

A Bird's Eye View on Requirements Engineering and Machine learning*

*Note: Sub-titles are not captured in Xplore and should not be used

1st Tahira Iqbal
fortiss GmbH
Munich, Germany
iqbal@fortiss.org

2nd Parisa Elahidoost
fortiss GmbH
Munich, Germany
elahidoost@fortiss.org

3rd Levi Lucio
fortiss GmbH
Munich, Germany
lucio@fortiss.org

Abstract—Machine learning (ML) has demonstrated practical impact in a variety of application domains. Software engineering is a fertile domain where ML is helping in automating different tasks. In this paper, our focus is the intersection of software requirement engineering (RE) and ML. To obtain an overview of how ML is helping RE and the research trends in this area, we have surveyed a large number of research articles. We found that the impact of ML can be observed in requirement elicitation, analysis and specification, validation and management. Furthermore, in these categories, we discuss the specific problem solved by ML, the features and ML algorithms used, and datasets, when available. We outline lessons learned and envision possible future directions for the domain.

Index Terms—component, formatting, style, styling, insert

I. INTRODUCTION

This is an introduction [10]

RQ1: How is ML used inside RE tools?

RQ2: Is using ML inside RE tools beneficial?

RQ3: What ML task(s) is (are) used for which purpose in RE?

RQ4: Which model types are used to perform the RE tasks?

RQ5: How are ML features extracted/selected to guide RE tasks?

A. Machine Learning

Tahira ▶ Describe terms for feature extraction that includes stopword removing, sentimental analysis, stemming and Bag of words, VSM etc ◀

Machine-Learning (ML) [?] is a range of algorithm to approximate functions and discover patterns in data. Historically, models and heuristics are human-built exhaustive prescriptions of how a system should behave. ML is grounded on different premises: rather than relying on humans to input all the possible cases the system can handle, the field attempts to extrapolate patterns from a representative set of examples that illustrate the expected behaviors. The way in which a learning algorithm operates attempts to emulate the way in which humans learn: from a set of examples, a general model for a behavior is induced.

Many learning algorithms exist, based on different visions of how learning happens in practice [10]. All these algorithms

have in common the notion of *features*. Features correspond to characteristics of what is being learned and provide the grounds for the algorithm to abstract from the complexities of the real world. Assume for example that an algorithm should learn, based on a brain scan of a medical patient, to decide whether that patient has brain cancer or not. A number of *features* such as for example the “number of irregular objects in the scan”, the “color of such objects”, the “disposition of such objects” would be provided to the algorithm. Additionally, the algorithm is fed a number of brain scans together with decisions previously taken on them (cancer found / cancer not found) – the *training data*. The learning algorithm then undergoes a *training phase*. It attempts to find an internal model that allows it to map the decisions to the brain scans, given the training data. The model obtained from the training step is useful if it performs well (generalizes) when applied to new data from outside the training set – in our example, when it can accurately diagnose brain cancer for new brain scans. Such generalization is based on the premise that inputs that are “closer”, in terms of the given *features*, should lead to “closer” outputs.

The formal notion of “closeness” is a characteristic of the learner algorithm being employed and determines how the algorithm generalizes the computation from the given examples. Achieving good generalizations is the cornerstone of machine learning and *overfitting* (performing very well on training inputs but very poorly on new inputs) is one of its major challenges. **Levi** ▶ cut these last couple of sentences ◀

More formally, in textbooks, courses and articles, Machine Learning is often defined following the definition of Tom Mitchell [?]:

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E.

Therefore, it is said that to classify some patients into classes (e.g healthy and unhealthy), the task T, one have to define an algorithm that provides a model, such as an artificial neural network. The quality of this model is quantified by a measure

P, for instance its accuracy while predicting the classes. This measure is then sent back to the algorithm, a new experience E, in order to choose or improve the model. A machine learning tasks can be discussed and subdivided based on the elements of the following equation:

$$f(\mathbf{X}) = \mathbf{y} + \xi$$

where \mathbf{X} is the $n \times d$ input matrix, containing n samples characterized by d features, \mathbf{y} is the $n \times 1$ target vector containing the classes of the n samples and ξ is a $n \times 1$ vector representing the noise. The goal is to approximate f in order to provide the best mapping between \mathbf{X} and \mathbf{y} , given some noise ξ . Indeed, the approximation has to map a \mathbf{X} containing noise, to a \mathbf{y} which may contain noise too. For instance, uncontrolled conditions such as the room temperature and the exposure time to this temperature can induce variations in the information contained in collected blood samples. Moreover, the ξ term also contains the approximation error when, for example, one tries to approximate a non-linear function with a linear function.

The problem presented by the above equation is called *supervised learning*, and can be roughly subdivided in two popular problems: *classification* and *regression*. When the target vector \mathbf{y} is composed of categorical values (i.e. classes), then we have a classification problem. The goal is to learn how to link instances or samples in \mathbf{X} to a certain class (e.g. healthy patient or unhealthy patient). However, if the target vector contains continuous values, we face a regression problem (e.g. predict the body temperature of a patient given some clinical features of the patient).

In some cases, \mathbf{y} is not given and we have to find patterns in \mathbf{X} “blindly.” This is called an *unsupervised* learning problem. Finding clusters in \mathbf{X} , i.e. finding a \mathbf{y} that has never been given, is such a problem. For instance, one may want to group patients based on symptoms they have.

Reinforcement Learning can be seen as an intermediate problem where \mathbf{y} is not given but the procedure is guided nevertheless. In RL, an agent has to find a sequence of actions leading to a success. The fact that the sequence leads to a success, i.e. what would be in \mathbf{y} , is not known in advance, but rewards are given to the agent in order for him to know if it follows a path to success. In other words, the goal is, by providing rewards along the way, to find the sequence of actions leading to the desired state. Typical examples can be found in gaming, where an agent receives a reward when he wins the game. From the different chains of actions that led him to a reward, the agent must generalize to find how to win that game.

B. Requirements Engineering

Software systems are developed over millions of lines of code, number of modules and documents. The primary goal of the software system is to satisfy users by developing the software that can meet their needs and expectations. This goal is achievable by applying different methodologies and engineering techniques. One of the key factor is to understand

and identify the needs of users, also known as, software requirements. Software requirement engineering is the process that helps to determine the requirements in a systematic way to know what functionalities the targeted system should have to fulfil user’s needs. Formally RE is defined as [47]:

“Requirements engineering is the branch of software engineering concerned with the real-world goals for, functions of, and constraints on software systems. It is also concerned with the relationship of these factors to precise specifications of software behavior, and to their evolution over time and across software families.”

Software requirements plays a key role in the success of a project. In the USA, a survey was conducted over 8380 projects by 350 companies to know the project failure rates. The report [40] results showed only 16.2% projects were completed successfully and one-half (52.7%) of projects met with challenges and were completed with partial functionalities, time delays and over budget. Almost 31% of the projects were never completed. The main cause told by the executive managers was poor requirements engineering. The main culprits were the lack of user involvement (13%), requirements incompleteness (12%), changing requirements (11%), unrealistic expectations (6%), and unclear objectives (5%).

Software requirement engineering has mainly four phases; *requirement elicitation*, *requirement analysis*, *requirement documentation* and *requirement verification* [26]. Requirement elicitation [6], [48] helps to understand the stakeholders needs, e.g. what features he wants in the software. Requirement elicitation techniques are mostly derived by the social sciences, organizational theory, knowledge engineering and practical experience. For requirements elicitation, different techniques exist in the literature that include interviews, questionnaires and ethnography etc to understand the stakeholder requirements. Requirement analysis [32] is the next step after requirement elicitation. In this phase, software requirements are analyzed to check for conflicts and consistency of requirements. It is also makes sure that the requirements are clear and complete. Furthermore, the agreed requirements are documented. This documentation has a clear and precise definition of the system functionalities. It also acts as an agreement between stakeholders and developers. These requirements are documented usually as diagrams, mathematically formulae or natural languages.

These documents are used and iterated upon until the end of the projects. System requirements are classified into business requirements, user requirements, functional requirements (FR) and non-functional requirements (NFR). Functional requirements are the system requirements that include the main features and characteristics of the desired system. Non-functional requirements are the system properties and constraints [7], [16]. NFRs set the criteria for judging the operation of the system e.g. performance, availability or reliability. Business

requirements are specified to address business objectives, vision, and goals. It is defined at a high level to keep the knowledge from organization or company side for designing the products. User requirements are the wish list for the system from users **Levi** ▶I don't think the intersection between FRs, NFRs, BRs and URs is empty◀ . User requirements are valuable for ensuring the system performs similarly as users wanted it to do.

II. CONTRIBUTIONS

Levi ▶summarize preprocessing steps and feature calculation◀ **Levi** ▶disclaimer: sometimes features are presented sometimes not◀ **Levi** ▶?◀ Stakeholders needs (features) **Levi** ▶features in this case are just words◀ are initially pre-processed by removing common (stop) terms, and computation of tf-idf **Levi** ▶explain tf-idf, maybe extract it to a common section, in general explain that features are words and the instances are tweets/reviews/...◀ for each term.

A. Requirements Elicitation and Discovery

The manual process of requirement elicitation is time-consuming and requires efforts and resources. A project's success majorly depends upon the precise identification of stakeholder expectations and requirements for their desired system. Poor elicitation can lead to over budget and time issues. A possibility to do requirements elicitation can be mining available datasets e.g. social media, documents. Also, gathering and analyzing information available on the internet, publically or privately available dataset. In short, gathering data from different sources and analyzing it for requirements elicitation. The latest trend for identifying user requirement is to mine social media coming from the twitter, google store, and app stores. These user reviews are not structured requirements and contain other information too such as praise, dislikes or bug reports. The unstructured, extra information and noise presence in the above-described sources makes manual requirements elicitation difficult and challenging. Automated requirement elicitation is helpful in these cases and can significantly reduce time, effort, and cost. This is ML classification and clustering task give the set of information and identifying it as a requirement or not.

Guzman *et al.* [18] have proposed the ALERTme approach for classifying, grouping and ranking tweets during the software evolution process. For identifying evolutionary requirement supervised classification has been performed on tweets to classify them as improvement requests or not, using the Naïve Bayes algorithm. It was the first this study of its kind performed on software tweets. For tweets preprocessing tokenization, conversion of all the texts into lower case, extracting n-gram with a one to three words length, removing stopwords and stemming were applied. After this classifier was trained with following steps: 1) conversion of preprocessed tweets into a VSM model, 2) train a classifier on a set of manually annotated tweets, 3) predict the tweets categories using a trained classifier. A data set of 68108 tweets i.e.

collection of two-month software tweets of Spotify, Dropbox, and Slack was used.

Williams *et al.* [45] performed a similar study on tweets for classifying them as user requirements or not. The study improved results in comparison to the existing studies using Naïve Bayes and VSM algorithm. Different techniques were implemented for feature extraction that includes stopword removing, sentimental analysis, stemming **Levi** ▶bring to machine learning section◀ , and Bag of words. For the learning process, manually annotated tweets were used. However, results showed that software tweets are neutral in nature and sentiment analysis did not influence the ML algorithm. It used 4000 randomly selected tweets from ten different software.

Jiang *et al.* [21] presented a study on mining user reviews in mobile application (app) stores for discovering the evolutionary requirements for developers. This method clustered opinion expressions in the form of a macro network topology, and combine polarized sentimental analysis with the economic factors to decide evolutionary requirements. It used two datasets of online reviews one was from software package Karplesky internet security 2011 from Amazon with 380 reviews, and mobile app of TuneIn Radio Pro V3.6 from the app store with 461 reviews.

Lange *et al.* [28] mapped software requirement elicitation process to an existing military tool skiweb that used for making the decision. The proposed methodology used the supervised Naive Bayes algorithm to classify the text document to find related requirements. Furthermore, the recommender system used topic modeling to identify the key stakeholder for which that requirement is important and allow for the analysis. This study used internal organizational dataset Skiweb Data i.e. wiki, blogs etc.

Jha *et al.* [20] discovered user requirements by mining feature requests from the app store. The data were classified into three categories bugs, features, and junk. For classification, SVM and Naïve Bayes algorithm were used. The study used frame semantic analysis that generalize more general contexts. This technique produced a slower dimensional model with a smaller number of features which help to enhance the prediction capabilities. It combined the existing dataset from the past studies and iOS apps reviews including CreditKarma, FitBit, and Gmail.

Maalej [31] presented a classification study of app reviews into bug reports, feature requests, user experiences, and ratings. The binary classifier Naive Bayes performed best than other algorithms also it showed binary classifier are suitable for predicting the than multi-class classifiers. Different review feature such as ratings, tense, sentiment scores, and length enhanced the classification performance. A dataset of 4400 manually annotated reviews from Google Play and the Apple app store was used.

Herrera *et al.* [4] build a semi-automated recommender system to manage broad stakeholder participation in the requirements elicitation and prioritization process. On these forums, stakeholders can work collaboratively to transform statements of need into sets of articulated and prioritized requirements. The proposed system analyzed online dataset i.e. gathered from stakeholders and automatically generated highly specialized topics for generating different forums. User profiles are created as a next step that depicts stakeholders interest. These profiles also help to create recommendations according to the interests of a community of similar stakeholders. For topics identification, an unsupervised clustering technique is applied to unstructured or semi-structured data. For the evaluation dataset was a collection of feature requests. 36 graduate-level students created these requests for an Amazon-like student web-portal system.

1) Text mining: Levi ►add to preprocessing◀

Hollis *et al.* [19] proposed an initial study to automate requirement elicitation in Agile environment by providing the list of words and loosely formatted list of requirements. The proposed methodology applied text mining technique on recorded conversation of the stakeholder and developer conversation. It was a short paper and did not provide details.

Dong *et al.* [11] also applied text mining on different form of document and resources from internet for gathering requirement. The system applied data preprocessing as word segmentation and stop words removal and build up the VSM model. Kaiya *et al.* [23] proposed a tool to improve the domain knowledge ontology for requirement elicitation by using web mining and NLP technique. It helped to mine the general concepts to ontology for requirements elicitation.

B. Requirements Specification and Analysis

1) *Identifying Non-Functional Requirements:* Non-functional requirements may not be explicitly mentioned in a formal specification requirements documents even though, all systems have them [39]. Moreover, freeform documents like interview notes, meeting minutes and scattered requirements specifications include non-functional requirements which need to be detected and classified. In order to support analyst in the error-prone task of manually discovering and classifying NFRs machine learning can be useful. Automatic detection can be used to quickly and more effectively analyze large and complex documents for searching the NFRs [5]. This is a classification problem as from a set of requirements we want to decide a class membership.

One of the studies is by Slankas *et al.* [39] where they automatically identified and classified sentences in natural language from use agreements, install manuals, regulations, request for proposals, requirements specifications, and user manuals output into 14 different NFRs categories: Access Control, Audit, Availability, Legal, Look and Feel, Maintenance, Operational, Privacy, Recoverability, Performance

and Scalability, Reliability, Security, Usability. Their two-step process: 1) parse natural language and turn sentences into graphs 2) classify sentences into categories with k-nearest neighbor algorithm led them into finding 20 keywords for each category of NFRs for their classifier. They trained the NFR classifier with a wide variety of open and closed source EHRs (Electronic Health Record), various industry standards (HL7, CCHIT), governmental regulations, and other document sources exist to elicit documentation.

Cleland-Huang *et al.* [5] provided the same approach and used k-nearest neighbor classification for grouping non-functional requirements: availability, look-and-feel, legal, maintainability, operational, performance, scalability, security, and usability. For training their classifier they used 15 requirements specifications developed as term projects by master students at DePaul University.

2) *Identifying Functional Requirements:* Software requirements specifications are usually stated in informal, imprecise and ambiguous natural language, thus analyzing them is a challenging task. However, for requirements reuse analyzing them is a vital task. Automatically extract structured information of functional requirements from Software Requirements Specifications and grouping them into different categories is a machine learning classification task [43].

Wang *et al.* [43] applied a combination of machine learning, natural language processing, and semantic analysis methods for automatically extract functional requirements and classify them into 10 different cases: Agentive(The main participant of activities), Action(The main functional operation), Objective(The objects affected by the Action), Agent mode(The modifier or property of the Agentive), Objmod(The modifier or property of the Agentive), Locational(The location or destination of an Action), Temporal(The occurrence time or frequency of an Action), Manner(The way or tool by which an Action is performed), Goal(The purpose to be achieved through the Action), Constraint(The conditions or constraints of the Action). Their framework employed techniques of semantic role labeling and machine learning and has four steps: corpus construction, NLP preprocessing, feature extraction and EFRF (Extended Functional Requirements Frame) functional cases extraction. which for NLP processing they did tokenization, lemmatization, part-of-speech tagging (POS tagging) and dependency parsing. They trained their bi-directional LSTM-CRF network which is a variant of Recurrent Neural Networks architecture model with E-commerce requirements dataset and test it on requirements of automaker systems. They proposed EFRF through analyzing the linguistic characterization of SRSs. EFRF consists of 10 mentioned functional cases for capturing the semantic information in the natural language functional requirements. For example, the “Agentive” case can be mapped to “User” or “System” and the “Constraint” is usually corresponding to the “Precondition”. Ultimately, they showed that their trained model on E-commerce requirements dataset can be used to extract semantic information from the

requirements of automaker systems.

3) *Distinguishing Functional from Non functional Requirements*: The success of a system does not solely depend on its functional requirements. Like functional requirements, it also significantly depends upon the adherence to non-functional requirements. However, the primary focus is generally more towards identification and specification of FRs. NFRs are usually identified and specified in later development stages that can increase the risks in development life cycle. FRs tend to be more straightforward e.g. store and retrieve transaction. On the other side, NFRs are complicated and challenging to implement e.g., making the design to meet NFRs or design test case for them. Different types of requirements analyzed in a different way, and it is useful to have a separate division to look at one particular division. That is why it is necessary to distinguish between FR and NFR and categorize NFRs into subcategories. This distinction helps to manage changes in requirements. The manual division is difficult and time consuming. Machine learning can be used for reducing the effort and categorizing the requirements based on the text segment analysis. This is ML classification task: give the set of requirements and identifying its category.

Mengmeng Lu *et al.* [29] automatically classify the user review text mobile from application (app) stores into FR, NFR, and others. It further classified NFRs into four categories including reliability, usability, portability, and performance. It used supervised machine learning algorithm (bagging) for training the classifier. The text trimming used stopwords elimination, lemmatization, stemming, and sentences split. Word2vec was used for augmenting the user review, which is a two layer neural network model to process text for finding the word embedding.

Deoxadez *et al.* [8] performed semi-supervised classification techniques for automated classification of FRs and NFRs on user reviews from the app store. This study dealt with two problems: 1) minimize annotation effort or label the big dataset of user reviews, and 2) classification of FRs and NFRs. The proposed solution to solve the first problem used semi-supervised self-labeling algorithm. Self-labeling algorithms required a small amount of dataset to get comparable results as supervised techniques. Naïve Bayes algorithm was selected because of high performance results for classification problem. Features were obtained by applying standard text mining technique and additional attribute embellishment. **Levi** ▶what is attribute embellishment?◀. For text mining technique Inverse Document Frequency (IDF) transform, Term Frequency (TF) transform, lowercase transformation, minimum term frequency, stemmer, and number of words were used. The second stage involved removing numbers, 2-letter words and other symbolic characters.

Kutranvoic *et al.* [27] performed automated analysis on software requirements and performed classification on FRs,

NFRs and subcategories of NFRs using supervised machine learning algorithm (support vector machine). Some NFR had negligible requirement in dataset due to this they were ignored. The additional dataset of user comments from Amazon was integrated into the main dataset to avoid the data imbalance problem in NFR. The predictor used text-preprocessing techniques such as removal of punctuations, removal of stop words, and lemmatization for feature extraction. All experiments were performed on dataset given by RE datatrack.


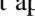
Abad *et al.* [1] targeted two problems: first is classification of FR and NFR and second classification of NFR into categories. Preprocessing of text performed by applying: POS tagging, Entity tagging and temporal tagging. As a next step feature co-occurrence and regular expression used to increase the weight of influential words in the dataset **Levi** ▶don't get it, too much detail◀. J.48 DT used for the classifying the FR and NFR. For achieving second goal topic Modeling unsupervised algorithm LDA and BTM applied. For topic generation different algorithms such as clustering, k-means, hybrid clustering and k-means, BNB were applied. The results showed BNB worked better out of clustering, k-means, LDA, BTM. learnt decision tree to classify the FRs and NFRs and improved results from 89% to 95% by analyzing the requirements in natural language. For NFRs subclassification, Binarized Naïve Bayes (BNB) achieved highest results. The preprocessing of text involved, Part of Speech (POS) tagging, entity tagging, and temporal tagging.

Garzoli [14] classified data into five types: FRs, NFRs, design and construction constraints, operator requirements, performance requirements. The dataset was taken from combat management system of a Naval Combat System in general complex system domain. However, the primary goal was to come up with a general architecture for large-scale and adaptive requirement analysis. For classification multi VSM model used and for lexical and grammatical features BoW + N-words + N-POS achieved highest accuracy.

Wieloch *et al.* [44] introduced methods that enhance the accuracy of tracing (classifying) requirements which occurred frequently across multiple projects and/or domains. In the paper, they presented Trace by Classification (TBC), a machine learning approach in which a classifier is trained to identify and classify requirements and/or other kinds of software artifacts which occur relatively frequently across different projects (they call this process generating trace links for software artifacts in their research). The first step in their approach was preprocessing to eliminate common stop words. Then training phase identified set of indicator terms for each NFR category and the classifier trained by the set of identified weighted indicator terms that can be used for the last step which is to classify additional artifacts into functional and non-functional (e.g. look-and-feel, performance, security, etc). A probability value represented the possibility that the new requirement belonged to a certain artifact type computed as a

function of the occurrence of indicator terms of that type in the requirement.

4) *Requirement Prioritization*: Complex software system generally has thousands of requirements with multiple stakeholders and customers. Each one of them has their own set of requirements and opinions and wants their requirements implementation accordingly. However, several factors make implementation process of all the requirements infeasible and hard. For example project resources, time, budget, different opinion among stakeholders, etc. all factors affect the implementation process. Therefore, it is important to make a proper decision for prioritizing requirements considering all the factors for the success of the project. Different models exist in the literature for prioritization of software requirements such as analytical hierarchical process (AHP) [38], Goal oriented [42], cost value approach [24] etc. In these techniques, human input is majorly involved. Qaddoura *et al.* [?] reviewed the prioritization techniques and also shed light on the ML contribution. ML can be used for automated analysis of these large set of software requirements prioritization, and it can also help to improve the existing techniques.

Dhingra *et al.* [9] predicted the most appropriate technique for software requirement prioritization process. The input from the user is taken as the characteristic value **Levi**  **what is the characteristic value?**  and output the most appropriate requirement prioritization method such as AGORA, AHP, etc. Fuzzy interface technique was used for prediction of prioritization method based on certain rules. The proposed system framework has three phases: training phase, fuzzing inference process, testing phase. The drawback of fuzzy approach was the wrong prediction for boundary values. It was overcome by adopting decision tree. It learned a model using the training dataset. The training dataset was gathered from the literature to obtain attributes. These attributes list included consistency, traceability, priority basis, rigorous/systematic, distributed stakeholder, cognitive aspects, and human experience. DT learned from the dataset and predicted the prioritization technique. Out of the 45 test samples, the framework was able to classify 43 tests accurately.

Avesani *et al.* [2] presented the study that dealt with the scalability problems that arises in managing the prioritization of a large number of requirements specified in AHP technique. The existing solution to scalability issues used heuristics. It helped to decide when to stop the pairwise elicitation process. The proposed framework outperformed AHP by giving an accurate approximation of the final ranking within a limited elicitation effort. It used rank-based learning algorithm and produced a ranking of all requirements. This technique builds up the solution by looking at examples. The input for the learning algorithm are a finite set of requirements, the ranking criteria, initial user preferences and density function.

The similar study was performed in [3] from the same group [2] for identification of decision-making issues related to the management of risks in Open Source Software adoption in medium and large organizations. A semi-automated system was proposed that used case based ranking classification algorithm. The input was priority elicitation of goals by the decision maker and risk goal ranking function (predefined ranking criteria ordering the goal). As an output, it ranked the final risk-based goals.

5) *Security Requirements*: Due to the orthogonal character of their impact on a system, *security* requirements are notoriously difficult to identify, objectify and quantify []. Also during requirement specification, it very often happens that security requirements are masked by functional requirements (but can be deduced from the context of the domain the system operates in) [37]. Because of this, it often happens in practice that security requirements are only marginally tackled during system construction [], paving the way to potentially catastrophic consequences. Machine learning can be of use here by aiding in the identification of segments of text that describe security requirements. This is a *classification* problem: given a text, identify which parts of it correspond to which type of security issues.

Jindalet *et al.* [22] automatically learn decision trees that can be used to classify security requirements as *authentication*, *access control*, *encryption* or *data integrity*. The *features* used are relevant terms found in the text. Such relevant terms are obtained by the following sequence of actions: 1) removing stop words from the text; 2) stemming the remaining words; and 3) ranking the stemmed words by their *info-gain* measure.

Riaz and her colleagues [37] use the k-nearest neighbors algorithm to classify sentences in requirements documents as *confidentiality*, *integrity*, *authentication*, *availability*, *accountability* or *privacy* requirements. In order to find adequate sentences and provide context to the classifier, the authors start by finding a type for each sentence among the possibilities *title*, *list start*, *list element* or *normal sentence*. For the classification the authors use a modified version of the Levenshtein distance [] based on the number of word transformations needed to go from one term in one sentence to a term in another sentence. The classifier is trained using requirements sentences from the healthcare domain that are manually classified. A particularity of the approach is that each security requirement type is associated to a template that helps in translating the security requirements into functional requirements in order to ease during the implementation of the final system.

[25]

6) *Contextual Requirements*:

C. Requirements Validation

1) *Traceability*: Validation is to guarantee that requirements are reflecting stakeholders' needs, confirm the quality of the

system, consistency, and traceability. One of the definitions for requirements traceability is given by [17]:

“Requirements traceability refers to the ability to describe and follow the life of a requirement, in both a forward and backward direction (i.e., from its origins, through its development and specification, to its subsequent deployment and use, and through periods of on-going refinement and iteration in any of these phases).”

Based on this definition the emphasis is on the ability to track the life of requirements and their established links within other artifacts. However, the main barrier assures traceability is the needed effort for building and maintaining the links between those artifacts. That is why many research has tried to apply machine learning and automated tools for facilitating the establishment of links [15]. Traceability tackled in the research mainly by the use of machine learning classification and reinforcement learning methods.

Gervasi *et al.* [15] investigate what can be learned from links that are already established. They build classifiers as a mean to develop models of tracing that can then be interpreted by humans to understand how requirement tracing is done in practice. Their purpose is to revise the existing models of hard-coded traceability tools such as VSM. They used a publicly-available dataset of requirements with traceability information, originally based on the CM-1 project by the NASA Metrics Data Program. Their approach has 5 steps: 1) tokenize and stem requirements and removing stopwords 2) derive two features from each term *t* in the vocabulary, one for the occurrence of *t* in a high-level requirement and one for the similar occurrence in a low-level requirement 3) transfer requirements into a vector of features 4) From these vectors derive set of classification cases by joining one high-level requirement and one low-level requirement and adding a classification of *link* or *nolink* based on whether that particular pair was a true link in the original dataset, or not 5) finally, use the dataset to train and test two different classifiers from the WEKA collection, a Naive Bayesian classifier, and the J48 decision-tree classifier.

Sultanov *et al.* [41] finds traceability candidates from high-level to low-level requirements by the use of reinforcement learning. They used textual high and low-level requirements documents as an input and try to find the candidate traces. Their technique demonstrated statistically significantly better results than the Information Retrieval technique.

D. Requirements Management

1) *Visualization*: Natural language requirement documents can be hard to comprehend and analyze. Similarly, stakeholders have to review and understand requirements for large and complex systems. In these scenarios, basic information visualizations, like charts and graphs have been used in requirements engineering. These visualizations are usually applied to improve textual requirements with

summarization that combined large amounts of information into a single representation for quick consumption by stakeholders [36]. Machine learning is useful in discovering visualized groups of large numbers of requirements. In research both clustering and classification methods were used for this purpose.

ReCVisu (Requirements Clustering Visualization) tool is presented in Reddivari *et al.* [36] paper. ReCVisu, an exploration tool based on quantitative visualizations helps requirements engineers understand the nature of the requirements in a visual form. In ReCVisu, the dependence graph consists of requirements artifacts as nodes and the textual similarities as edges. The automatic grouping of requirements into clusters can help in areas such as uncovering the requirements structure, navigating around the requirements space, modularizing crosscutting concerns, and understanding requirements interactions and evolution.

Pinqui *et al.* [34] recognize the enormous volume of requirements as big data with which companies struggle to make strategic decisions early on. Therefore, they built a complete visual framework to filter requirements from stakeholders in a way that architects can make better insightful decisions. They suggest training a multi-class SVM model from domain-specific (mechanics, electronics, etc.) dictionaries and handbooks. Overall, the paper proposes a framework to go from management-oriented requirements to architecture-oriented requirements in which SVM is only applied in a small part of it.

Software requirements are mostly stated in natural text notations such as user stories which is making it hard for people to develop an accurate mental image of the most relevant entities and relationships. Lucassen *et al.* [30] introduced an automated method for visualizing requirements at different levels of granularity. Their visualization method from user stories consists of 1) the generation of an overview which provides a general context for understanding the dataset: **Parisa** ► suggestions for which part to remove for summarizing it?!◀

- Extract a set of relevant concepts from the user stories and their relationships
- Calculate the semantic similarity by using skip-gram implementation word2vec
- Utilize Ward's clustering algorithm to group all the concepts according to their similarity
- Identify the concept which is most similar to the collection of concepts in a cluster
- Generate inter-cluster relationships matrix
- Visualization Drawing

2) zooming in and out mechanisms and 3) filtering techniques to reduce the complexity of the data presentation. Possible anticipated applications of this visualization are: discovering missing relationships between clusters that may result in further user stories, teaching system functionality by exploring

simplified, manageable chunks, and analyzing expected system changes after introducing new sets of user stories.

2) *Structuring Documents*: Requirements of the system are usually presented in natural language documents. These documents require to be properly structured for a better overall understanding of the requirements. For this purpose, the document should be organized with independent sections which each one contains conceptually connected requirements [13]. Moreover, technical review is a usual way to guarantee the quality in natural language specifications. However, extensive and comprehensive specifications make it problematic for reviewers to find defects, especially consistency or completeness ones. Therefore, use of machine learning algorithms can support reviewers with their work by automatically classifying and clustering the information that is spread over many sections of many documents [33].

Duan *et al.* [12] used hierarchical automated clustering technique for detecting cross-cutting concerns as it is beneficial for the process of requirements analysis and architectural design. The reported experiments in this paper were supported by two tool sets, Poirot, a web-based tool designed to generate traces between various software engineering artifacts which was applied to compute similarity scores between requirements and a developed prototype tool, capable of reading structured requirements specification and generated similarity scores and then clustering requirements.

Requirements engineering process results are usually documented in the natural language specifications. In most cases, these documents not only contain requirements but also some additional information such as explanations, summaries, and figures. As it is important to differentiate between relevant requirements and other auxiliary content it is often the case that requirements engineers manually label each element of the specification document. Winkler *et al.* [46] applied convolutional neural networks to automatically classify content elements of a natural language requirements specification as “requirement” or “information”. **Parisa** ►requirements elicitation or management?◀ Their approach increases the quality of requirements specifications as it distinguishes important content for activities like test and etc. For converting natural language into a vector representation word2vec method is used. A set of 10000 content elements extracted from 89 requirements specifications of an industry partner used for training the network through the use of Tensorflow library using stochastic gradient descent.

For having a better understanding the natural language requirements specification documents should be properly structured. two quality characteristics of such a document are requirements relatedness which is each requirement is conceptually connected with the requirements in the same section and sections independence which is each section is conceptually separated from the others. based on Ferrari *et al.* [13] automatically recognizing the sections

in the document that need requirements relatedness and sections independence may help enhance the document structure. The authors defined a novel algorithm named Sliding Head-Tail Component (S-HTC) for clustering the requirements according to relatedness (the algorithm is based on known distance - Jaccard similarity metric, Levenshtein distance and, the convex combination between σ_{jac} and σ_{lev}). The algorithm groups together similar requirements that are contiguous in the requirements document. The effectiveness of the algorithm was evaluated with a test on requirements standard of a railway domain (583 requirements).

Based on Rauf *et al.* [35] software specification documents usually contain instances of logical structures, such as business rules, use cases, and functional requirements. Automated identification and extraction of these instances will benefit requirements management features, like automated traceability, template conformance checking, and guided editing. The authors planned a framework that gets requirements documents as an input and tries to develop a template for the general structure of it by specifying logical structures in terms of their content, textual rendering, and variability and then the extracting the instances of such structures from rich-text documents.

Ott *et al.* [33] automatically classified and extracted requirements with related information which are spread over many sections over many documents by the use of Multinomial Naive Bayes and Support Vector Machines classification algorithms as it will be helpful for reviewers with their work. As their input, they have used two German automotive specifications (Mercedes-Benz) which describe the functional and non-functional requirements of a Doors Closure Module (DCU). A specification and its referenced documents often sum up to 3,000 pages at Mercedes-Benz. Their method collects requirements of related information into classes, which they call topic landscape and later they built a tool, ReCaRe(Review with Categorized Requirements) which is the realization of the topic landscape based on eclipse with a data connection to IBM Rational DOORS.

III. DISCUSSION

- Not better than SOA, comparisons with humans, comparison with non-ML techniques
- There are a number of datasets available
- Clustering is on par with classification
- There are more contributions on Elicitation and Analysis

Levi ►internal vs external◀

IV. THREATS TO VALIDITY

The validity of the study might be affected by the coverage of the search results, bias on study selection, and inaccuracy of data extraction.

| | Themes | Contributions | ML Task | ML Model Types | Datasets Used |
|---|---|--|---|--|---|
| S | External | [18](○) [45] [21](+) [28] [20] [4] | Classification / Clustering | (Multinomial) Naïve Bayes; Support Vector Machines | Online reviews Skiweb data |
| | Non-Functional | [39](○) [5](○) | Classification | k-Nearest Neighbors | Open Source E-commerce Siemens Logistics |
| | Functional | [43](○) | Classification | Bi-Directional Long Short-Term Memory Conditional Random Field Network | E-commerce S |
| | Functional & Non-Functional Prioritization | [29] [8] [27] [1] [14] [44] [9] [2] [3] | Classification Classification | Case Based Ranking J48 DT | Software Requ |
| | Security Contextuals | | | | |
| V | Traceability | [15](○) [41](○) | Classification / Reinforcement Learning | Naive Bayesian / J48 Decision-Tree | Open Source C Open Source F |
| M | Visualization | [36](○) [34](○) [30](○) | Classification / Clustering | Support Vector Machines / Ward's method | - |
| | Structuring | [12](○) [46](○) [35](○) [13](○) [33](○) | Classification / Clustering | Multinomial Naive Bayes / Support Vector Machines Convolutional Neural Networks Sliding Head-Tail Componen Clustering Hierarchical Clustering | International U Mercedes-Benz |

Legend: (+) improves the state of the art; (-) comparable to or worse than state of the art; (○) no information on how the approach relates to the state of the art

TABLE I: Contributions and ML tasks related to each theme within each RE approach.

a) *Study Coverage.*: Some relevant studies could be missing in our study due to inadequate search strings or missing databases. To cope with this threat, the data preparation was based on a systematic method.

b) *Study Selection Bias.*: Study assessment might be biased by researchers. To mitigate for this threat, a set of include and exclude criteria was predefined and researchers assessed the title and abstract of the papers based on them to steer the assessment.

c) *Inaccuracy of Data Extraction.*: Also the data extraction process might be biased by researchers. To mitigate for this threat, the selection of data items was strictly driven by the research questions. Moreover, assignments were marked by the researchers depending on their confidence level. Low-confidence assignments were discussed between the authors until a consensus was reached.

V. CONCLUSION

REFERENCES

- [1] Z. S. H. Abad, O. Karras, P. Ghazi, M. Glinz, G. Ruhe, and K. Schneider. What works better? A study of classifying requirements. *CoRR*, abs/1707.02358, 2017.
- [2] P. Avesani, C. Bazzanella, A. Perini, and A. Susi. Facing scalability issues in requirements prioritization with machine learning techniques. In *13th IEEE International Conference on Requirements Engineering (RE'05)*, pages 297–305, Aug 2005.
- [3] P. Avesani, A. Perini, A. Siena, and A. Susi. Goals at risk? machine learning at support of early assessment. In *2015 IEEE 23rd International Requirements Engineering Conference (RE)*, pages 252–255, Aug 2015.
- [4] C. Castro-Herrera, C. Duan, J. Cleland-Huang, and B. Mobasher. A recommender system for requirements elicitation in large-scale software projects. In *Proceedings of the 2009 ACM Symposium on Applied Computing, SAC '09*, pages 1419–1426, New York, NY, USA, 2009. ACM.
- [5] J. Cleland-Huang, R. Settini, X. Zou, and P. Solc. Automated classification of non-functional requirements. *Requirements Engineering*, 12(2):103–120, Apr 2007.
- [6] J. Coughlan and R. D. Macredie. Effective communication in requirements elicitation: A comparison of methodologies. *Requir. Eng.*, 7(2):47–60, June 2002.
- [7] A. M. Davis. *Software Requirements: Objects, Functions, and States*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1993.
- [8] R. Deocadez, R. Harrison, and D. Rodriguez. Automatically classifying requirements from app stores: A preliminary study. In *2017 IEEE 25th International Requirements Engineering Conference Workshops (REW)*, pages 367–371, Sept 2017.
- [9] S. Dhingra, S. G. M. Madan, and M. R. Selection of prioritization technique for software requirement using fuzzy logic and decision tree. In *2016 Online International Conference on Green Engineering and Technologies (IC-GET)*, pages 1–11, Nov 2016.
- [10] P. Domingos. *The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World*. Basic Books, 2015.
- [11] L. Dong, X. Zhang, N. Ye, and X. Wan. Research on user requirements elicitation using text association rule. In *Intelligence Information Processing and Trusted Computing (IPTC), 2010 International Symposium on*, pages 357–359. IEEE, 2010.
- [12] C. Duan and J. Cleland-Huang. A clustering technique for early detection of dominant and recessive cross-cutting concerns. In *Proceedings of the Early Aspects at ICSE: Workshops in Aspect-Oriented Requirements Engineering and Architecture Design, EARLYASPECTS '07*, pages 1–, Washington, DC, USA, 2007. IEEE Computer Society.
- [13] A. Ferrari, S. Gnesi, and G. Tolomei. Using clustering to improve the structure of natural language requirements documents. In J. Doerr and A. L. Opdahl, editors, *Requirements Engineering: Foundation for Software Quality*, pages 34–49, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [14] F. Garzoli, D. Croce, M. Nardini, F. Ciambra, and R. Basili. Robust requirements analysis in complex systems through machine learning. In A. Moschitti and B. Plank, editors, *Trustworthy Eternal Systems via Evolving Software, Data and Knowledge*, pages 44–58, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [15] V. Gervasi and D. Zowghi. Mining requirements links. In D. Berry and X. Franch, editors, *Requirements Engineering: Foundation for Software Quality*, pages 196–201, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [16] M. Glinz. On non-functional requirements. In *15th IEEE International Requirements Engineering Conference (RE 2007)*, pages 21–26, Oct 2007.
- [17] O. C. Z. Gotel and C. W. Finkelstein. An analysis of the requirements traceability problem. In *Proceedings of IEEE International Conference on Requirements Engineering*, pages 94–101, Apr 1994.
- [18] E. Guzman, M. Ibrahim, and M. Glinz. A little bird told me: Mining tweets for requirements and software evolution. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 11–20, Sept 2017.
- [19] C. Hollis and T. Bhowmik. Automated support to capture verbal just-in-time requirements in agile development: A practitioner view. In *2017 IEEE 25th International Requirements Engineering Conference Workshops (REW)*, pages 419–422, Sept 2017.
- [20] N. Jha and A. Mahmoud. Mining user requirements from application store reviews using frame semantics. In P. Grünbacher and A. Perini, editors, *Requirements Engineering: Foundation for Software Quality*, pages 273–287, Cham, 2017. Springer International Publishing.
- [21] W. Jiang, H. Ruan, L. Zhang, P. Lew, and J. Jiang. For user-driven software evolution: Requirements elicitation derived from mining online reviews. In V. S. Tseng, T. B. Ho, Z.-H. Zhou, A. L. P. Chen, and H.-

- Y. Kao, editors, *Advances in Knowledge Discovery and Data Mining*, pages 584–595, Cham, 2014. Springer International Publishing.
- [22] R. Jindal, R. Malhotra, and A. Jain. Automated classification of security requirements. In *2016 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, pages 2027–2033, Sept 2016.
- [23] H. Kaiya, Y. Shimizu, H. Yasui, K. Kaijiri, and M. Saeki. Enhancing domain knowledge for requirements elicitation with web mining. In *2010 Asia Pacific Software Engineering Conference*, pages 3–12, Nov 2010.
- [24] J. Karlsson and K. Ryan. A cost-value approach for prioritizing requirements. *IEEE Softw.*, 14(5):67–74, Sept. 1997.
- [25] E. Knauss, S. Houmb, K. Schneider, S. Islam, and J. Jürjens. Supporting requirements engineers in recognising security issues. In D. Berry and X. Franch, editors, *Requirements Engineering: Foundation for Software Quality*, pages 4–18, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [26] G. Kotonya and I. Sommerville. *Requirements Engineering: Processes and Techniques*. Wiley Publishing, 1st edition, 1998.
- [27] Z. Kurtanović and W. Maalej. Automatically classifying functional and non-functional requirements using supervised machine learning. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 490–495, Sept 2017.
- [28] D. S. Lange. Text classification and machine learning support for requirements analysis using blogs. In *Innovations for Requirement Analysis. From Stakeholders’ Needs to Formal Designs*, pages 182–195, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.
- [29] M. Lu and P. Liang. Automatic classification of non-functional requirements from augmented app user reviews. In *Proceedings of the 21st International Conference on Evaluation and Assessment in Software Engineering*, EASE’17, pages 344–353, New York, NY, USA, 2017. ACM.
- [30] G. Lucassen, F. Dalpiaz, J. M. E. M. van der Werf, and S. Brinkkemper. Visualizing user story requirements at multiple granularity levels via semantic relatedness. In I. Comyn-Wattiau, K. Tanaka, I.-Y. Song, S. Yamamoto, and M. Saeki, editors, *Conceptual Modeling*, pages 463–478, Cham, 2016. Springer International Publishing.
- [31] W. Maalej and H. Nabil. Bug report, feature request, or simply praise? on automatically classifying app reviews. In *2015 IEEE 23rd International Requirements Engineering Conference (RE)*, pages 116–125, Aug 2015.
- [32] B. Nuseibeh and S. Easterbrook. Requirements engineering: A roadmap. In *Proceedings of the Conference on The Future of Software Engineering*, ICSE ’00, pages 35–46, New York, NY, USA, 2000. ACM.
- [33] D. Ott. Automatic requirement categorization of large natural language specifications at mercedes-benz for review improvements. In J. Doerr and A. L. Opdahl, editors, *Requirements Engineering: Foundation for Software Quality*, pages 50–64, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [34] R. Pinquicé, P. Véron, F. Segonds, and N. Croué. A collaborative requirement mining framework to support oems. In Y. Luo, editor, *Cooperative Design, Visualization, and Engineering*, pages 105–114, Cham, 2015. Springer International Publishing.
- [35] R. Rauf, M. Antkiewicz, and K. Czarnecki. Logical structure extraction from software requirements documents. In *2011 IEEE 19th International Requirements Engineering Conference*, pages 101–110, Aug 2011.
- [36] S. Reddivari, Z. Chen, and N. Niu. Recvisu: A tool for clustering-based visual exploration of requirements. In *2012 20th IEEE International Requirements Engineering Conference (RE)*, pages 327–328, Sept 2012.
- [37] M. Riaz, J. King, J. Slankas, and L. Williams. Hidden in plain sight: Automatically identifying security requirements from natural language artifacts. In *2014 IEEE 22nd International Requirements Engineering Conference (RE)*, pages 183–192, Aug 2014.
- [38] T. L. Saaty. Decision making with the analytic hierarchy process. *International journal of services sciences*, 1(1):83–98, 2008.
- [39] J. Slankas and L. Williams. Automated extraction of non-functional requirements in available documentation. In *2013 1st International Workshop on Natural Language Analysis in Software Engineering (NaturalLiSE)*, pages 9–16, May 2013.
- [40] C. Standish Group. The standish group report, 2014.
- [41] H. Sultanov and J. H. Hayes. Application of reinforcement learning to requirements engineering: requirements tracing. In *2013 21st IEEE International Requirements Engineering Conference (RE)*, pages 52–61, July 2013.
- [42] A. Van Lamsweerde. Goal-oriented requirements engineering: A guided tour. In *Proceedings of the Fifth IEEE International Symposium on Requirements Engineering*, RE ’01, pages 249–, Washington, DC, USA, 2001. IEEE Computer Society.
- [43] Y. Wang and J. Zhang. Experiment on automatic functional requirements analysis with the efrs semantic cases. In *2016 International Conference on Progress in Informatics and Computing (PIC)*, pages 636–642, Dec 2016.
- [44] M. Wieloch, S. Amornborvornwong, and J. Cleland-Huang. Trace-by-classification: A machine learning approach to generate trace links for frequently occurring software artifacts. In *2013 7th International Workshop on Traceability in Emerging Forms of Software Engineering (TEFSE)*, pages 110–114, May 2013.
- [45] G. Williams and A. Mahmoud. Mining twitter feeds for software user requirements. In *2017 IEEE 25th International Requirements Engineering Conference (RE)*, pages 1–10, Sept 2017.
- [46] J. Winkler and A. Vogelsang. Automatic classification of requirements based on convolutional neural networks. In *2016 IEEE 24th International Requirements Engineering Conference Workshops (REW)*, pages 39–45, Sept 2016.
- [47] P. Zave. Classification of research efforts in requirements engineering. *ACM Comput. Surv.*, 29(4):315–321, Dec. 1997.
- [48] D. Zowghi and C. Coulin. *Requirements Elicitation: A Survey of Techniques, Approaches, and Tools*, pages 19–46. Springer Berlin Heidelberg, Berlin, Heidelberg, 2005.