Therefore, the proposed model has been designed to solve a regression problem where the exact amount of a specific variable needs to be estimated. More specifically, as multiple partial shipments can occur for each order, the proposed model has been designed to solve a multi-output regression problem. In particular, assuming that for the majority of the cases, no more than two partial shipments originate from a single order, for each order, the selected target variables predicted by the proposed models are represented by:

1. The number of days of delay or advance reported in the first partial shipment related to a specific order of a specific component

2. The number of days of delay or advance reported in the second partial shipment related to a specific order of a specific component

3. The amount of quantity delivered in the first partial shipment related to a specific order of a specific component

The quantity delivered in the second partial shipment related to a specific order of a specific component has been assumed to be an unnecessary prediction as it can be deduced by subtracting the original ordered quantity from the target variable defined at point 3.

### 2.1.2 Proposed input variables

Based on the defined problem, the following variables have been proposed to be adopted as input to predict the three targets defined in Section 2.1.1:

1. The year, month, week, and day of the week related to the planned date on which the ordered component should be delivered.

2. The quantity planned to be received for the specific component.

3. The last historical value recorded for a specific component of the three target variables defined in Section 2.1.1.

4. The supplier in charge of the specific delivery.

### 2.1.3 Proposed model architecture

A gradient-boosting machine learning architecture called CatBoost has been proposed as potential architecture for the investigated problem. The decision to utilize CatBoost, a non-linear gradient boosting model, for the predictive module was made for several reasons. While linear models are typically quicker to train, non-linear machine learning models such as CatBoost excel at capturing complex relationships and patterns, particularly in large datasets. Secondly, gradient boosting models offer higher

explainability compared to non-linear black box models like neural networks or support vector machines. Lastly, among gradient boosting algorithms, CatBoost, despite requiring longer training times than LightGBM, has demonstrated superior predictive performance when applied to tabular data (Dorogush, Ershov & Gulin 2018).

## 2.2 Experimental design

The case study selected to investigate the proposed approach and the data collected in relation to it are presented in this section. Moreover, the benchmark model against which the performance of the proposed approach has been compared, the metrics adopted for the comparison and the experimental setup adopted to perform the computation are presented.

### 2.2.1 Data collection

A real case study in the automotive sector has been selected to test the performances of the proposed model. 2248 different components supplied by 159 different suppliers over a period of 4 years where partial shipments occurred have been considered. In particular, the delivery punctuality observed in the first and in second partial shipments has been recorded for each component. Moreover, the amount of components delivered in each partial shipment has been tracked. The collected data can be available upon request.

### 2.2.2 Benchmark model

In line with the models proposed in the literature, a model able to estimate only one delivery punctuality value for each order of a specific component has been considered as the benchmark of the proposed approach. In particular, the value predicted by the benchmark approach in the presence of a partial shipment is represented by the mean of the delivery punctuality value reported in the two partial shipments. The same CatBoost algorithm adopted for the proposed model has been considered for building the benchmark model. This choice has been made to avoid differences in results related to the selection of different algorithms. Moreover, for the same reason, the same features described in Section 3.1 have been adopted as predictors also in this model.

### 2.2.3 Evaluation metrics

Two different metrics have been considered to evaluate the proposed model's performance and compare it against the benchmark.

First, two widely adopted accuracy metrics have been considered to evaluate the predictive capability of the proposed model: the Mean Absolute Error (MAE) and the Mean Absolute Percentage Error (MAPE). In particular the errors have been computed with respect of each component i:

$$MAE_i = \frac{1}{T}\sum_{t=1}^{T}|y_{it} - \hat{y}_{it}| \quad (1)$$

$$MAPE_i = \frac{1}{T}\sum_{t=1}^{T}\left|\frac{y_{it} - \hat{y}_{it}}{y_{it}}\right| * 100 \quad (2)$$

Where T is the number of deliveries recorded for the specific component i, $y_{it}$ represent the real delivery

punctuality experienced by component i in the delivery t, while $\hat{y}_{it}$ represent the delivery punctuality estimated by the proposed model for the same component i and delivery t.

In addition, a new tailored cost metrics has been introduced to assess the capability of both the proposed and the benchmark models to estimate the cost impact generated by non-punctual deliveries. Specifically, the introduced Delivery Cost Impact Error (DCIE) metric has been defined as:

$$DCIE_i = \left|\frac{\sum_{t=1}^{T} C_i^{HOLDING} * q_{it} * \hat{y}_{it}^{ADVANCE} + C_i^{SHORTAGE} * q_{ik} * \hat{y}_{it}^{DELAY}}{\sum_{t=1}^{T} C_i^{HOLDING} * q_{it} * y_{ik}^{ADVANCE} + C_{ik}^{SHORTAGE} * q_{ik} * y_{it}^{DELAY}} * 100\right|$$
(3)

Where:

- $C_i^{HOLDING}$ is the unitary holding cost of component i expressed in $[euro/pieces * days]$,

- $C_i^{SHORTAGE}$ is the unitary shortage cost of component i expressed in $[euro/pieces * days]$,

- $q_{it}$ is the quantity of component i delivered in the delivery t expressed in $[pieces]$

- $\hat{y}_{it}^{ADVANCE}$ represent the estimated amount of days of advance with which the component i has been delivered in delivery t

- $\hat{y}_{it}^{DELAY}$ represent the estimated amount of days of delay with which the component i has been delivered in delivery t

- $y_{it}^{ADVANCE}$ represent the true amount of days of advance with which the component i has been delivered in delivery t

- $y_{it}^{DELAY}$ represent the true amount of days of delay with which the component i has been delivered in delivery t

Compared with the MAE and MAPE metrics the DCIE thus allows to compare the results obtained from the proposed and the benchmark model when evaluating the capability of each model to estimate the real impact that deliveries in delay or advance generate in the presence of partial shipments. Indeed, a fair comparison was impossible to obtain based on the accuracy metrics as the two models differ in predicted output. Indeed, in presence of partial shipments, for a specific component, the benchmark model predicts the mean punctuality over the recorded partial shipments. Contrarywise, the proposed model punctually estimates the delivery punctuality and the respective percentage of the overall quantity delivered in each partial shipment. On the contrary, even if the two models estimate different values, the DCIE allows to compare in economic terms the capability of the two models to estimate which of the two models commits the lower error in estimating the real impact generated by component delivery issues. In particular, the denominator of Equation 3 always refers to the real situation where partial shipments occur. On the contrary, based on the considered model, the numerator is computed considering the values related to each partial

shipment (when implementing the proposed model) or an approximation considering the mean delay or advance over the partial shipments composing the ordered quantity (when considering the benchmark model).

### 2.2.4 Experimental set up

For each of the 2248 considered components, the same experimental setup was followed.

First, the historical records of the selected target variables and input features of the proposed model and benchmark approach have been split over time in three subsets: the training set, the validation set and the test set. An exemplification of this division is reported in Figure 1. Here, a red line illustrates the three distinct set of data for the historical record of the target variables predicted by the models. In particular Figure 1.a reports the split over the historical record of the mean days of advance or delay reported by a specific component which represents the variable predicted by the benchmark model. Contrarywise, Figures 1.b and 1.c reported the historical records related to the days of delay or advance reported in each partial shipment and the quantity respectively delivered in each shipment for the same component.



**Figure 1: Target variables split into training, validation and test subsets**

Based on this split, the training set, has thus been initially adopted to train the proposed and the benchmark models and to allow them to learn the hidden relationship between the input variables described in Section 2.1.2 and the target variables to estimate.

The validation set has been instead adopted to identify the best hyperparameters to adopt for each model. A summary of the considered hyperparameters and the research space within which the value of each hyperparameter has been searched is reported in Table 1.

**Table 1: Hyperparameters research space**

| Hyperparameters | Hyperparameters research space |
| --- | --- |
| Learning rate | $0.01 - 0.1$ |
| Depth | $4 - 8$ |
| L2 leaf regularization | 1e-6 – 1e-2 |

In particular, for both the proposed and the benchmark models, a time limit of 8 hours has been considered for the identification of the best hyperparameter values. Within this time, a Bayesian optimization strategy has been adopted to progressively sample new combinations of the investigated hyperparameters within their respective research space to find the hyperparameters combination reporting the lowest Root Mean Squared Error (RMSE) over the predictions generated for the validation set, where the RMSE has been defined as:

$$RMSE = \frac{1}{N} \sum_{i=1}^{N} \sqrt{\frac{\sum_{t=1}^{T}(y_{it} - \hat{y}_{it})^2}{T}} \ (4)$$

Here N is the number of considered components, T is the number of deliveries recorded for the specific component i, $y_{it}$ represent the real days of delay or advance experienced by component i in the delivery t, while $\hat{y}_{it}$ represent the days of delay or advance estimated by the proposed model for the same component i and delivery t.

The optimal number of iterations to consider in training the model has been identified based on an early stop strategy. More specifically, an initial number of 1000 iterations has been considered for each model. Afterward, as reported in Figure 2, by monitoring the RMSE error reported by both the models over the training and the validation set, the training of the models has been stopped when 20 consecutive iterations with no improvements in terms of the RMSE reported in the validation set have been found.



**Figure 2: Early stopping procedure**

Lastly, when the optimal value of the investigated hyperparameters have been found, both the models have been retrained considering the identified hyperparameters on both the historical data contained in the training and in

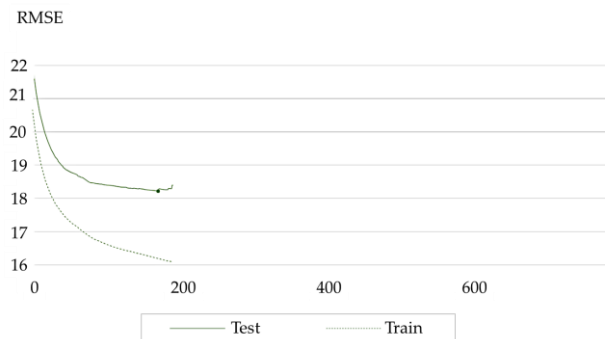the validation set and the MAE, MAPE and DCIE have been computed with respect of the predictions and the true value reported in the test set.

The optimal value of the selected hyperparameters for the proposed and the benchmark model are reported in Table 2.

**Table 2: Selected hyperparameter values**

| Hyperparameters | Proposed model | Benchmark model |
|---|---|---|
| Learning rate | 0.086 | 0.084 |
| Depth | 8 | 4 |
| L2 leaf regularization | 4.9 e-5 | 1.2 e-6 |
| Iterations | 55 | 65 |

In the study, all the experiments were performed on computer, with a processor Intel(R) Core(TM) i7-8550U running at 1.80GHz. The system utilized 16 GB of RAM. All experiments were performed with Windows 10 Pro as the operating system, and all the code has been written in Python. The CatBoost library has been adopted to implement the models, while the Bayesian optimization has been implemented through the Optuna library.

### 3. Results

In this Section first the accuracy reached by the proposed approach in terms of MAE and MAPE is presented. Afterward, a comparison between the DCIE reported by the proposed approach and the benchmark approach is made.

### 3.1 Proposed approach prediction accuracy

The boxplot in Figure 3 reports the distribution of the MAE when comparing the predictions provided by the proposed model with the real value observed over the test set for each of the 2248 considered components.



**Figure 3: Mean Absolute Error reported by the proposed approach**

According to the chart, a median value of 8.2 and 1.9 days has been observed respectively when considering the difference between the predictions reported by the proposed approach and the true delivery punctuality of the first and the second partial shipment. A median value of 1.5 pieces has been instead found between the values

estimated by the proposed approach and the true amount of pieces delivered in the first delivery. The error's interquartile spans from a minimum of 0 to a maximum of 28.5 days for the predictions related to the first partial shipment delivery punctuality. On the contrary, the interquartile related to the predictions of the second partial shipment delivery punctuality has been observed to vary between 0 to 10.7 days. Lastly the error's interquartile related to the prediction of the quantity delivered in the first partial shipment have been observed to vary from 0 to 14.8 pieces.

The boxplot reported in Figure 4, reports instead the distribution of the MAPE over the test set.



**Figure 4: Mean Absolute Percentage Error reported by the proposed approach**

According to the chart a median value of 122%, 97%, 9% has been observed respectively when considering the difference between the predictions reported by the proposed approach and the true delivery punctuality of the first partial shipment, the true delivery punctuality of the second partial shipment and the true quantity delivered in the first partial shipment. The error's interquartile has been observed to span respectively from a minimum of 1% to a maximum of 352% for the first target, from 22% to 162% for the second target and from 1% to 112% for the third target.

### 3.2 Delivery cost impact error comparison

The results of the comparison between the results produced by the proposed approach and the benchmark model in terms of DCIE are reported in Figure 5. In particular the chart reports the distribution of the variable $\Delta DCIE_i$ defined as:

$$\Delta DCIE_i = DCIE_i^{proposed} - DCIE_i^{benchamrk} \quad (5)$$

For each component i the variable $\Delta DCIE_i$ thus allow to understand if the proposed or the benchmark model represent a better estimation of the real impact generated by non punctual delivery. In particular, a negative value of the variable $\Delta DCIE_i$ indicates that the proposed approach resulted in a better estimation of the delivery cost impact compared to the benchmark approach. Conversely, a value greater than zero means that the benchmark approach resulted in a better model to estimate the impact

generated by non punctual deliveries compared to the proposed model.



**Figure 5: Delivery cost impact error difference distribution between the proposed and the benchmark approach**

According to the chart for 75 % of the considered components a negative $\Delta DCIE_i$ value has been reported. In particular, for 50 % of the considered components a $\Delta DCIE_i$ between 0 % and -20 % has been observed. On the contrary, only for the remaining 25 % positive values of the $\Delta DCIE_i$ variable has been observed, meaning that in this case the proposed approach not represent a better estimation of the real cost impact compared to the benchmark approach.

**3.3 Training time comparisons**

Lastly, the results of the comparison in terms of training time (expressed in seconds) required by the proposed model and by the benchmark approach to perform their training is reported in Figure 6.



**Figure 6: Training time comparison between the proposed and the benchmark approach**

According to the Figure no significant difference has been found in the training times required by the proposed and the benchmark approach, which both take only a few seconds to be completed.

**4. Discussion and conclusions**

Artificial intelligence and machine learning models have been recently identified as fundamental tools to proactively estimate supply chain risks. For this reason, over time, several studies have proposed models to forecast future risks. In particular, considerable attention has been recently obtained by models predicting component delivery punctuality due to the impact that nonpunctual delivery can produce in manufacturing systems.

However, while these models are able to estimate components' delivery punctuality under normal conditions, they fail to model situations where orders are not only delivered late or in advance but the originally ordered quantity is also delivered and split over multiple partial shipments. For this reason, the present study thus proposed a new model able to predict the delivery punctuality of purchased components also under partial shipment circumstances.

The accuracy of the proposed model in terms of Mean Absolute Error (MAE) and Mean Absolute Percentage Error (MAPE) has been investigated through an experimental design based on dat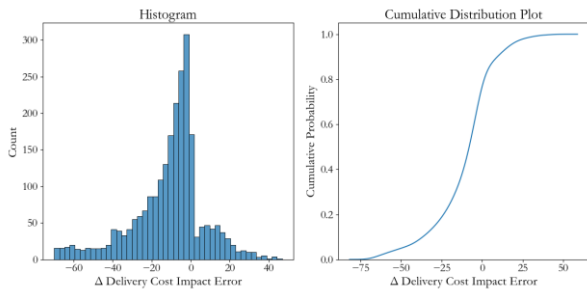a related to a real automotive case study. Moreover, a new metric, the Delivery Cost Impact Error (DCIE) has been introduced to compare the capability of the proposed approach to estimate the real impact that non-punctual delivery can generate with a benchmark model not explicitly designed to address partial shipment conditions.

In particular, considering the RQ1: Which predictive accuracy can the proposed approach reach in estimating the days of delay and the delivered quantity in each partial shipment? The empirical investigations reported that for half of the considered component error of 122%, 97%, 9% in predicting the days of delay or advance of the first and the second partial shipment, and the amount of quantity delivered in the first partial shipment has been reported by the proposed model in terms of MAPE. Results thus highlight that, while predicting the quantity delivered in each partial shipment can be easily forecasted, the estimation of the delivery delays reported in each shipment is not easy also for advanced machine learning models. However, it has also to be noticed that for 25 % of the considered components the MAPE error when estimating the delivery punctuality reported in the first and in the second partial shipments stay in a lower range spanning from 0% to and from 90% and from 15% to 85% in terms of MAPE respectively. Moreover, when considering the MAE error reported by the proposed model, it has to be considered that similar errors have also been reported in other studies (Steinberg et al. 2023).

Considering the RQ2: Which advantages can the proposed model lead in the estimation of the cost impact of non-punctual delivery compared to a model not specifically designed to consider partial shipments? The comparison performed between the benchmark model and a generic model not explicitly designed for partial shipments revealed that for 75% of the considered components, the proposed model resulted in a better estimation of the cost impact that components' non-punctual delivery can generate.

Lastly, considering the RQ3: Does the proposed model require more computational time to be trained than a model not specifically designed to consider partial shipments? No significant difference in terms of computational times has been observed between the proposed model and the benchmark. This result, thus

suggest that the cost impact estimation advantages reported by the proposed model don't come at the cost of a significantly increased computational effort to train the model.

These results need to be however considered subjected to some limitations. In particular, only one case study has been adopted to investigate the capability of the proposed approach. Moreover, the hyperparameter tuning phase has been limited to 8 hours.

Based on these limitations, future research can be thus directed to investigate the capability of the proposed approach on multiple case studies and to eventually extend the time considered in the hyperparameter tuning phases. In addition, incorporating the developed predictive models to design effective optimization models to solve typical decision-making problems arising in the context of supply chain management like supplier selection and order allocation or inventory management could represent and interesting future research direction.

## References

Akbari, M & Do, TNA 2021, 'A systematic review of machine learning in logistics and supply chain management: current trends and future directions', *Benchmarking: An International Journal*, vol. 28, no. 10, pp. 2977–3005.

Banerjee, S, Banerjee, A, Burton, J & Bistline, W 2001, 'Controlled partial shipments in two-echelon supply chain networks: a simulation study', *International Journal of Production Economics*, vol. 71, no. 1–3, pp. 91–100.

Baryannis, G, Dani, S & Antoniou, G 2019, 'Predicting supply chain risks using machine learning: The trade-off between performance and interpretability', *Future Generation Computer Systems*, vol. 101, pp. 993–1004.

Baryannis, G, Validi, S, Dani, S & Antoniou, G 2019, *Supply chain risk management and artificial intelligence: state of the art and future research directions*, International Journal of Production Research, vol. 57, no. 7, pp. 2179–2202.

Bodendorf, F, Sauter, M & Franke, J 2023, 'A mixed methods approach to analyze and predict supply disruptions by combining causal inference and deep learning', *International Journal of Production Economics*, vol. 256, p. 108708.

Brintrup, A, Pak, J, Ratiney, D, Pearce, T, Wichmann, P, Woodall, P & McFarlane, D 2020, 'Supply chain data analytics for predicting supplier disruptions: a case study in complex asset manufacturing', *International Journal of Production Research*, vol. 58, no. 11, pp. 3330–3341.

Cannas, VG, Ciano, MP, Saltalamacchia, M & Secchi, R 2023, 'Artificial intelligence in supply chain and operations management: a multiple case study research', *International Journal of Production Research*, pp. 1–28.

Cavalcante, IM, Frazzon, EM, Forcellini, FA & Ivanov, D 2019, 'A supervised machine learning approach to data-driven simulation of resilient supplier selection in digital manufacturing', *International Journal of Information Management*, vol. 49, pp. 86–97.

Dorogush, AV, Ershov, V & Gulin, A 2018, 'CatBoost: gradient boosting with categorical features support'.

Gabellini, M, Calabrese, F, Civolani, L, Regattieri, A & Mora, C 2024, 'A Data-Driven Approach to Predict Supply Chain Risk Due to Suppliers' Partial Shipments', in SG Scholz, RJ Howlett & R Setchi (eds), *Smart Innovation, Systems and Technologies*, vol. 377, Springer Science and Business Media Deutschland GmbH, pp. 227–237,

Gabellini, M, Calabrese, F, Regattieri, A & Ferrari, E 2022, 'Multivariate multi-output LSTM for time series forecasting with intermittent demand patterns', in *Proceedings of the Summer School Francesco Turco*, AIDI - Italian Association of Industrial Operations Professors,

Gabellini, Matteo, Civolani, L, Calabrese, F & Bortolini, M 2024, 'A Deep Learning Approach to Predict Supply Chain Delivery Delay Risk Based on Macroeconomic Indicators: A Case Study in the Automotive Sector', *Applied Sciences (Switzerland)*, vol. 14, no. 11.

Gabellini, M, Civolani, L, Regattieri, A & Calabrese, F 2023, 'A Data Model for Predictive Supply Chain Risk Management', *Lecture Notes in Mechanical Engineering*, pp. 365 – 372.

Ganesh, AD & Kalpana, P 2022, 'Future of artificial intelligence and its influence on supply chain risk management A systematic review', *Computers &amp Industrial Engineering*, vol. 169, p. 108206.

Helo, P & Hao, Y 2021, 'Artificial intelligence in operations management and supply chain management: an exploratory case study', *Production Planning &amp Control*, vol. 33, no. 16, pp. 1573–1590.

Ni, D, Xiao, Z & Lim, MK 2020, 'A systematic review of the research trends of machine learning in supply chain management', *International Journal of Machine Learning and Cybernetics*, vol. 11, no. 7, pp. 1463–1482.

Regattieri, A, Gabellini, M, Calabrese, F, Civolani, L & Galizia, FG 2024, 'Balancing Data Acquisition Benefits and Ordering Costs for Predictive Supplier Selection and Order Allocation', *Applied Sciences (Switzerland)*, vol. 14, no. 10.

Steinberg, F, Burggräf, P, Wagner, J, Heinbach, B, Saßmannshausen, T & Brintrup, A 2023, 'A novel machine learning model for predicting late supplier deliveries of low-volume-high-variety products with application in a German machinery industry', *Supply Chain Analytics*, vol. 1, p. 100003.

Zheng, G, Kong, L & Brintrup, A 2023, 'Federated machine learning for privacy preserving, collective supply chain risk prediction', *International Journal of Production Research*, pp. 1–18.