

Article

Help me learn! Architecture and strategies to combine recommendations and active learning in manufacturing.

Patrik Zajec^{2,†}0000-0002-6630-3106, Jože M. Rožanec^{1,2,3,†*}0000-0002-3665-639X, Elena Trajkova^{2,4}0000-0001-5342-1085, Inna Novalija²0000-0003-2598-0116, Klemen Kenda^{1,2,3}0000-0002-4918-0650, Blaž Fortuna^{2,3}0000-0002-8585-9388, Dunja Mladenec²0000-0003-4480-082X

- ¹ Jožef Stefan International Postgraduate School, Jamova 39, 1000 Ljubljana, Slovenia
² Jožef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia
³ Qlector d.o.o., Rovšnikova 7, 1000 Ljubljana, Slovenia
⁴ University of Ljubljana, Faculty of Electrical Engineering, Tržaška 25, 1000 Ljubljana, Slovenia
* Correspondence: joze.rozanec@ijs.si (J.M.R.)
† Current address: Affiliation 3
‡ These authors contributed equally to this work.

Abstract: This research work describes an architecture for building a system that guide a user from a forecast generated by a machine learning model through a sequence of decision-making steps. The system is demonstrated in manufacturing demand forecasting use case and can be extended to other domains. In addition, the system provides means for knowledge acquisition by gathering data from users. Finally, it implements an active learning component and compares multiple strategies to recommend media news to the user. Such media news aims to provide additional context to demand forecasts and enhance judgment on decision-making.

Keywords: artificial intelligence; machine learning; active learning; knowledge acquisition; explainable artificial intelligence; manufacturing; demand forecasting; smart assistant

0. Introduction

The decreased cost of sensors and connectivity [1], along with the development of the Internet of Things, Cloud Computing, Big Data Analytics, and Blockchain technologies [2] have enabled an increasing digitalization of manufacturing and the introduction of new paradigms, such as Cyber-Physical Systems (CPS) [3,4], and Digital Twins (DTs) [5–7]. Moreover, they bring extensive added value to the Industry 4.0 [8], enabling more effective operations, cost saving, and better product quality [9].

While an explosive growth of data available in the manufacturing industry has been observed [10], captured through sensors or made available from software, such as Enterprise Resource Planning (ERP) or Manufacturing Execution Systems (MES), much collective, semantic, and tacit knowledge the employees are aware of, is not digitalized. Furthermore, much of the digitalized data is not labeled, and thus no supervised learning algorithms can be applied to it. It is thus essential to identify how informative the newly collected data instances are to make good decisions regarding data management and machine learning models.

Much of the missing information can be introduced into the digital domain by asking users specific questions. Users can be queried regarding missing labels, asked for feedback on particular entries, or missing domain knowledge. The collection of locally observed collective knowledge can be achieved through a specialized solution [11,12]. The particular case of querying a user for labels given a large pool of unlabeled data is addressed by a sub-field of machine learning known as Active learning (AL) [13]. Active learning attempts to identify the most informative data instances, which are presented to the *oracle* (e.g., a human expert) asking for a label, reducing the data annotation effort. Newly labeled data is incorporated into the existing dataset and can be fed to the machine learning models. Batch machine learning models require regular deployments to make available the last trained version to the manufacturing software.



Citation: Zajec P., Rožanec J.M., Trajkova E., Novalija I., Kenda K., Fortuna B., Mladenec D. Help me learn! Architecture and strategies to combine recommendations and active learning in manufacturing.. *Preprints* 2021, 1, 0. <https://doi.org/>

Received:

Accepted:

Published:

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Active learning reduces the labeling stress posed on the user and provides a solution to the users' reticence to provide information and feedback [14]. Though, active learning alone does not solve the data labeling issue: a good user experience is key to the success of such a system [15], impacting conversion rates (amount of labeled samples) and user satisfaction (users won't abandon the feature or application) [16]. Therefore, we designed a user interface considering users' feedback can be implicit [17] or explicit. Assuming that the quality of our entries is acceptable (implicit feedback, if no other feedback is provided), we provide means to the user to signal disagreement (explicit feedback). This approach can be later used to implement co-active learning [18]. When providing recommendations to the users, candidate data instances identified by an active learning strategy do not guarantee their quality and the consequent good user experience [19]. A compromise is required to balance exploration and exploitation while delivering good results. Furthermore, we ranked the unlabeled data entries to ensure entries whose high-quality is most probable are displayed first, and those that do not meet a certain quality threshold are not shown at all. For particular cases, such as when collecting feedback on decision-making options suggested to the user, we allowed the user to provide their own input. This way, we gather additional knowledge when the options provided so far do not satisfy the user. Such knowledge can be later incorporated into the application, promoting continuous knowledge gathering and learning.

This paper evolves previous work done in [20]. The scientific contributions of this paper are twofold. First, it describes an architecture we developed to realize a system that combines semantic technologies, machine learning, and explainable artificial intelligence to provide forecasts, explanations, and contextual information while guiding users' decision-making. Second, it compares nine active learning scenarios to understand the learning versus recommendation trade-off. Then, we evaluate them implementing a prototype application and recommending four categories of media news that enhance planners' awareness in a demand forecasting setting. In addition, we describe the implementation of a knowledge-based decision-making options recommender system implemented to advise logisticians regarding transport scheduling based on demand forecasts.

The media news we recommend to the users relates to four aspects influencing the demand for automotive engine components produced by a European original equipment manufacturer selling its products worldwide. First, the demand forecasting models were trained using real-world data provided by manufacturing partners of the European Horizon 2020 project FACTLOG [21–23]. Data we used included three years of shipment information daily, a month of demand forecasts for material and clients at a daily level, feature relevance for every prediction, forecast explanations created based on those feature rankings, and decision-making options created based on demand forecasts and heuristics.

We evaluate the outcomes of the machine learning models across different active learning scenarios assessing two metrics: area under the receiver operating characteristic curve (AUC ROC) [24] and Mean Average Precision (MAP) [25]. AUC ROC is widely adopted as a classification metric due to its desirable properties, such as being threshold independent and invariant to *a priori* class probabilities. We measure AUC ROC considering prediction scores cut at a threshold of 0.5. On the other side, MAP is a popular metric in the information retrieval domain, computing the precision of the recommendation set with the size associated with the relevant item's rank. Both metrics are used to assess the performance of recommender systems [26].

The rest of this paper is structured as follows: Section 1 presents related work, and Section 2 details the architecture we designed to satisfy the requirements described above. Section 3 describes the demand forecasting use case we considered to build test the concept architecture and system. Section 4 describes the decision-making recommender system implementations, while Section 5 details the experiments and results obtained when applying active learning for media news categorization and recommendation. Finally, Section 6 provides the conclusions and outlines future work.

1. Related Work

In this section we first briefly introduce scientific literature describing demand forecasting models related to the automotive industry. We then describe related work regarding Explainable Artificial Intelligence (XAI), and conclude with an overview of scientific works related to the active learning field.

1.1. Demand Forecasting

Products' demand forecasting requires the application of different approaches conditioned by the demand characteristics. Widely adopted criteria to characterize the demand relate to the demands' lead times variance [27], the average demand interval magnitude [28], or the coefficient of variation (see Eq. 1) [29].

$$CV = \frac{\text{Demand Standard Deviation}}{\text{Demand Mean}} \quad (1)$$

Demand is closely related to the product's characteristics and is influenced by the economic context, market type, and customer expectations. Among factors affecting the demand in the automotive industry we find personal income [30], fuel prices [31,32], gross domestic product [33], inflation and unemployment rates [34,35]. This information can be collected and encoded to datasets used to train machine learning models, which learn to predict future demand based on past data.

Statistical and machine learning models were successfully applied to provide accurate car, and car components demand forecasts. Among the most frequent machine learning algorithms used to train the models we find the Support Vector Machine (SVM) [34], Multiple Linear Regressor (MLR) [36,37], and Artificial Neural Networks (ANN) [38–40]. Popular statistical forecasting methods include the autoregressive integrated moving average (ARIMA) [30,41], autoregressive moving average (ARMA) [31], and moving average models [42].

While the accuracy of the demand forecasting models is critical for their adoption, given the influence on decision-making, it is imperative to provide details on the rationale followed by the model. Such insights help the user understand the reasons behind the forecast and decide whether it can trust it or not [43]. Furthermore, it has been argued that including domain context can further aid the planners assess the forecasts' soundness, and eventually correct it before making a decision [44–46].

1.2. Explainable Artificial Intelligence

While the Industry 4.0 paradigm represents a great potential for the manufacturing industry [47], risks associated with its implementation, such as the complexity of integration or the perceived risks of novel technologies [48] must be mitigated. One such perceived risk is the difficulty of providing an intelligible explanation regarding the machine learning models' predictions. Usual reasons behind models' opaqueness are (i) the complexity of the formal structure of the model, which can be beyond human comprehension [49], or alien to human reasoning, and (ii) intentional hiding of the inner workings of the model (e.g., to avoid exposing some trade secret, or sensitive information) [50]. Research on how to provide intelligibility on the reasons behind the forecast and transparency regarding the machine learning forecasting model is known as explainable artificial intelligence [45]. Such insights and explanations increase the trust in AI models and provide additional information to assist users' decision-making.

Best practices on how to convey the insights regarding the models' reasoning process require the explanation to resemble a logic explanation [51], and take into account relevant context. Among context elements, [52] considers three related to the explainee: (i) the user profile to whom the explanation is given, (ii) the goal of the explanation, and (iii) if the explanation is either global (describes the average AI model forecast), or local (describes a specific forecast instance). Common explanation types include feature rankings, prototype (local) explanations, and counterfactual explanations. Multiple techniques were developed

to compute feature rankings, which convey information on which features exercised most influence on a given forecast (local explanation) [53–55], or forecasts in general (global explanation). Prototype explanations are data instances obtained from the train set, which are similar to the feature vector used to issue the prediction [56]. Such samples help to understand which instances most likely influenced the model learning to provide a particular forecast. Finally, counterfactual explanations provide perturbed data samples that produce a different forecasting outcome than the original data instance [57–59]. Such samples allow the user to understand what values need to be changed to change a forecast outcome. Ideally, the perturbed features correspond to actionable aspects, on which the user can be advised to take action to influence future outcomes [60].

In the context of manufacturing, XAI technologies have been tested in several scenarios such as predictive maintenance [61], real-time process management [62], and quality monitoring [63]. One of our research goals is to highlight the models' explainability in smart manufacturing processes, aligning XAI technologies with human interaction. We also aim to collect feedback on the quality of such explanations since there are few validated measurements for user evaluations on explanations' quality [64].

1.3. Active Learning

Active learning is a sub-field of machine learning that studies how to improve the learners' performance by asking questions to an *oracle* (e.g., a human annotator), under the assumption that unlabeled data is abundant, while the labels are expensive to obtain [13]. Since users are usually reluctant to provide information and feedback, AL can be used to identify a set of data instances on which the users' input conveys the most valuable information to the system [14]. While active learning in itself helps to reduce the labeling effort focusing on the data that provides new information, it has been demonstrated that explainable artificial intelligence can provide meaningful information to the user, increasing the accuracy of the labels provided [65]. Furthermore, feedback on the explanations can be used to enhance them in the future further. A framework of three components can be used to gather feedback, considering a forecasting engine, an explanation engine, and a feedback loop to learn from the users [66].

The scientific literature describes multiple approaches towards the realization of active learning [13,67]. Regarding how the unlabeled data instances are obtained, we distinguish three scenarios: (i) membership query synthesis, (ii) stream-based selective sampling, and (iii) pool-based active learning. Membership query synthesis requires some mechanism (e.g., adversarial generative sampling [68]) to synthesize new data instances for the specific label they were requested. Stream-based selective sampling assumes a stream of unlabeled data instances is available. A decision must be made for each data instance regarding whether it should be discarded or provided to the oracle for labeling. Such a decision can be made based on an informativeness measure or determining a region of uncertainty, querying the data instances within it. Finally, pool-based active learning assumes a pool of unlabeled data from which data instances are selected greedily based on an informativeness measure, which enables to rank the entire pool before selecting the best candidate data instance.

While we envision that active learning can be applied to enhance the explanations provided by XAI, and the decision-making options recommendations we provide to the users regarding manufacturing-related operations [14], in this work, we only compare different active learning strategies to classify and recommend media news to the users. We extend the approach proposed by [66] to collect feedback from forecasts, forecast explanations, media news related to demand forecasting, and decision-making options we recommend to the users. When recommending media news to the users, we evaluate our approaches against baselines described in [19]. Those baselines allow us to understand the exploration and exploitation trade-off required to learn from promising unlabeled data instances while providing good recommendations to the users.

1.4. Active Learning for Text Classification

Text classification is a procedure of assigning predefined labels to the text and is considered as one of the most fundamental tasks in natural language processing [69]. Most classical machine learning approaches follow the two steps, where in the first step (hand-crafted), features are extracted from the input texts and in the second step, the features are fed to a classifier that makes predictions. The choices of features include the bag-of-words (BoW) approach with various extensions, such as BoW with TF-IDF weighting [70], while the choices of classifiers include logistic regression and support vector machines [71]. In some tasks, such approaches can still provide competitive baselines.

To address the limitations of hand-crafted features, neural approaches have been explored, where the model learns to map the input text to a low-dimensional continuous feature vector [71,72]. Feature extraction from text can be done using the approaches, such as word2vec [73], doc2vec [74], universal sentence encoder [75], or by using transformer-based models, such as BERT [76,77] and RoBERTa [78]. In some approaches, there are multiple ways to obtain a single feature vector for the input text. E.g., this can be done, by using only the vector of a specific word from text, for example the classification token, or by averaging the feature vectors of all the words. Different techniques might yield different performances on a given task [77,79]. A neural feature extractor can be used to produce fixed feature vectors that are fed to the classifier as in the classical two-step approach, or the neural model can be trained end-to-end on the given task.

To achieve satisfying performance, text classification models need a large number of annotated examples to learn from. As manual labeling is a resource-intensive task, active learning can alleviate some of the efforts. Different feature extraction techniques, classification models and query strategies might be used [72,79]. The prediction uncertainty-based query strategies are widely adopted and used with both single model or committees [80,81] approaches. We are primarily interested in evaluating the strategies that tackle the trade-off between learning and recommendation, so we follow the conclusions from [79] to select the feature extraction method and classification model.

2. Proposed Architecture

To realize the system described in Section 0, we first drafted and iterated an architecture, which requires the following components: (see Fig. 1A):

- **Database**, stores operational data from the manufacturing plant. Data can be obtained from ERP, MES, or other manufacturing platforms;
- **Knowledge Graph**, stores data ingested from a database or external sources and connects it, providing a semantic meaning. To map data from the database to the knowledge graph, virtual mapping procedures can be used, built considering ontology concepts and their relationships;
- **Active Learning Module**, aims to select data instances whose labels are expected to be most informative to a machine learning model and thus are expected to contribute most to its performance increase when added to the existing dataset. Obtained labels are persisted to the knowledge graph and database;
- **AI model**, aims to solve a specific task relevant to the use case, such as classification, regression, clustering, or ranking;
- **XAI Library**, provides some insight into the AI models' rationale used to produce the output for the input instance considered at the task at hand. E.g., in the case of a classification task, it may indicate the most relevant features for a given forecast or counterfactual examples;
- **Decision-Making Recommender System** recommends decision-making options to the users. Recommended decision-making options can vary depending on the users' profile, specific use case context, and feedback provided in the past;
- **Feedback module**, collects feedback from the users and persists it into the knowledge graph. The feedback can correspond to predetermined options presented to the users (including labels for a classification problem) or custom feedback written by the users;

- **User Interface**, provides relevant information to the user through a suitable information medium. The interface must enable user interactions to create two-way communication between the human and the system.

The knowledge graph is a central component of the system. Instantiated from an ontology (see Fig. 1B), it relates forecasts, forecast explanations, decision-making options, and feedback provided by the users. To ensure context regarding decision-making options and feedback provided is preserved, different relationships are established. The feedback entity directly relates to a forecast, forecast explanation, and decision-making option. While a chain of decisions can exist for a given forecast, there is a need to model the decision-making options available at each stage and the sequence on which they are displayed. To that end, the decision-making snapshot entity aims to capture a list of decision-making options provided at a given point in time. A relationship between decision-making option snapshots (*followedBy*) provides information on such a sequence. For each decision-making snapshot, a *selectedOption* relationship is created to the user's selected decision-making option. A *suggestsActionFor* relationship is created between the forecast entity and entities that correspond to the first decision-making options displayed for that particular forecast. Since the decision-making options are linked to decision-making option snapshots and preserve a sequential relationship, all decision-making options can be traced back to the forecast that originated them.

3. Use Case

Demand forecasting is a key component of supply chain management since it directly affects production planning and order fulfillment. Accurate forecasts enable operational and strategic decisions regarding manufacturing and logistics for deliveries. We developed a model to forecast demand on a material and client level daily. The model was trained on three years of data for 516 time-series corresponding to 279 materials and 149 clients of a European automotive original equipment manufacturer's daily demand. We used a subset of demand forecasts to evaluate the application. We generated forecast explanations using the LIME library [53]. We implemented two strategies for decision-making options recommendations, that allow to select a new transport or chose among existing ones. The first one consisted of a set of heuristics that satisfy certain criteria (e.g., have enough capacity to satisfy the expected demand for a given client), while the second one was a knowledge-based recommender. To enhance the context understanding related to demand forecasting, we display media entries for predetermined topics (*Automotive Industry*, *Global Economy*, *Unemployment*, and *Logistics*) obtained from a media event retrieval system for that day. Media events are queried based on a set of keywords. We developed machine learning models to classify media entries as interesting or not to the users and then gather labels from the users for new media entries.

Finally, we developed a user interface to display forecasts, forecast explanations, media news, and decision-making options (see Fig. 2). In the user interface, we identify five distinct parts:

- A **Media news panel**: displays media news regarding the automotive industry, global economy, unemployment, and logistics. The user can provide explicit feedback on them (if they are suitable or not), acting as an oracle for the active learning classifier. Once feedback is provided, a new piece of news is displayed to the user.
- B **Forecast panel**: given the date and material, it displays the forecasted demand for different clients. For each forecast, three options are available: edit the forecast (providing explicit feedback on the forecast value), display the forecast explanation, and display the decision-making options. The lack of editing on displayed forecasts is considered implicit feedback approving the forecasted demand quantities.
- C **Forecast explanation panel**: displays the forecast explanation for a given forecast. Our implementation displays the top three features identified by the LIME algorithm as relevant to the selected forecast. If users consider that some of the displayed

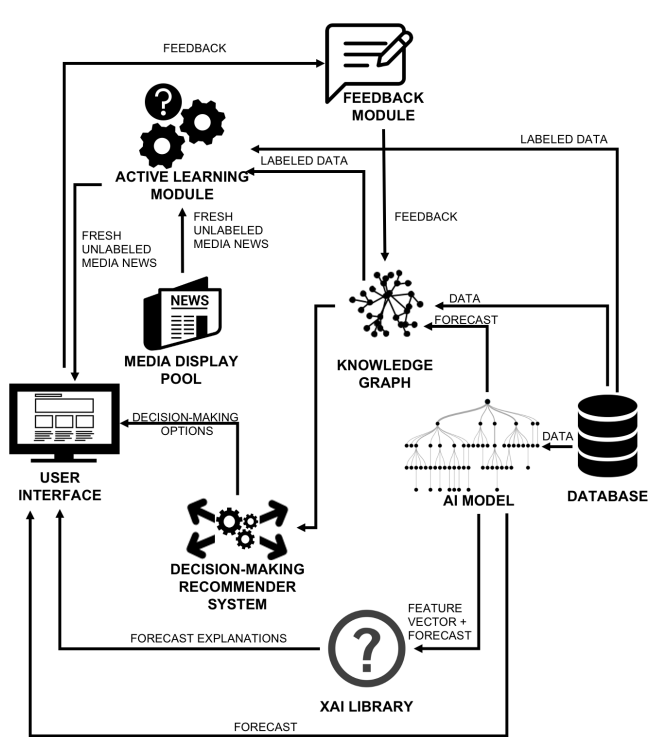


Fig. 1A

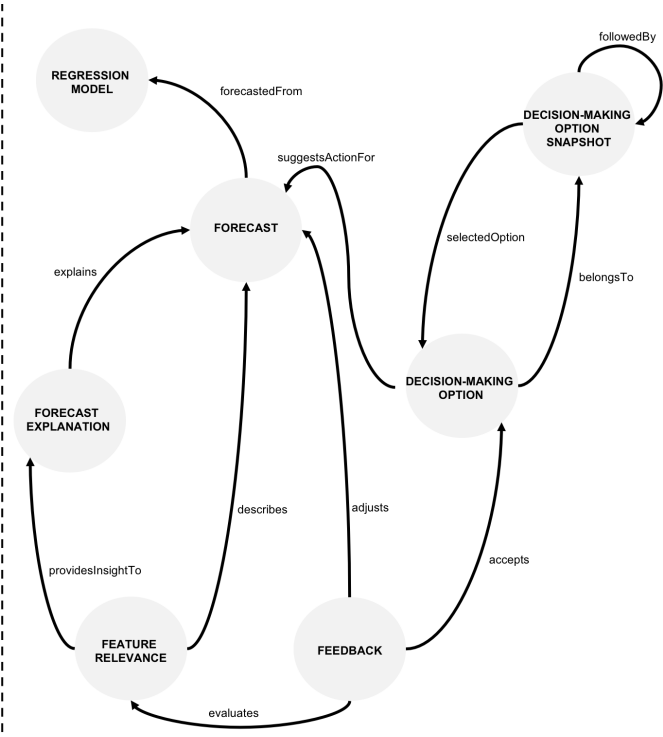


Fig. 1B

Figure 1. Fig. 1A displays a diagram of the system components and their interaction. Fig. 1B shows the main ontology concepts we considered, and their relationships.

- features do not explain the given forecast, they can provide feedback by removing it from the list.
- D **Decision-making options panel:** displays possible decision-making options for a given forecast or step in the decision-making process. In particular, the decision-making options relate to possible shipments. If no good option exists, the user can create its own.
 - E **Feedback panel:** gathers feedback from the user to understand the reasons behind the chosen decision-making option. While some pre-defined are shown to the user, we always include the user’s possibility to add their reasons and enrich the existing knowledge base. Furthermore, such data can be used to expand feedback options displayed to the users in the future.

We implemented two decision-making options recommender systems for this research work: one based on heuristics and a knowledge-based recommender system. We describe both in Section 4.

4. Decision-making options recommendation

Demand forecasts influence decision-making on a wide variety of scenarios: from raw material orders to workers hiring and upskilling to logistics arrangements to meet the required deadlines. Decision-making recommender systems can alleviate such decision-making by suggesting to the user appropriate actions based on the projected demand. In particular, we implemented a decision-making options recommender system considering the logistics use case. We consider two possible scenarios. The first scenario refers to the user who schedules a new transport for a given demand, material, client, and date. Here, the decision-making options are the possible transports, differing in transport type, delivery time, and price. The second possible scenario relates to the user who decides to use an existing transport. Here each decision-making option selects one of the existing

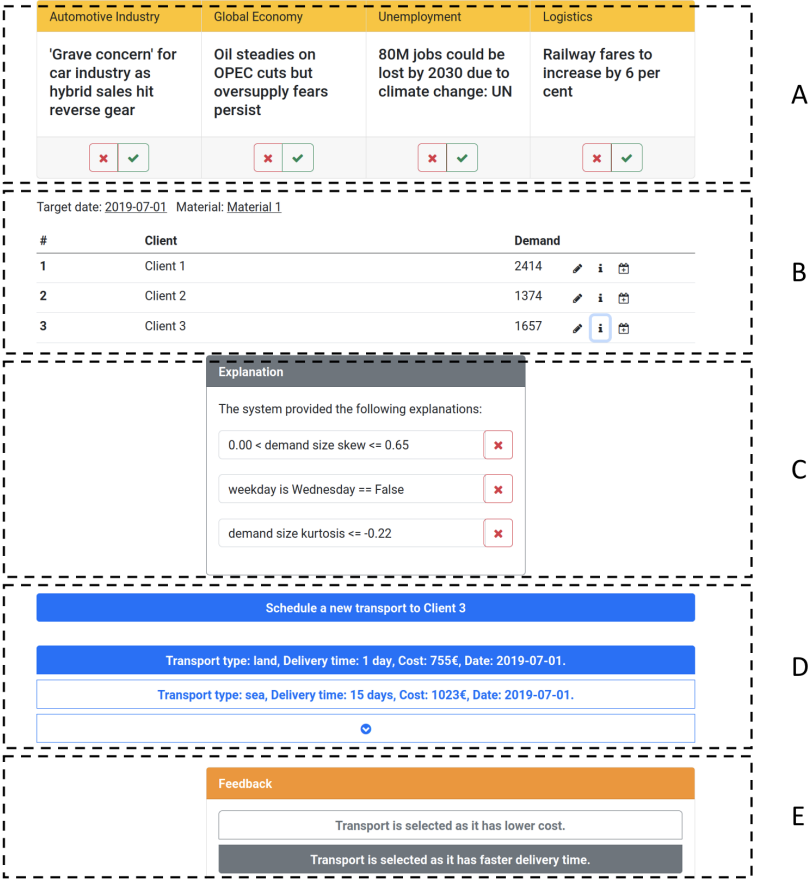


Figure 2. User interface, displaying contextual media news, forecasts, forecast explanations, and recommended decision-making options.

transports. In both steps, the recommendation module ranks the decision-making options from most to least relevant.

We developed two recommendation strategies: a heuristic-based and a knowledge-based approach. The heuristic-based recommender system follows simple rules, hand-crafted either by the domain expert or simply by the system’s developer based on his incomplete knowledge about the problem. At each step, the user should have the possibility to select any of the possible options regardless of their ranking. Such a system has no learning capacity, and therefore has little potential to improve the users’ experience. The recommendation quality directly depends on the quality of the designed rules. An example of such a heuristic rule is consistently ranking the transports according to the price or keeping only existing transports delivering in the client’s proximity and ranking them according to the remaining capacity.

The knowledge-based approach provides recommendations based on the feature vectors’ similarity to a target vector describing users’ requirements. To that end, each decision-making option at the given step is represented as a vector v . The representation captures all necessary information for the ranking, encoding the context up to the current step, the corresponding decision-making option, and its relation to all other possible decision-making options (the decision-making options snapshot). The model assigns the relevance score to each option based on v . The ranking is determined by sorting the scores from highest to lowest.

The representation and the underlying model should be expressive enough to cover the scenarios encountered in the use case. As with the heuristic-based strategy, domain knowledge heavily influences the design of features, but the content-based strategy provides greater flexibility. The features directly capture the context, which in our use case

includes the forecasted demand, date, material, and client; the decision-making option, which in the case of scheduling a new transport includes the transport type, time of delivery, capacity, and price; and the relation of the decision-making option to the whole set of available options to capture, how this option is different from others and why should it be preferred.

Among the constraints of our recommender system, we must mention that we had no data regarding the physical characteristics of each product we created demand forecasts for. In addition, while we had no information regarding the specific addresses of the clients ordering such products, we had information of the destination country. To mitigate these constraints, we collected pricing and delivery time information for air, land, and sea shipments considering single standard forty feet containers from Slovenia to fourteen countries. Such data was retrieved from two specialized web-pages¹. Finally, given the application was not deployed to a production environment yet, we lack data regarding logisticians' interaction and choices, which would enable recommender systems' performance evaluation. We envision that more complex models can be developed in the future once data regarding users' interaction with decision-making options is obtained.

5. Active Learning for Media News Categorization and Recommendation

When providing a demand forecast and the explanation that conveys an intuition regarding the reasons behind the forecast, the user can be interested in getting media news on events that can influence demand. In particular, when forecasting engine parts for the automotive industry, the user can be interested in news regarding the automotive industry, the global economy, unemployment, or logistics. While media news can be retrieved from some news intelligence platforms, keywords based queries can issue many false positives. It is thus imperative to develop a recommender system capable of discriminating and prioritizing good quality news over those considered false positives. Furthermore, it is desired that such a model improve the quality of discrimination over time and require as little manual labeling effort as possible. To realize this, we built a set of active learning binary classifiers, each one informing if the media news considered does fit into a specific media news category or not. We consider the end-user is at the same time the news consumer and the active learning's *oracle*, providing feedback regarding unlabeled instances. In our design, we display the news and collect feedback regarding them in the same user interface. This poses an exploration-exploitation dilemma since the same user interface space must be optimized to provide high-quality media news balancing between those entries where high confidence on the category exists but provide little additional information to the existing dataset, and those entries where the confidence is lower, but can provide a higher degree of novelty to the dataset[82]. In particular, each day, at most, ten pieces of news per each of the four categories are displayed to the user. The user can then provide positive or negative feedback (label) regarding each piece of news. The news should be informative for the system as the goal is to achieve good classification performance as soon as possible. On the other hand, the displayed news events should also be relevant so that the system is usable to the users after the first few iterations. The set of displayed news events on each day should therefore balance the learning vs. recommendation (exploration vs. exploitation). In this research, we do not deal with the cold-start problem since we consider it can be mitigated by pre-training the models with a set of manually annotated instances before starting the active learning dynamics. We have evaluated nine strategies (see Table 1), balancing learning and recommendation.

Different measures can be used to measure the classification certainty of the model, which is needed in the 5 out of 9 strategies presented in the table 1. We use the uncertainty of classification which is defined for a single sample x as $U(x) = 1 - \max_y P(y|x)$ where higher value of $U(x)$ means higher uncertainty. In the case of the SVM model, the distance

¹ We collected data regarding pricing and shipment time from World Freight Rates (<https://worldfreightrates.com/freight>), and SeaRates (<https://www.searates.com/freight>). We retrieved the data between July 12th and July 16th 2021.

STRATEGY	DESCRIPTION
Random	Selects the k random instances at each step.
Uncertain	Selects k instances with highest uncertainty score at each step.
Certain	Selects k instances with lowest uncertainty score, that is, most certain examples.
Positive uncertain	Select at most k instances that were labeled as positive by the classifier and have the highest uncertainty scores.
Positive certain	Select at most k instances that were labeled as positive by the classifier and have the lowest uncertainty scores.
Positive certain and uncertain	Select at most $k/2$ positive points with lowest and at least $k/2$ points with highest uncertainty score.
Alpha trade-off ($\alpha = 0.5, 0.75, 1.0$)	We adapt the strategy proposed by [19]

Table 1: Active learning and recommendation strategies.

to the separating hyper-plane is an indicator of uncertainty, with the example having the lowest distance being most uncertain [83].

Strategy *Uncertain* straightly implements the uncertainty assumption that labels of the instances with the highest classification uncertainty are the most informative. It solely focuses on learning as such instances tend not to be the most relevant for the recommendation. The *Random* strategy is included as a baseline, and so is the *Certain* strategy, which only selects the least uncertain instances whose labels should provide the most negligible value for the system according to the uncertainty assumption. To also address the recommendation, the *Positive uncertain* strategy selects the instances labeled as positive by the model as this already signals that the instance is likely to be relevant for the recommendation. At the same time, it might still provide some value for learning due to uncertainty. On the other hand, the *Positive certain* strategy selects only the positive instances. Therefore, it ranks them according to the certainty, which should highly favor the recommendation while providing little value for learning. The *Positive certain and uncertain* strategy tries to include both recommendation and learning by following the *Positive certain* strategy for the first $k/2$ instances (or less if there are not enough positive instances) to provide relevant recommendations and next following the *uncertainty* strategy to select at least the $k/2$ instances relevant for learning.

The *Alpha trade-off* strategy is adapted from [19] and has a parameter α , used to control between learning and recommendation. It selects the instances according to the formula

$$x_\alpha = \arg \min_x |P_\alpha - P(y = 1 | x)|$$

with P_α being the $(100\alpha)^{th}$ percentile of the distribution of predictive probabilities of positive class induced on the pool of new examples. For example $P_{0.5}$ equals to the median probability and $P_{1.0}$ equals to the maximum probability of the positive class assigned to an instance from the pool. According to [19], $\alpha = 0.5$ selects the instance with highest uncertainty from the pool and thus favours learning while $\alpha = 1.0$ selects most certainly positive instance and thus favours recommendation. Setting $\alpha = 0.75$ could therefore provide a trade-off between learning and recommendation. The k instances closest to P_α are selected to form the pool.

5.1. Active Learning Experiments

The active learning experiments were performed on a dataset of media news events classified into four categories: *Automotive Industry*, *Global Economy*, *Unemployment*, and *Logistics*. The dataset was manually annotated by three human annotators, based on the

specific keywords used in each category to retrieve them (see Table 2). The media news events were retrieved daily for a period of six months (from July 2019 to December 2019) from *Event Registry* [84], a well-established media events monitoring platform that has monitored mainstream media since 2014. The first month of the dataset was reserved for training the initial version of the models and for tuning the models' hyperparameters. The last month of the dataset was reserved for testing the classification performance of the models at each active learning step. The remaining data was used to execute the active learning experiments and evaluate the recommendation performance at each step.

We executed the following procedure (see Fig. 3). For each day, we retrieved all available events for that given day, and for each media news entry, we assessed whether it should be displayed to the user to gather feedback (label the instance) or not. Such decision was made based on a strategy (see Table 1) that considered how informative the news entries were to the existing dataset, and their quality towards the target category, given the requirement that the events should be both relevant for the user (recommendation quality) and informative for the model (improvement of classification). For each day, we selected at most k events, which were then shown to the user. We set $k = 10$, based on the median number of events per day, and acknowledging it is a common practice to query a fixed number of instances at each step according to the literature[79]. Once the media entry was displayed to the user, it was incorporated into the existing dataset if it provided an annotation. The models were retrained once a day, incorporating newly labeled instances into the dataset.

There are cases where the model can recommend less than k events. For example, this could be due to not enough events for a particular category exist that day or that only $k' < k$ of them are relevant or need a label. Thus, fewer events of that category are displayed to the user.

We use a separate test set to measure the active learning models' classification performance to evaluate the models. In contrast, the recommendations' quality is measured at each step of active learning using the gold labels of the displayed news events. To measure the models' discrimination power, we adopted the AUC ROC, a widely used metric due to its invariance to *a priori* class probabilities. On the other hand, to measure the quality of the recommendations, we adopted the MAP metric, which computes the precision of the recommendation set, and is not affected by the number of entries considered in each particular case (the desired property when $k' < k$ media news events are shown to the user).

Only the titles of news events were considered in classification. We experimented with three text representation techniques: TF-IDF weighted BoW representation, which is a classical representation technique used for text classification and serves as a strong baseline in our experiments; an average of token embeddings from the RoBERTa model² which proved to be most effective for text classification based on the results obtained by [79], and representations obtained from the Universal sentence encoder [75]³.

We use three different classification models in the single model setting, namely logistic regression, support vector machine, and random forest. The selection of the models follow the related work [72,79] where the SVM model was identified as a frequent choice for active learning for text classification. The models were also selected based on their fast training since we trained them from scratch for each AL iteration.

5.2. Results

In this section, we present the results we obtained when conducting experiments regarding different AL strategies. Strategies and models were evaluated in the AL setting by following the procedure explained in section 5.1. We report the classification performance as the AUC ROC score obtained in the last iteration of active learning, while recommendation

² We have used the pre-trained version of "RoBERTa-base" model implemented in the Huggingface library [85].

³ We have to use the model available at <https://tfhub.dev/google/universal-sentence-encoder/4>

CATEGORY	KEYWORDS	# INSTANCES	MISSING DATA	MEPD
(A) Automotive Industry	car sales demand, new car sales, vehicle sales, car demand, automotive industry	3865	10 days	20
(B) Global Economy	global GDP projection, global economic outlook, economic forecast	853	29 days	5
(C) Unemployment	unemployment rate, unemployment numbers, unemployment report, employment growth, long-term unemployment	3801	8 days	22
(D) Logistics	logistics, maritime transport, railroad transport, freight, cargo transport, supply chain	28231	0 days	133

Table 2: Active learning dataset categories, keywords used to query them, the number of instances per category, the number of days without entries for a given category, and the median of events per day for that category (MEPD).

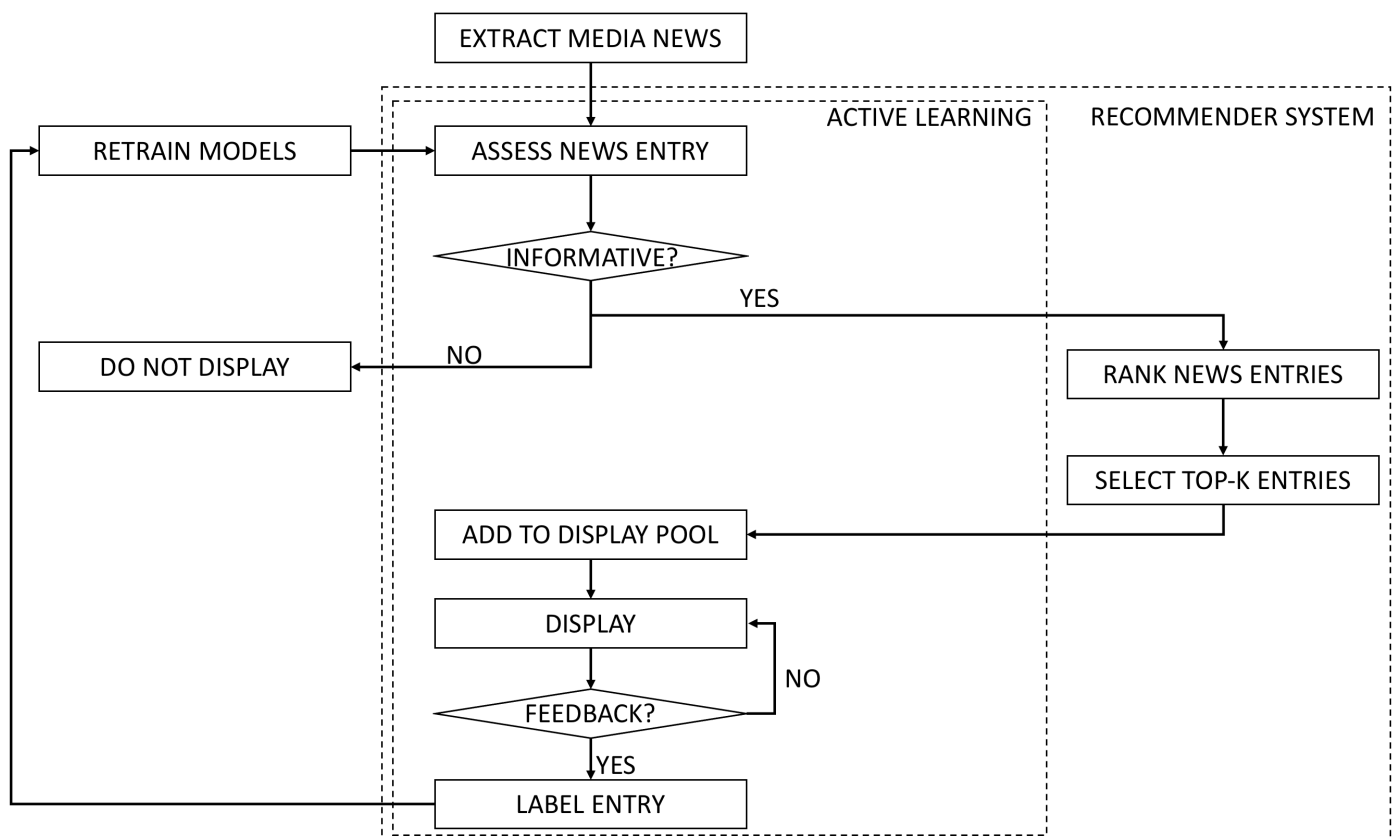


Figure 3. Fluxogram, showing how active learning and recommendations are implemented for media news event entries.

performance is reported with the MAP metric for all iterations. We provide the results of all experiments in the Tables A1, A2 and A3, in Appendix ???. Further, we compare different active learning strategies to determine the most successful in tackling the learning versus recommendation trade-off.

5.2.1. Evaluating the classification baselines.

Before conducting the experiments, we established a baseline by training multiple supervised machine learning models on all available labeled data, excluding the test set. In the baseline, we also included a fine-tuned RoBERTa model. This set of models aims to understand the maximum expected performance achieved with this dataset and its features. We report the baseline AUC ROC scores in Table 3.

From the baseline results shown in Table 3, we observe that RoBERTa model achieves the best or at least competitive performance on all but a single dataset. This is expected as fine-tuned language models are known to achieve state-of-the-art results on many text

MODEL	REPRESENTATION	A	B	C	D
LR	TF-IDF	0.8575	0.8592	0.9856	0.9456
	RoBERTa	0.8788	0.8681	0.9769	0.9297
	USE	0.8654	0.8681	0.9875	0.9195
SVM	TF-IDF	0.8639	0.8744	0.9846	0.9494
	RoBERTa	0.8889	0.8702	0.9693	0.8916
	USE	0.8828	0.8920	0.9799	0.9314
RF	TF-IDF	0.8506	0.8345	0.9733	0.8987
	RoBERTa	0.8720	0.8850	0.9179	0.8235
	USE	0.9197	0.8756	0.9854	0.8899
Fine-tuned RoBERTa	RoBERTa	0.8854	0.9081	0.9865	0.9531

Table 3: Classification performance of the models trained on all labeled examples excluding the test set. Best score for each dataset is shown in bold. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

classification tasks. Still, we can observe that the performance of the second-best model on each dataset is very close, thus providing a good alternative to the RoBERTa model in our case since other models usually require less time to train.

Fine-tuning the RoBERTa model is shown to be almost always better than using fixed RoBERTa representations with a classifier on our datasets. There is no clear winner in terms of representations, although universal sentence encoder (USE) appears to be a strong competitor (if not better) to RoBERTa based representations recommended by [79].

An unexpected finding was that models based on the TF-IDF-based representations achieve very competitive performance. Namely, on text classification tasks, TF-IDF-based models usually lag in performance behind neural-based approaches.

5.2.2. Evaluating the classification performance of AL strategies.

As an aggregation of the results from Table 1, we report the mean value and standard deviation, aggregated over models and representations on strategy and dataset level, for AUC ROC score in the Table 4. This gives us insight into the actual classification performance of the strategies on each of the datasets.

We observe little difference in final classification performance among the strategies in Table 4, although they have many different policies for selecting the instances. For example, the *Uncertain* and *Certain* strategies favor different (and, in a sense, complementary) subsets of instances while their performance appears not to differ much.

As we aim to find the strategies suitable for learning regardless of the model and representation used, we further compare the classification performance of the strategies across all datasets. First, we group the results by model, representation, and dataset. Then, inside each group, we sort and rank the strategies by their AUC ROC score. Finally, we report the mean rank for each strategy in the table 5. Additionally, for each active learning strategy, we compute the mean AUC ROC ratio towards the best strategy in the group (see Table 5). The mean rank gives us the ordering of the strategies. Finally, we determine the significance of differences between them using the Wilcoxon signed-rank test [86] on AUC ROC scores from all experiments, at a p-value = 0.005.

According to the results from Table 5 there is little difference between the best seven strategies in terms of mean rank. Furthermore, we have observed no significant difference among those strategies. We attribute this result to the large enough number of queried instances at each step ($k = 10$ in our experiments) which, for our datasets, allows us to cover a diverse set of instances regardless of the instance selection strategy. We observed, however, a significant difference between the top seven strategies and the *Positive certain* and *Positive uncertain* strategies. We attribute this difference to the two strategies limiting only to the instances with a positive label assigned by the model, which might noticeably limit the labeled set obtained during active learning. In comparison, other strategies always request the label for k instances at the given step.

STRATEGY	A	B	C	D
Random	0.8600 ± 0.0263	0.8706 ± 0.0255	0.9569 ± 0.0327	0.8643 ± 0.0511
Uncertain	0.8502 ± 0.0339	0.8684 ± 0.0215	0.9641 ± 0.0259	0.8867 ± 0.0362
Certain	0.8548 ± 0.0276	0.8701 ± 0.0237	0.9514 ± 0.0472	0.8590 ± 0.0510
Positive uncertain	0.8191 ± 0.0370	0.8467 ± 0.0336	0.9385 ± 0.0338	0.8693 ± 0.0361
Positive certain	0.8116 ± 0.0474	0.8490 ± 0.0307	0.9436 ± 0.0295	0.8680 ± 0.0419
Positive certain and uncertain	0.8438 ± 0.0332	0.8718 ± 0.0219	0.9667 ± 0.0254	0.8853 ± 0.0349
Alpha trade-off ($\alpha = 0.5$)	0.8540 ± 0.0328	0.8719 ± 0.0204	0.9579 ± 0.0352	0.8618 ± 0.0495
Alpha trade-off ($\alpha = 0.75$)	0.8478 ± 0.0306	0.8714 ± 0.0207	0.9630 ± 0.0318	0.8679 ± 0.0415
Alpha trade-off ($\alpha = 1.0$)	0.8414 ± 0.0437	0.8733 ± 0.0188	0.9649 ± 0.0250	0.8800 ± 0.0423

Table 4: Mean AUC ROC score with standard deviation, aggregated over used models and representations, for each strategy, reported on each of the four datasets. Best score for each dataset is shown in bold. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

STRATEGY	MEAN RANK	MEAN RATIO TO BEST
Alpha trade-off ($\alpha = 1.0$)	3.5000	0.9513
Uncertain	3.5833	0.9539
Positive certain and uncertain	3.6111	0.9534
Alpha trade-off ($\alpha = 0.75$)	4.6389	0.9487
Alpha trade-off ($\alpha = 0.5$)	4.7500	0.9476
Random	4.8056	0.9493
Certain	5.1944	0.9449
Positive certain	7.2222	0.9278
Positive uncertain	7.6944	0.9283

Table 5: AL strategies sorted according to mean rank of AUC.

5.2.3. Evaluating the recommendation performance of AL strategies.

To evaluate one aspect of the strategies' recommendation performance, we aggregate the results from Table 2 and report the mean value and standard deviation, aggregated over models and representations on strategy and dataset level, for MAP score in Table 6. MAP enables us to quantify, for each dataset, the strategies' performance on how accurate the recommended entries are while penalizing their ordering within the *topK* entries.

We observe that the performance of strategies that focus on the positively labeled instances, such as *Positive certain* or *Alpha trade-off* ($\alpha = 1.0$), far exceeds the performance of uncertainty focused strategies, such as *Uncertain* or *Alpha trade-off* ($\alpha = 0.5$). This is especially evident on the strongly imbalanced datasets C and D, where there is a large number of negatives, that is, irrelevant news events. A large number of negatives also explains the poor performance of *Certain* strategy, as classifying negatives appear to be more certain.

Further, to find the strategies which are good in terms of MAP score regardless of the model and representation used, we compare them across all datasets by following the same procedure as in Table 5 for the classification performance. The metric under consideration is not the AUC ROC but the MAP score in this particular case. Results are reported in Table 7, where the mean rank is used to order the strategies. We determine the significance of differences between the strategies using the Wilcoxon signed-rank test on MAP scores from all experiments, at a p-value = 0.005.

We can observe that *Positive certain* strategy achieves significantly better performance than others. Moreover, despite showing worse classification performance, according to the results from Table 5, and thus yielding less capable classification models, it displays the most relevant instances to the user at each step. However, it has to be noted that such a strategy displays much fewer instances than others and thus might miss many relevant recommendations, achieving low recommendation recall. The *Alpha trade-off* ($\alpha = 1.0$),

STRATEGY	A	B	C	D
Random	0.2729 \pm 0.0070	0.4736 \pm 0.0138	0.1218 \pm 0.0129	0.0345 \pm 0.0072
Uncertain	0.4560 \pm 0.0701	0.5332 \pm 0.0184	0.4790 \pm 0.0928	0.3098 \pm 0.0651
Certain	0.1697 \pm 0.0439	0.4130 \pm 0.0120	0.0368 \pm 0.0207	0.0108 \pm 0.0010
Positive uncertain	0.5762 \pm 0.0305	0.6692 \pm 0.0305	0.8566 \pm 0.0830	0.6642 \pm 0.2601
Positive certain	0.6651 \pm 0.0483	0.7253 \pm 0.0235	0.8862 \pm 0.0701	0.6678 \pm 0.2494
Positive certain and uncertain	0.6587 \pm 0.0413	0.6772 \pm 0.0153	0.6092 \pm 0.0370	0.3761 \pm 0.0330
Alpha trade-off ($\alpha = 0.5$)	0.3238 \pm 0.0149	0.6476 \pm 0.0106	0.0898 \pm 0.0066	0.0212 \pm 0.0021
Alpha trade-off ($\alpha = 0.75$)	0.5681 \pm 0.0219	0.6677 \pm 0.0174	0.1236 \pm 0.0131	0.0267 \pm 0.0057
Alpha trade-off ($\alpha = 1.0$)	0.6915 \pm 0.0430	0.6758 \pm 0.0185	0.6122 \pm 0.0372	0.3826 \pm 0.0419

Table 6: Mean MAP score with standard deviation, aggregated over used models and representations, for each strategy, reported on each of the four datasets. Best score for each dataset is shown in bold. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

STRATEGY	MEAN RANK	MEAN RATIO TO BEST
Positive certain	1.4167	0.8549
Alpha trade-off ($\alpha = 1.0$)	2.5556	0.7030
Positive uncertain	3.1389	0.7991
Positive certain and uncertain	3.2778	0.6905
Alpha trade-off ($\alpha = 0.75$)	5.6389	0.4418
Uncertain	5.6667	0.5267
Random	7.1667	0.2845
Alpha trade-off ($\alpha = 0.5$)	7.1667	0.3462
Certain	8.9722	0.2028

Table 7: AL strategies sorted according to mean rank of MAP.

Positive uncertain and *Positive certain and uncertain* strategies follow with significantly worse performance. We can further observe a drop in performance after the first four strategies focused on positively labeled instances. The performance of *Alpha trade-off* ($\alpha = 0.75$), which is meant to balance between the learning and recommendation, is not significantly different from the performance of *Uncertain* strategy. The *Random*, *Alpha trade-off* ($\alpha = 0.5$) and *Certain* strategy follow, again all with significantly worse performance, with *Certain* strategy being significantly the worst-performing.

Another relevant dimension of recommender systems' performance is the recall. Recall evaluates how many of the relevant instances were actually recommended and displayed to the user. While MAP score measures how many of the k (or less) displayed instances are relevant and whether the relevant instances are shown first, the recall score measures the ratio of shown relevant instances versus all relevant instances. We aggregate the results from Table 3 and report the mean value and standard deviation, aggregated over models and representations on strategy and dataset level, for Recall score in Table 8.

The *Alpha trade-off* ($\alpha = 1.0$) strategy achieves the best mean recall score on all datasets and is closely followed by the *Positive certain and uncertain* strategy. Although the *Positive certain* strategy was ranked first according to the MAP score (see Table 7), it is evident that it performs well in terms of precision by trading the recall.

To further compare the strategies regardless of the model and representation used, we follow the same procedure as for the MAP score (see Table 7). Results are reported in Table 9, where the mean rank is used to order the strategies. The significance of differences between the strategies is determined using the Wilcoxon signed-rank test on recall scores from all experiments, at a p-value = 0.005.

We found the *Alpha trade-off* ($\alpha = 1.0$) displayed the best performance with significant difference to the second best, *Positive certain and uncertain* strategy. The *Uncertain* strategy follows with significantly better results than the remaining strategies. It can be observed

STRATEGY	A	B	C	D
Random	0.5382 \pm 0.0129	0.9657 \pm 0.0060	0.4544 \pm 0.0401	0.0810 \pm 0.0158
Uncertain	0.6944 \pm 0.0537	0.9807 \pm 0.0044	0.8820 \pm 0.0691	0.5460 \pm 0.0577
Certain	0.4061 \pm 0.0510	0.9555 \pm 0.0051	0.1564 \pm 0.0403	0.0121 \pm 0.0036
Positive uncertain	0.6768 \pm 0.0482	0.7726 \pm 0.0371	0.5346 \pm 0.1009	0.2464 \pm 0.0487
Positive certain	0.6950 \pm 0.0584	0.7663 \pm 0.0441	0.5322 \pm 0.1033	0.2454 \pm 0.0474
Positive certain and uncertain	0.8208 \pm 0.0280	0.9926 \pm 0.0033	0.9512 \pm 0.0402	0.5966 \pm 0.0524
Alpha trade-off ($\alpha = 0.5$)	0.4206 \pm 0.0234	0.9518 \pm 0.0056	0.1414 \pm 0.0288	0.0279 \pm 0.0078
Alpha trade-off ($\alpha = 0.75$)	0.7119 \pm 0.0200	0.9804 \pm 0.0056	0.2189 \pm 0.0489	0.0495 \pm 0.0173
Alpha trade-off ($\alpha = 1.0$)	0.8546 \pm 0.0258	0.9929 \pm 0.0031	0.9526 \pm 0.0340	0.6018 \pm 0.0664

Table 8: Mean recall score with standard deviation, aggregated over used models and representations, for each strategy, reported on each of the four datasets. Best score for each dataset is shown in bold. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

STRATEGY	MEAN RANK	MEAN RATIO TO BEST
Alpha trade-off ($\alpha = 1.0$)	1.2778	0.9434
Positive certain and uncertain	1.7778	0.9318
Uncertain	3.3333	0.8580
Positive certain	5.4444	0.6072
Alpha trade-off ($\alpha = 0.75$)	5.4722	0.5152
Random	5.8333	0.5345
Positive uncertain	5.9167	0.6047
Alpha trade-off ($\alpha = 0.5$)	7.9167	0.4002
Certain	8.0278	0.3953

Table 9: AL strategies sorted according to mean rank of recall.

from the Table 8 that the score of *Uncertain* strategy is in range with the scores of the best two strategies on all datasets, and it even does not decrease as much as the score of others, worse-performing strategies, on dataset *D*. It might be that the uncertain instances are frequently from the positive class in our datasets. Next, we can observe the decrease in performance with the differences between the following four strategies not being significant, and the *Alpha trade-off* ($\alpha = 0.5$) and *Certain* strategies at the tail.

Through the classification and recommendation results, we have evaluated how well each strategy performs in terms of learning and recommendation and how does its performance compares to others. As just a single strategy is implemented in the active learner, it has to be such that it best balances the learning and recommendation for the best user experience. Based on the results, we consider the *Alpha trade-off* ($\alpha = 1.0$) strategy to be the best choice, followed by the *Positive certain and uncertain* strategies. The classification results (see Table 5) showed no statistically significant difference in performance between the best strategies. Although based on the precision of recommendation (MAP score) results (see Table 7), the *Positive certain* is the best performing strategy, it only performs well in one aspect of recommendation and ignores the recall. Both *Alpha trade-off* ($\alpha = 1.0$) and *Positive certain and uncertain* are second-tiers in terms of MAP score with *Alpha trade-off* ($\alpha = 1.0$) strategy being slightly better, while they rank first and second in terms of recommendation recall.

6. Conclusions and Future Work

The current work presents an architecture designed to acquire and encapsulate complex knowledge using semantic technologies and artificial intelligence. The system was instantiated for the demand forecasting use case in the manufacturing domain, using real-world data from partners from the EU H2020 projects STAR and FACTLOG. In particular, the system provides forecasts and explanations, enriches users' domain knowledge through

a set of media news, recommends decision-making options, and collects users’ feedback. Furthermore, the system uses active learning to reduce manual labeling effort and better discriminates between good and bad media news event entries reporting relevant events related to the forecast domain. Multiple strategies were assessed to understand the best exploration and exploitation trade-off required to learn from unlabeled media news entries while providing good recommendations to the users. We consider the best performance was achieved by the *Alpha trade-off* ($\alpha = 1.0$) and *Positive certain and uncertain*, which displayed a strong performance in terms of MAP score and recall. Future work will explain the models’ criteria for classifying the media news events and the associated unlabeled entry uncertainty. We expect such explanations will enhance users’ understanding of the underlying model and ease their labeling effort. Furthermore, such explanations will be extended towards the decision-making recommendations to increase the transparency behind such recommendations.

Author Contributions: Conceptualization, J.M.R.; methodology, J.M.R. and P.Z.; software, P.Z., J.M.R. and E.T; validation, P.Z., and J.M.R.; formal analysis, P.Z. and J.M.R.; investigation, J.M.R., P.Z., K.K., and I.N.; resources, K.K., B.F., and D.M.; data curation, E.T., P.Z., and J.M.R.; writing—original draft preparation, J.M.R. and P.Z.; writing—review and editing, J.M.R., K.K., I.N. and D.M.; visualization, J.M.R. and P.Z.; supervision, J.M.R., K.K., I.N., B.F. and D.M.; project administration, K.K., I.N., B.F. and D.M.; funding acquisition, K.K., B.F., and D.M. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported by the Slovenian Research Agency and the European Union’s Horizon 2020 program projects FACTLOG under grant agreement H2020-869951 and STAR under grant agreement number H2020-956573.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement:

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this manuscript:

AI	Artificial Intelligence
AL	Active Learning
ANN	Artificial Neural Networks Neural Networks
ARIMA	AutoRegressive Integrated Moving Average
ARMA	AutoRegressive Moving Average
AUC ROC ROC	Area Under the Receiver Operating Characteristic Curve
BoW	Bag-Of-Words
CPS	Cyber-Physical System
DT	Digital Twin
ERP	Enterprise Resource Planning Resource Planning
LIME	Local Interpretable Model-agnostic Explanations
MAP	Mean Average Precision
MES	Manufacturing Execution System
MLR	Multiple Linear Regression
SVM	Support Vector Machine
TF-IDF	Term Frequency - Inverse Document Frequency
USE	Universal Sentence Encoder
XAI	Explainable Artificial Intelligence

Table with results obtained across all experiments.

STRATEGY	MODEL	FEATURES	A	B	C	D
Random	LR	TF-IDF	0.8416	0.8763	0.9692	0.8835
		RoBERTa	0.8638	0.9015	0.9527	0.9114
		USE	0.8545	0.8655	0.9812	0.8996
	SVM	TF-IDF	0.8353	0.8653	0.9722	0.8703
		RoBERTa	0.8348	0.8609	0.9453	0.8935
		USE	0.8783	0.8936	0.9837	0.9023
	RF	TF-IDF	0.8491	0.8132	0.9414	0.8313
		RoBERTa	0.8637	0.8860	0.8816	0.7511
		USE	0.9185	0.8728	0.9852	0.8356
Uncertain	LR	TF-IDF	0.8391	0.8564	0.9777	0.9019
		RoBERTa	0.8209	0.8754	0.9655	0.9208
		USE	0.8595	0.8697	0.9818	0.9257
	SVM	TF-IDF	0.8264	0.8629	0.9694	0.9043
		RoBERTa	0.8035	0.8747	0.9455	0.8552
		USE	0.8803	0.8957	0.9801	0.9195
	RF	TF-IDF	0.8332	0.8199	0.9642	0.8553
		RoBERTa	0.8828	0.8857	0.9038	0.8230
		USE	0.9060	0.8753	0.9885	0.8743
Certain	LR	TF-IDF	0.8377	0.8613	0.9702	0.8716
		RoBERTa	0.8579	0.8957	0.9494	0.9105
		USE	0.8500	0.8726	0.9805	0.8925
	SVM	TF-IDF	0.8305	0.8677	0.9755	0.8719
		RoBERTa	0.8400	0.8732	0.9395	0.8919
		USE	0.8827	0.8924	0.9808	0.8960
	RF	TF-IDF	0.8553	0.8142	0.9544	0.8374
		RoBERTa	0.8266	0.8807	0.8323	0.7552
		USE	0.9129	0.8733	0.9797	0.8042
Positive uncertain	LR	TF-IDF	0.7730	0.8311	0.9442	0.8644
		RoBERTa	0.8076	0.8793	0.9285	0.9191
		USE	0.8308	0.8466	0.9685	0.9090
	SVM	TF-IDF	0.7956	0.8031	0.9536	0.8712
		RoBERTa	0.7851	0.8669	0.9087	0.8728
		USE	0.8462	0.8700	0.9626	0.9011
	RF	TF-IDF	0.7998	0.7849	0.9167	0.8327
		RoBERTa	0.8430	0.8725	0.8779	0.8133
		USE	0.8912	0.8663	0.9861	0.8403
Positive certain	LR	TF-IDF	0.7894	0.8311	0.9442	0.8661
		RoBERTa	0.8047	0.8899	0.9305	0.9192
		USE	0.8161	0.8477	0.9685	0.9086
	SVM	TF-IDF	0.7352	0.8457	0.9536	0.8713
		RoBERTa	0.7768	0.8716	0.9136	0.8670
		USE	0.8465	0.8781	0.9626	0.9011
	RF	TF-IDF	0.7917	0.7869	0.9476	0.8260
		RoBERTa	0.8455	0.8331	0.8876	0.7850
		USE	0.8989	0.8573	0.9842	0.8677
Positive certain and uncertain	LR	TF-IDF	0.8246	0.8609	0.9774	0.9022
		RoBERTa	0.8181	0.8978	0.9657	0.9235

Alpha trade-off ($\alpha = 0.5$)	SVM	USE	0.8410	0.8677	0.9826	0.9260
		TF-IDF	0.8220	0.8657	0.9688	0.8956
		RoBERTa	0.8270	0.8704	0.9624	0.8585
	RF	USE	0.8712	0.8959	0.9838	0.9179
		TF-IDF	0.8182	0.8253	0.9707	0.8586
		RoBERTa	0.8545	0.8861	0.9029	0.8351
	LR	USE	0.9180	0.8765	0.9863	0.8507
		TF-IDF	0.8363	0.8791	0.9773	0.8888
		RoBERTa	0.8341	0.8973	0.9616	0.9171
Alpha trade-off ($\alpha = 0.75$)	SVM	USE	0.8503	0.8644	0.9819	0.9115
		TF-IDF	0.8405	0.8609	0.9627	0.8710
		RoBERTa	0.8014	0.8698	0.9570	0.8350
	RF	USE	0.8873	0.8969	0.9803	0.9000
		TF-IDF	0.8556	0.8295	0.9575	0.8403
		RoBERTa	0.8662	0.8749	0.8679	0.7646
	LR	USE	0.9147	0.8747	0.9751	0.8275
		TF-IDF	0.8303	0.8653	0.9864	0.8766
		RoBERTa	0.8303	0.8969	0.9672	0.9134
Alpha trade-off ($\alpha = 1.0$)	SVM	USE	0.8480	0.8707	0.9800	0.9048
		TF-IDF	0.8353	0.8644	0.9762	0.8740
		RoBERTa	0.7992	0.8688	0.9630	0.8821
	RF	USE	0.8680	0.8978	0.9789	0.9029
		TF-IDF	0.8374	0.8274	0.9477	0.8347
		RoBERTa	0.8811	0.8792	0.8844	0.7860
	LR	USE	0.9003	0.8719	0.9830	0.8366
		TF-IDF	0.8150	0.8609	0.9773	0.9027
		RoBERTa	0.8210	0.8978	0.9658	0.9251
	SVM	USE	0.8421	0.8677	0.9813	0.9260
		TF-IDF	0.8324	0.8657	0.9533	0.8969
		RoBERTa	0.7770	0.8704	0.9632	0.8252
	RF	USE	0.8702	0.8959	0.9827	0.9179
		TF-IDF	0.8124	0.8373	0.9688	0.8480
		RoBERTa	0.8777	0.8838	0.9049	0.8226
	LR	USE	0.9248	0.8803	0.9872	0.8559

Table 1: AUC ROC scores for all experiments. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

STRATEGY	MODEL	FEATURES	A	B	C	D
Random	LR	TF-IDF	0.2741	0.4717	0.1265	0.0307
		RoBERTa	0.2675	0.4631	0.1005	0.0271
		USE	0.2787	0.4584	0.1062	0.0449
	SVM	TF-IDF	0.2699	0.4939	0.1238	0.0312
		RoBERTa	0.2799	0.4874	0.1433	0.0349
		USE	0.2694	0.4530	0.1175	0.0298
	RF	TF-IDF	0.2690	0.4753	0.1283	0.0357
		RoBERTa	0.2849	0.4741	0.1309	0.0476
		USE	0.2629	0.4854	0.1193	0.0286
Uncertain	LR	TF-IDF	0.5056	0.5227	0.5576	0.4220
		RoBERTa	0.3885	0.5212	0.4328	0.2808

Certain	SVM	USE	0.4059	0.5235	0.5186	0.2804
		TF-IDF	0.5428	0.5155	0.3572	0.1969
		RoBERTa	0.3618	0.5144	0.3313	0.2850
	RF	USE	0.3805	0.5385	0.4531	0.2747
		TF-IDF	0.4965	0.5687	0.5629	0.3595
		RoBERTa	0.5088	0.5468	0.4990	0.3350
	LR	USE	0.5132	0.5476	0.5985	0.3538
		TF-IDF	0.1791	0.4219	0.0274	0.0103
		RoBERTa	0.1623	0.4087	0.0328	0.0101
Positive uncertain	SVM	USE	0.1676	0.4158	0.0236	0.0106
		TF-IDF	0.2661	0.4084	0.0904	0.0110
		RoBERTa	0.1651	0.4308	0.0347	0.0101
	RF	USE	0.1993	0.4227	0.0289	0.0102
		TF-IDF	0.1285	0.3887	0.0258	0.0102
		RoBERTa	0.1390	0.4091	0.0394	0.0131
	LR	USE	0.1203	0.4109	0.0283	0.0114
		TF-IDF	0.5575	0.6555	0.9191	0.6422
		RoBERTa	0.5554	0.6876	0.7333	0.3689
Positive certain	SVM	USE	0.5642	0.6423	0.8603	0.3739
		TF-IDF	0.5796	0.6649	0.9378	0.9355
		RoBERTa	0.5288	0.7256	0.7001	0.4127
	RF	USE	0.6066	0.7058	0.8721	0.5820
		TF-IDF	0.5618	0.6373	0.8929	0.6623
		RoBERTa	0.6160	0.6476	0.8878	1.0000
	LR	USE	0.6156	0.6565	0.9059	1.0000
		TF-IDF	0.6145	0.7174	0.9240	0.6379
		RoBERTa	0.6448	0.7196	0.7862	0.3922
Positive certain and uncertain	SVM	USE	0.6655	0.7233	0.8978	0.4124
		TF-IDF	0.6636	0.7012	0.9383	0.9355
		RoBERTa	0.5900	0.7363	0.7493	0.4423
	RF	USE	0.7129	0.7610	0.8933	0.5675
		TF-IDF	0.6742	0.6837	0.9403	0.6226
		RoBERTa	0.6675	0.7439	0.9057	1.0000
	LR	USE	0.7528	0.7417	0.9406	1.0000
		TF-IDF	0.6449	0.6672	0.6296	0.4259
		RoBERTa	0.6140	0.6660	0.6093	0.3900
Alpha trade-off ($\alpha = 0.5$)	SVM	USE	0.6584	0.6942	0.6312	0.3249
		TF-IDF	0.6325	0.6679	0.6250	0.3826
		RoBERTa	0.6012	0.6610	0.5887	0.3730
	RF	USE	0.6912	0.7059	0.6408	0.4162
		TF-IDF	0.6710	0.6672	0.6260	0.3731
		RoBERTa	0.6807	0.6786	0.5192	0.3370
	LR	USE	0.7348	0.6867	0.6134	0.3621
		TF-IDF	0.3076	0.6420	0.0794	0.0192
		RoBERTa	0.3177	0.6467	0.0874	0.0206
	SVM	USE	0.3145	0.6538	0.0857	0.0205
		TF-IDF	0.3246	0.6337	0.0837	0.0224
		RoBERTa	0.3566	0.6388	0.0975	0.0238
	RF	USE	0.3319	0.6678	0.0872	0.0188
		TF-IDF	0.3226	0.6422	0.0931	0.0188

Alpha trade-off ($\alpha = 0.75$)	LR	RoBERTa	0.3296	0.6450	0.0973	0.0242
		USE	0.3095	0.6580	0.0970	0.0223
		TF-IDF	0.5520	0.6495	0.1073	0.0237
	SVM	RoBERTa	0.5645	0.6620	0.1036	0.0232
		USE	0.5707	0.6871	0.1143	0.0269
		TF-IDF	0.5351	0.6537	0.1237	0.0199
	RF	RoBERTa	0.5546	0.6512	0.1386	0.0291
		USE	0.5965	0.6985	0.1377	0.0248
		TF-IDF	0.5577	0.6589	0.1211	0.0401
Alpha trade-off ($\alpha = 1.0$)	LR	RoBERTa	0.5777	0.6664	0.1366	0.0247
		USE	0.6038	0.6818	0.1291	0.0282
		TF-IDF	0.6691	0.6671	0.6317	0.4211
	SVM	RoBERTa	0.6767	0.6665	0.6103	0.3914
		USE	0.6769	0.6976	0.6317	0.3257
		TF-IDF	0.6693	0.6679	0.6269	0.4259
	RF	RoBERTa	0.6280	0.6619	0.5907	0.3820
		USE	0.7341	0.7087	0.6412	0.4462
		TF-IDF	0.6856	0.6559	0.6279	0.3605
		RoBERTa	0.7076	0.6649	0.5212	0.3458
		USE	0.7765	0.6915	0.6282	0.3444

Table 2: MAP scores for all experiments. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

STRATEGY	MODEL	FEATURES	A	B	C	D
Random	LR	TF-IDF	0.5221	0.9610	0.4113	0.0659
		RoBERTa	0.5267	0.9611	0.3750	0.0755
		USE	0.5453	0.9670	0.4424	0.0909
	SVM	TF-IDF	0.5517	0.9642	0.4943	0.0734
		RoBERTa	0.5537	0.9669	0.4837	0.0905
		USE	0.5270	0.9699	0.4430	0.0722
	RF	TF-IDF	0.5341	0.9738	0.4774	0.0815
		RoBERTa	0.5540	0.9721	0.4844	0.1146
		USE	0.5293	0.9551	0.4781	0.0644
Uncertain	LR	TF-IDF	0.7179	0.9804	0.9161	0.6631
		RoBERTa	0.6448	0.9778	0.8304	0.5545
		USE	0.6692	0.9759	0.9546	0.5790
	SVM	TF-IDF	0.7383	0.9767	0.7693	0.4523
		RoBERTa	0.6217	0.9893	0.8141	0.5252
		USE	0.6283	0.9795	0.9022	0.5392
	RF	TF-IDF	0.7215	0.9841	0.9465	0.5653
		RoBERTa	0.7582	0.9778	0.8457	0.5009
		USE	0.7494	0.9845	0.9589	0.5345
Certain	LR	TF-IDF	0.4300	0.9472	0.1337	0.0114
		RoBERTa	0.4016	0.9576	0.1657	0.0091
		USE	0.4285	0.9633	0.1111	0.0110
	SVM	TF-IDF	0.4970	0.9549	0.2469	0.0139
		RoBERTa	0.3926	0.9569	0.1659	0.0091
		USE	0.4495	0.9594	0.1317	0.0101
	RF	TF-IDF	0.3565	0.9550	0.1287	0.0101

Positive uncertain	LR	RoBERTa	0.3662	0.9572	0.1796	0.0139
		USE	0.3333	0.9482	0.1444	0.0205
		TF-IDF	0.6374	0.7670	0.4826	0.2527
	SVM	RoBERTa	0.6254	0.7617	0.6415	0.1638
		USE	0.7107	0.7544	0.6419	0.1787
		TF-IDF	0.6392	0.7235	0.4333	0.2545
	RF	RoBERTa	0.6173	0.7966	0.6665	0.2538
		USE	0.7167	0.8332	0.5974	0.3245
		TF-IDF	0.6772	0.7300	0.4565	0.2807
Positive certain	LR	RoBERTa	0.7239	0.7688	0.4117	0.2545
		USE	0.7436	0.8180	0.4800	0.2545
		TF-IDF	0.6550	0.7670	0.4826	0.2565
	SVM	RoBERTa	0.6451	0.7544	0.6426	0.1638
		USE	0.7307	0.7507	0.6419	0.1805
		TF-IDF	0.6776	0.6928	0.4333	0.2545
	RF	RoBERTa	0.6002	0.7948	0.6685	0.2538
		USE	0.7253	0.8167	0.5974	0.3245
		TF-IDF	0.7207	0.7241	0.4683	0.2647
Positive certain and uncertain	LR	RoBERTa	0.6999	0.7615	0.4356	0.2545
		USE	0.8001	0.8344	0.4194	0.2561
		TF-IDF	0.8142	0.9910	0.9650	0.6678
	SVM	RoBERTa	0.7988	0.9942	0.9531	0.6057
		USE	0.8258	0.9929	0.9885	0.6036
		TF-IDF	0.7972	0.9929	0.9506	0.6395
	RF	RoBERTa	0.7718	0.9927	0.9483	0.5771
		USE	0.8373	0.9986	0.9744	0.6509
		TF-IDF	0.8415	0.9869	0.9633	0.5781
Alpha trade-off ($\alpha = 0.5$)	LR	RoBERTa	0.8381	0.9896	0.8493	0.5066
		USE	0.8628	0.9942	0.9681	0.5401
		TF-IDF	0.4041	0.9464	0.1189	0.0243
	SVM	RoBERTa	0.4108	0.9488	0.1389	0.0195
		USE	0.4133	0.9455	0.1139	0.0285
		TF-IDF	0.4336	0.9491	0.1309	0.0254
	RF	RoBERTa	0.4720	0.9630	0.1606	0.0261
		USE	0.4122	0.9520	0.1185	0.0200
		TF-IDF	0.4241	0.9536	0.1396	0.0272
Alpha trade-off ($\alpha = 0.75$)	LR	RoBERTa	0.4266	0.9577	0.2072	0.0446
		USE	0.3888	0.9498	0.1444	0.0352
		TF-IDF	0.6982	0.9722	0.1676	0.0355
	SVM	RoBERTa	0.7277	0.9838	0.1654	0.0381
		USE	0.7243	0.9852	0.1843	0.0390
		TF-IDF	0.6729	0.9721	0.2031	0.0290
	RF	RoBERTa	0.7212	0.9844	0.2548	0.0636
		USE	0.7272	0.9879	0.2372	0.0414
		TF-IDF	0.6892	0.9784	0.2241	0.0823
Alpha trade-off ($\alpha = 1.0$)	LR	RoBERTa	0.7219	0.9787	0.3215	0.0634
		USE	0.7245	0.9813	0.2122	0.0530
		TF-IDF	0.8444	0.9910	0.9661	0.6668
		RoBERTa	0.8431	0.9942	0.9554	0.6111
		USE	0.8541	0.9951	0.9896	0.6036

SVM	TF-IDF	0.8309	0.9929	0.9494	0.6894
	RoBERTa	0.8192	0.9949	0.9406	0.5789
	USE	0.8720	0.9986	0.9767	0.6820
RF	TF-IDF	0.8487	0.9893	0.9569	0.5325
	RoBERTa	0.8732	0.9891	0.8707	0.5271
	USE	0.9055	0.9907	0.9683	0.5246

Table 3: Recall scores for all experiments. A-D correspond to the four datasets we used to conduct the experiments, which are described in Table 2.

References

1. Benbarrad, T.; Salhaoui, M.; Kenitar, S.B.; Arioua, M. Intelligent machine vision model for defective product inspection based on machine learning. *Journal of Sensor and Actuator Networks* **2021**, *10*, 7.
2. Raut, R.D.; Gotmare, A.; Narkhede, B.E.; Govindarajan, U.H.; Bokade, S.U. Enabling technologies for Industry 4.0 manufacturing and supply chain: concepts, current status, and adoption challenges. *IEEE Engineering Management Review* **2020**, *48*, 83–102.
3. Lee, E.A. Cyber physical systems: Design challenges. 2008 11th IEEE international symposium on object and component-oriented real-time distributed computing (ISORC). IEEE, 2008, pp. 363–369.
4. Rajkumar, R.; Lee, I.; Sha, L.; Stankovic, J. Cyber-physical systems: the next computing revolution. Design automation conference. IEEE, 2010, pp. 731–736.
5. Rosen, R.; Von Wichert, G.; Lo, G.; Bettenhausen, K.D. About the importance of autonomy and digital twins for the future of manufacturing. *IFAC-PapersOnLine* **2015**, *48*, 567–572.
6. Grieves, M.; Vickers, J. Digital twin: Mitigating unpredictable, undesirable emergent behavior in complex systems. In *Transdisciplinary perspectives on complex systems*; Springer, 2017; pp. 85–113.
7. Grieves, M.W. Virtually intelligent product systems: digital and physical twins; 2019.
8. Grangel-González, I. A knowledge graph based integration approach for industry 4.0. PhD thesis, Universitäts-und Landesbibliothek Bonn, 2019.
9. Mogos, M.F.; Eleftheriadis, R.J.; Myklebust, O. Enablers and inhibitors of Industry 4.0: results from a survey of industrial companies in Norway. *Procedia Cirp* **2019**, *81*, 624–629.
10. Tao, F.; Qi, Q.; Liu, A.; Kusiak, A. Data-driven smart manufacturing. *Journal of Manufacturing Systems* **2018**, *48*, 157–169.
11. Preece, A.; Webberley, W.; Braines, D.; Hu, N.; La Porta, T.; Zaroukian, E.; Bakdash, J. SHERLOCK: Simple Human Experiments Regarding Locally Observed Collective Knowledge. Technical report, US Army Research Laboratory Aberdeen Proving Ground, United States, 2015.
12. Bradeško, L.; Witbrock, M.; Starc, J.; Herga, Z.; Grobelnik, M.; Mladenčić, D. Curious Cat—Mobile, Context-Aware Conversational Crowdsourcing Knowledge Acquisition. *ACM Transactions on Information Systems (TOIS)* **2017**, *35*, 1–46.
13. Settles, B. Active learning literature survey **2009**.
14. Elahi, M.; Ricci, F.; Rubens, N. A survey of active learning in collaborative filtering recommender systems. *Computer Science Review* **2016**, *20*, 29–50.
15. Konstan, J.A.; Riedl, J. Recommender systems: from algorithms to user experience. *User modeling and user-adapted interaction* **2012**, *22*, 101–123.
16. Gualtieri, M. Best practices in user experience (UX) design. *Design Compelling User Experiences to Wow your Customers* **2009**, pp. 1–17.
17. Oard, D.W.; Kim, J.; others. Implicit feedback for recommender systems. Proceedings of the AAAI workshop on recommender systems. WoUongong, 1998, Vol. 83, pp. 81–83.
18. Shivaswamy, P.; Joachims, T. Coactive learning. *Journal of Artificial Intelligence Research* **2015**, *53*, 1–40.
19. Yang, S.C.; Rank, C.; Whritner, J.A.; Nasraoui, O.; Shafto, P. Unifying recommendation and active learning for information filtering and recommender systems, 2020. doi:10.31234/osf.io/jqa83.
20. Zajec, P.; Rožanec, J.M.; Novalija, I.; Fortuna, B.; Mladenčić, D.; Kenda, K. Towards Active Learning Based Smart Assistant for Manufacturing. Advances in Production Management Systems. Artificial Intelligence for Sustainable and Resilient Production Systems; Dolgui, A.; Bernard, A.; Lemoine, D.; von Cieminski, G.; Romero, D., Eds.; Springer International Publishing: Cham, 2021; pp. 295–302.
21. Rožanec, J.M.; Kažič, B.; Škrjanc, M.; Fortuna, B.; Mladenčić, D. Automotive OEM Demand Forecasting: A Comparative Study of Forecasting Algorithms and Strategies. *Applied Sciences* **2021**, *11*, 6787.
22. Rožanec, J.M.; Mladenčić, D. Reframing demand forecasting: a two-fold approach for lumpy and intermittent demand. *arXiv preprint arXiv:2103.13812* **2021**.
23. Rožanec, J. Explainable Demand Forecasting: A Data Mining Goldmine. *Companion Proceedings of the Web Conference 2021 (WWW '21 Companion)*, April 19–23, 2021, Ljubljana, Slovenia **2021**. doi:10.1145/3442442.3453708.

24. Bradley, A.P. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition* **1997**, *30*, 1145 – 1159. doi:https://doi.org/10.1016/S0031-3203(96)00142-2.
25. Robertson, S. A new interpretation of average precision. Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval, 2008, pp. 689–690.
26. Schröder, G.; Thiele, M.; Lehner, W. Setting goals and choosing metrics for recommender system evaluations. UCERSTI2 workshop at the 5th ACM conference on recommender systems, Chicago, USA, 2011, Vol. 23, p. 53.
27. Williams, T. Stock control with sporadic and slow-moving demand. *Journal of the Operational Research Society* **1984**, *35*, 939–948.
28. Johnston, F.; Boylan, J.E. Forecasting for items with intermittent demand. *Journal of the operational research society* **1996**, *47*, 113–121.
29. Syntetos, A.A.; Boylan, J.E.; Croston, J. On the categorization of demand patterns. *Journal of the operational research society* **2005**, *56*, 495–503.
30. Wang, F.K.; Chang, K.K.; Tzeng, C.W. Using adaptive network-based fuzzy inference system to forecast automobile sales. *Expert Systems with Applications* **2011**, *38*, 10587–10593.
31. Gao, J.; Xie, Y.; Cui, X.; Yu, H.; Gu, F. Chinese automobile sales forecasting using economic indicators and typical domestic brand automobile sales data: A method based on econometric model. *Advances in Mechanical Engineering* **2018**, *10*, 1687814017749325.
32. Ubaidillah, N.Z. A STUDY OF CAR DEMAND AND ITS INTERDEPENDENCY IN SARAWAK. *International Journal of Business & Society* **2020**, *21*.
33. Dargay, J.; Gately, D. Income's effect on car and vehicle ownership, worldwide: 1960–2015. *Transportation Research Part A: Policy and Practice* **1999**, *33*, 101–138.
34. Brühl, B.; Hülsmann, M.; Borscheid, D.; Friedrich, C.M.; Reith, D. A sales forecast model for the german automobile market based on time series analysis and data mining methods. Industrial Conference on Data Mining. Springer, 2009, pp. 146–160.
35. Vahabi, A.; Hosseininia, S.S.; Alborzi, M. A Sales Forecasting Model in Automotive Industry using Adaptive Neuro-Fuzzy Inference System (Anfis) and Genetic Algorithm (GA). *management* **2016**, *1*, 2.
36. Dwivedi, A.; Niranjan, M.; Sahu, K. A business intelligence technique for forecasting the automobile sales using Adaptive Intelligent Systems (ANFIS and ANN). *International Journal of Computer Applications* **2013**, *74*.
37. Farahani, D.S.; Momeni, M.; Amiri, N.S. Car sales forecasting using artificial neural networks and analytical hierarchy process. *DATA ANALYTICS 2016* **2016**, p. 69.
38. Sharma, R.; Sinha, A.K. Sales forecast of an automobile industry. *International Journal of Computer Applications* **2012**, *53*.
39. Henkelmann, R. A Deep Learning based Approach for Automotive Spare Part Demand Forecasting.
40. Chandriah, K.K.; Naraganahalli, R.V. RNN/LSTM with modified Adam optimizer in deep learning approach for automobile spare parts demand forecasting. *Multimedia Tools and Applications* **2021**, pp. 1–15.
41. Matsumoto, M.; Komatsu, S. Demand forecasting for production planning in remanufacturing. *The International Journal of Advanced Manufacturing Technology* **2015**, *79*, 161–175.
42. Hanggara, F.D. Forecasting Car Demand in Indonesia with Moving Average Method. *Journal of Engineering Science and Technology Management (JES-TM)* **2021**, *1*, 1–6.
43. Rozanec, J.M. Explainable demand forecasting: A data mining goldmine. Companion Proceedings of the Web Conference 2021, 2021, pp. 723–724.
44. Biran, O.; McKeown, K.R. Human-Centric Justification of Machine Learning Predictions. *IJCAI*, 2017, Vol. 2017, pp. 1461–1467.
45. Arrieta, A.B.; Díaz-Rodríguez, N.; Del Ser, J.; Bennetot, A.; Tabik, S.; Barbado, A.; García, S.; Gil-López, S.; Molina, D.; Benjamins, R.; others. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion* **2020**, *58*, 82–115.
46. Ferreira, J.J.; Monteiro, M. The human-AI relationship in decision-making: AI explanation to support people on justifying their decisions. *arXiv preprint arXiv:2102.05460* **2021**.
47. Büchi, G.; Cugno, M.; Castagnoli, R. Smart factory performance and Industry 4.0. *Technological Forecasting and Social Change* **2020**, *150*, 119790. doi:https://doi.org/10.1016/j.techfore.2019.119790.
48. Micheler, S.; Goh, Y.M.; Lohse, N. Innovation landscape and challenges of smart technologies and systems – a European perspective. *Production & Manufacturing Research* **2019**, *7*, 503–528. doi:10.1080/21693277.2019.1687363.
49. Müller, V.C. Deep Opacity Undermines Data Protection and Explainable Artificial Intelligence. *Overcoming Opacity in Machine Learning* **2021**, p. 18.
50. Chan, L. Explainable AI as Epistemic Representation. *Overcoming Opacity in Machine Learning* **2021**, p. 7.
51. Samek, W.; Müller, K.R. Towards explainable artificial intelligence. In *Explainable AI: interpreting, explaining and visualizing deep learning*; Springer, 2019; pp. 5–22.
52. Henin, C.; Le Métayer, D. A multi-layered approach for tailored black-box explanations **2021**.
53. Ribeiro, M.T.; Singh, S.; Guestrin, C. "Why should i trust you?" Explaining the predictions of any classifier. Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, 2016, pp. 1135–1144.
54. Lundberg, S.; Lee, S.I. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874* **2017**.
55. Ribeiro, M.T.; Singh, S.; Guestrin, C. Anchors: High-Precision Model-Agnostic Explanations. *AAAI*, 2018, Vol. 18, pp. 1527–1535.
56. Rüping, S.; others. Learning interpretable models **2006**.
57. Artelt, A.; Hammer, B. On the computation of counterfactual explanations—A survey. *arXiv preprint arXiv:1911.07749* **2019**.

58. Mothilal, R.K.; Sharma, A.; Tan, C. Explaining machine learning classifiers through diverse counterfactual explanations. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020, pp. 607–617.
59. Verma, S.; Dickerson, J.; Hines, K. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596* **2020**.
60. Singh, R.; Dourish, P.; Howe, P.; Miller, T.; Sonenberg, L.; Velloso, E.; Vetere, F. Directive explanations for actionable explainability in machine learning applications. *arXiv preprint arXiv:2102.02671* **2021**.
61. Hrnjica, B.; Softic, S. Explainable AI in Manufacturing: A Predictive Maintenance Case Study. *Advances in Production Management Systems. Towards Smart and Digital Manufacturing*; Lalic, B.; Majstorovic, V.; Marjanovic, U.; von Cieminski, G.; Romero, D., Eds.; Springer International Publishing: Cham, 2020; pp. 66–73.
62. Rehse, J.R.; Mehdiyev, N.; Fettke, P. Towards explainable process predictions for industry 4.0 in the dfki-smart-lego-factory. *KI-Künstliche Intelligenz* **2019**, *33*, 181–187.
63. Goldman, C.V.; Baltaxe, M.; Chakraborty, D.; Arinez, J. Explaining Learning Models in Manufacturing Processes. *Procedia Computer Science* **2021**, *180*, 259–268.
64. van der Waa, J.; Nieuwburg, E.; Cremers, A.; Neerincx, M. Evaluating XAI: A comparison of rule-based and example-based explanations. *Artificial Intelligence* **2021**, *291*, 103404.
65. Ghai, B.; Liao, Q.V.; Zhang, Y.; Bellamy, R.; Mueller, K. Explainable active learning (xal) toward ai explanations as interfaces for machine teachers. *Proceedings of the ACM on Human-Computer Interaction* **2021**, *4*, 1–28.
66. Tulli, S.; Wallkötter, S.; Paiva, A.; Melo, F.S.; Chetouani, M. Learning from Explanations and Demonstrations: A Pilot Study. 2nd Workshop on Interactive Natural Language Technology for Explainable Artificial Intelligence, 2020, pp. 61–66.
67. Settles, B. From theories to queries: Active learning in practice. *Active Learning and Experimental Design workshop In conjunction with AISTATS 2010. JMLR Workshop and Conference Proceedings*, 2011, pp. 1–18.
68. Zhu, J.J.; Bento, J. Generative adversarial active learning. *arXiv preprint arXiv:1702.07956* **2017**.
69. Li, Q.; Peng, H.; Li, J.; Xia, C.; Yang, R.; Sun, L.; Yu, P.S.; He, L. A Survey on Text Classification: From Shallow to Deep Learning. *ArXiv* **2020**, *abs/2008.00364*.
70. Sparck Jones, K., A Statistical Interpretation of Term Specificity and Its Application in Retrieval. In *Document Retrieval Systems*; Taylor Graham Publishing: GBR, 1988; p. 132–142.
71. Minaee, S.; Kalchbrenner, N.; Cambria, E.; Nikzad, N.; Chenaghlu, M.; Gao, J. Deep Learning-Based Text Classification: A Comprehensive Review. *ACM Comput. Surv.* **2021**, *54*. doi:10.1145/3439726.
72. Schröder, C.; Niekler, A. A Survey of Active Learning for Text Classification using Deep Neural Networks, 2020, [[arXiv:cs.CL/2008.07267](https://arxiv.org/abs/2008.07267)].
73. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient Estimation of Word Representations in Vector Space, 2013, [[arXiv:cs.CL/1301.3781](https://arxiv.org/abs/1301.3781)].
74. Le, Q.V.; Mikolov, T. Distributed Representations of Sentences and Documents, 2014, [[arXiv:cs.CL/1405.4053](https://arxiv.org/abs/1405.4053)].
75. Cer, D.; Yang, Y.; yi Kong, S.; Hua, N.; Limtiaco, N.; John, R.S.; Constant, N.; Guajardo-Cespedes, M.; Yuan, S.; Tar, C.; Sung, Y.H.; Strope, B.; Kurzweil, R. Universal Sentence Encoder, 2018, [[arXiv:cs.CL/1803.11175](https://arxiv.org/abs/1803.11175)].
76. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, 2019, [[arXiv:cs.CL/1810.04805](https://arxiv.org/abs/1810.04805)].
77. Reimers, N.; Gurevych, I. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, 2019, [[arXiv:cs.CL/1908.10084](https://arxiv.org/abs/1908.10084)].
78. Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; Stoyanov, V. RoBERTa: A Robustly Optimized BERT Pretraining Approach, 2019, [[arXiv:cs.CL/1907.11692](https://arxiv.org/abs/1907.11692)].
79. Lu, J.; MacNamee, B. Investigating the Effectiveness of Representations Based on Pretrained Transformer-based Language Models in Active Learning for Labelling Text Datasets, 2020, [[arXiv:cs.IR/2004.13138](https://arxiv.org/abs/2004.13138)].
80. Liere, R.; Tadepalli, P. Active Learning with Committees for Text Categorization. *AAAI/IAAI*, 1997.
81. Schröder, C.; Niekler, A.; Potthast, M. Uncertainty-based Query Strategies for Active Learning with Transformers, 2021, [[arXiv:cs.CL/2107.05687](https://arxiv.org/abs/2107.05687)].
82. Yang, S.C.H.; Rank, C.; Whritner, J.; Nasraoui, O.; Shafto, P. Unifying recommendation and active learning for information filtering and recommender systems **2020**.
83. Bloodgood, M. Support Vector Machine Active Learning Algorithms with Query-by-Committee Versus Closest-to-Hyperplane Selection. *2018 IEEE 12th International Conference on Semantic Computing (ICSC)* **2018**. doi:10.1109/icsc.2018.00029.
84. Leban, G.; Fortuna, B.; Brank, J.; Grobelnik, M. Event registry: learning about world events from news. *Proceedings of the 23rd International Conference on World Wide Web*, 2014, pp. 107–110.
85. Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; Davison, J.; Shleifer, S.; von Platen, P.; Ma, C.; Jernite, Y.; Plu, J.; Xu, C.; Scao, T.L.; Gugger, S.; Drame, M.; Lhoest, Q.; Rush, A.M. Transformers: State-of-the-Art Natural Language Processing. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*; Association for Computational Linguistics: Online, 2020; pp. 38–45.
86. Wilcoxon, F. Individual comparisons by ranking methods. In *Breakthroughs in statistics*; Springer, 1992; pp. 196–202.