



23 offices, nor being accurately represented by man-made queries.

24 **Keywords:** ICT in construction; NLP; Deep learning; Information management.

## 25 **1. Introduction**

### 26 **1.1. Research background**

27 Information and communication technology (ICT) has been recognized as a key determinant to  
28 improve the level of coordination and collaboration in the architectural, engineering, and  
29 construction (AEC) industry (Davies and Harty, 2013; Wu et al., 2017). Yet, compared with other  
30 industries, the overall adoption rate of ICT in the AEC industry is low (Ahuja et al., 2009), and  
31 only a few number of regular and conventional ICTs such as 2D drawings are widely adopted.  
32 Regardless of the widely recognized benefits, most of the advanced and novel ICTs applications  
33 such as GPS, 4D modelling, BIM and mobiles are still incidentally employed in the industry  
34 (Ahuja et al., 2010; Dehlin and Olofsson, 2008; Frits, 2007; Li et al., 2019). One of the major  
35 barriers is that construction practitioners always lack technological knowledge about ICTC  
36 (Adriaanse et al., 2010; Sardroud, 2015).

37

38 Up to 80% technological information is exclusively provide by patents - recognized as one of the  
39 most valuable resources for technical analysis (Chiarello et al., 2018; Hoetker and Agarwal, 2007;  
40 Terragno, 1979). The content archived in a patent document normally expresses scientific and  
41 technological information for the technology application in terms of main machines and  
42 approaches involved, basic functions of the application, process whereby the application  
43 implements, and solutions to problems (Intarakumnerd and Charoenporn, 2015). Therefore, a  
44 corpus of patents that widely covers the inventions of ICTC is a valuable database, not only  
45 providing a dictionary for accessing ICTC, but also identifying problems to be solved by the state

46 of art ICTC inventions and recognizing all possible specific embodiments of ICTC (El-Ghandour  
 47 and Al-Hussein, 2004).

48  
 49 However, such a corpus of ICTC patents does not exist. Table 1 provides the existing patent classes  
 50 for the AEC industry in the three major patent offices, including World Intellectual Property  
 51 Organization (WIPO), European Patent Office (EPO), and United States Patent and Trademark  
 52 Office (USPTO). Table 1 shows that none of the patent offices provide a searchable classification  
 53 for ICTC. Two offices, WIPO and USPTO provide a specific category of patents that are relevant  
 54 to the AEC industry, namely E and D25 respectively. The two classes focus on inventions about  
 55 building materials and fixed construction rather than information and communication technologies.

56

57 **Table 1 Existing classification schemes in the three major patent offices**

Classification Scheme	Organizations	The specific classification of patents related to the AEC industry
<b>International Patent Classification (IPC)</b>	World Intellectual Property Organization (WIPO)	E: Fixed Constructions E01. Construction of roads, railways, or bridges E02. Hydraulic engineering; foundations; soil-shifting E03. Water supply; sewerage E04. Building E05. Locks; keys; window or door fittings; safes E06. Doors, windows, shutters, or roller blinds, in general; ladders E21. Earth or rock drilling; mining E99. Subject matter not otherwise provided for in this section
<b>Cooperative Patent Classification (CPC)</b>	European Patent Office (EPO)	None D25: Building units and construction elements <ol style="list-style-type: none"> <li>1. Structure</li> <li>2. Prefabricated unit</li> <li>3. Stair, ladder, scaffold, or similar support</li> <li>4. Trellis or treillage unit</li> <li>5. Architectural stock material</li> </ol>
<b>United States Patent Classification (USPC)</b>	USPTO	

58

## 59 **1.2. The problem of retrieving ICTC patents by using query-based** 60 **methods**

61 In the absence of a searchable classification for ICTC in the patent offices, query-based methods  
62 (including the searching engines and other patent retrieval methods) became a possible way for  
63 users to retrieve the patents. These query-based methods aim to retrieve all documents that are  
64 relevant to a given patent application according to a query. Hence, the accuracy and coverage of  
65 retrieval results highly depends on the query (Zhang et al., 2018), which can be formed by a variety  
66 of items, such as keywords, citations, authors, granted year, application date, or combinations of  
67 them. The core technique lies in the query-based methods is query reformulation - converting the  
68 input query into new and more searchable queries (Alberts et al., 2017; Shalaby and Zadrozny,  
69 2019). The query reformulation, including query reduction techniques (Bouadjenek et al., 2015;  
70 Mahdabi et al., 2011), query expanding method (mainly by external dictionary and corpus or  
71 ontologies) (Azad and Deepak, 2019; Enesi et al., 2018; Tannebaum and Rauber, 2014), semantic-  
72 based methods (Girthana and Swamynathan, 2018), metadata-based methods (citations and  
73 classification) (Azad and Deepak, 2019; Giachanou et al., 2015; Mahdabi and Crestani, 2014), and  
74 interactive methods (Shalaby and Zadrozny, 2018), enriched the query-based methods and have  
75 obtained performance improvement in recent studies.

76  
77 However, gathering the corpus of ICTC patents may not achieve good performance by using the  
78 query-based methods, because it is extremely challenging to accurately represent and widely cover  
79 the ICTC patents by man-made queries. The patent retrieval tasks, including *prior-art search*,  
80 *patentability search* and *infringement search* (Zhang et al., 2018), aim to return a wide coverage  
81 of patent documents that are relevant to a patent application according to a query, helping potential  
82 patentees check and analyze relevant information before the patent application is granted.

83 Therefore, the queries are frequently used to represent a specific patent application rather than a  
 84 set of patents. ICTC, standing for a set of ICT applications that were invented with major  
 85 embodiments in the AEC industry (Ahuja et al., 2009; Alsafouri and Ayer, 2018), incorporates a  
 86 number of technologies which may vary with each other (for example, both BIM and RFID are  
 87 important ICT applications in the AEC industry, but they are totally different technologies).  
 88 Therefore, completely representing all the ICTC patents by a man-made query leaves a tough task  
 89 to return accurate results. Moreover, using a query combined by a number of items to represent all  
 90 the ICTC patents increases the irrelevant instances returned due to polysemy (the same spellings  
 91 may have two or more different meanings).

92  
 93 This study performs two trails for retrieving ICTC patents from the USPTO website (USPTO,  
 94 2007) based on a query combined by a number of items. Table 2 shows the results using this query-  
 95 based method, and 50 patents were randomly selected to manually check the accuracy - the  
 96 proportion the ICTC patents occupy all the retrieved patents. Even though a complicated  
 97 combination of items were used to search ICTC patents, the accuracy is low. Moreover, most of  
 98 the latent users in the construction practice are non-experts, who may not be able to perform such  
 99 a searching task that is complex and time-consuming (Liu et al., 2011).

100  
 101 **Table 2 Searching results by using the search engine in USPTO**

Querying strategies	Matched results	Accuracy
<b>Query items:</b> CPC Classification Class and topic (matching input keywords within patent titles, abstracts and descriptions)		
Strategy 1 CPC Classification Class: ICT-related classes, including <b>H04</b> - <i>electric communication technique</i> ; <b>G06</b> - <i>computing, calculating or counting</i> ; <b>H01P</b> - <i>waveguides, resonators, lines, or other devices of the waveguide type</i> ; <b>H01Q</b> - <i>antennas, i.e. radio aeriels</i> ; <b>G01S</b> - <i>radio direction-finding, radio navigation,</i>	Collection 1: 5311 patents	8%

---

*determining distance or velocity by use of radio waves, locating or presence-detecting by use of the reflection or reradiation of radio waves or analogous arrangements using other waves; G08B - signalling or calling systems, order telegraphs or alarm systems; G08C - transmission systems for measured values, control or similar signals; G11B - information storage based on relative movement between record carrier and transducer.*

Keywords: AEC domain terms, including *construction project, project management, infrastructure project, civil engineering and transportation project.* (Flyvbjerg, 2014; Greiman, 2013; Levitt, 2007; Mok et al., 2015; Zidane et al., 2013)

---

**Query item:** Topic

Strategy 2	Keywords: ICT-related terms ( <i>Radio frequency identification (RFID), 3D laser scanning, Quick response, NFC, Augmented reality (AR), Mobile computing, Wireless connection (Wi-Fi) and Robotics(drones)</i> ) (Ahuja et al., 2009; Alsafouri and Ayer, 2018; Li et al., 2016) and AEC domain terms	Collection 2: 922 patents	12%
------------	---	------------------------------	-----

---

102

### 103 **1.3. Research Objectives**

104 Given the aforementioned constraints of existing query-based methods, this study develops a  
105 binary classifier to automatically identify whether a patent is relevant to ICTC, and thus accurately  
106 screening a corpus of ICTC patents from the primarily searched results containing a number of  
107 irrelevant patents. Therefore, this study treats the task of screening ICTC patents as a classification  
108 task rather than a retrieval task. A large number of studies have investigated patent classification,  
109 and most of them emphasized the use of traditional machine learning (i.e. SVM and Bayes) and  
110 text mining techniques (i.e. n-gram and stop-words removal) (Li et al., 2012; Wu et al., 2010).  
111 Alternatively, this study resorts to the techniques from the realm of NLP and deep learning. On  
112 one hand, NLP techniques provide a smart way to process textual data (Kurdi, 2017), saving time  
113 and avoiding personal bias in analysis processes (Agrawal and Henderson, 2002; Bell et al., 2009;  
114 Cassetta et al., 2017; Choi et al., 2012; Gwak and Sohn, 2018), especially when the volume of a  
115 text is large (Shekarpour et al., 2015; Silva et al., 2016). On the other hand, a deep learning model,  
116 Multi-Layer Perceptron (MLP) is developed to learn the relations between the input features and

117 outputs. Deep learning is the most state-of-the-art approach with significant performance  
118 improvement in NLP tasks. Compared with the algorithms and statistics of the machine learning  
119 models, the deep learning models are organized by multiple layers of neural networks. Each layer  
120 consists of neurons, receiving signals from the former layer and passing converted signals by  
121 activation functions to the subsequent layer (Riedmiller, 1994). With the multiple layers of neural  
122 networks, the whole deep learning model can address highly non-linear associations between the  
123 representations and the outputs (Wang et al., 2016), whereas the machine learning algorithms can  
124 only examine linear relations.

## 125 **2. Related work**

126 Several attempts have been made to establish classifiers for automatic patent classification  
127 (Chakrabarti et al., 1998; Smith, 2002; Venugopalan and Rai, 2015). Most of the studies, at the  
128 beginning, extracted the features from the structured data or metadata, such as keywords and  
129 citations (Michel and Bettels, 2001; Perez-Molina, 2018). In the recent decade, scholars prefer  
130 using unstructured data (Cambria and White, 2014; Collobert et al., 2011; Gimpel et al., 2011).  
131 These studies typically have three key steps: processing the textual data, vectorizing the patents  
132 and using machine learning methods to train the models. Focusing on the three key steps, this  
133 section describes a synopsis of the relevant literature.

### 134 **2.1. NLP techniques for processing textual data**

135 A large volume of unformatted texts exist in the current web 2.0 era (Ittoo et al., 2016). The chunk  
136 information is mainly unstructured and thus cannot be processed by machine-tractable ways  
137 (Cambria and White, 2014) that structured data can be. Processing the unstructured data is regarded  
138 as one of the most time-consuming step of text classification tasks (Munková et al., 2013). The

139 major object of the processing is to clean and format the raw textual data, which can largely  
140 eliminate noisy features for further vectorization (Haddi et al., 2013). Many NLP tools have been  
141 introduced, such as stop and common words removal, tokenization, lemmatization and stemming  
142 (Aggarwal and Reddy, 2013). Two typical tools are part-of-speech (POS) and Named Entity  
143 Recognition (NER), which can recognize and process syntax information (grammatical meanings)  
144 (Collobert et al., 2011; Gimpel et al., 2011) and named entities respectively (Nadeau and Sekine,  
145 2007).

## 146 **2.2. Vectorizing methods**

147 With regard to the vectorization, a number of algorithms have been developed to convert the  
148 textual data into vectors. Bag-of-words (BOW), topic models and subject–action–object (SAO)  
149 have been used in recent patent classification studies (Li et al., 2018; Venugopalan and Rai, 2015).  
150 Traditional BOW models typically construct the feature space vectors in which each position is  
151 occupied by a term or a phrase (Forman, 2002). Its measurements include n-grams, bi-grams, and  
152 word frequency to identify phrases from the texts (Onan et al., 2016), depending on how the  
153 phrases were counted. Although BOW models are simple and may generate a large number of  
154 features, they remain the most effective feature selection method (Mirończuk and Protasiewicz,  
155 2018; Onan et al., 2016). Topic model and subject–action–object (SAO) were mainly developed  
156 to solve the high dimension problem, replacing the BOW features by latent topics (Kaplan and  
157 Vakili, 2013) or SAO structures (Gerken, 2012).

## 158 **2.3. Supervised learning models for patent classification**

159 To date, most of the patent classifiers are trained by machine learning models, such as SVM, Naive  
160 Bayes, and k-nearest-neighbor (KNN). The accuracy was relatively low in earlier studies (around  
161 70%) (Saiki et al., 2006), but an increasing trend has been observed when feature selection models



162 and NLP techniques were used to extract useful features from unstructured texts, achieving  
163 accuracy around 85% (Venugopalan and Rai, 2015; Wu et al., 2010). In the recent decade, the  
164 widespread use of deep learning models has led to notable success. Deep learning models have  
165 been developed and adopted in a variety of fields, such as natural language understanding, video  
166 and image recognition and game of GO (Al Rahhal et al., 2018; Cocarascu and Toni, 2018; Silver  
167 et al., 2016). Those deep learning models were always developed with complex and elaborate  
168 architecture in which multiple layers of neural networks were well structured. However, only the  
169 artificial neural network (ANN) with one layer neurons was applied to patent classification, and  
170 the accuracy was relatively low, with 75% (Li et al., 2012).

171  
172 In addition, most of the research have sought to classify patents into pre-defined classes that  
173 already existed in the classification schemes in patent offices, such as IPC and European  
174 classification code hierarchy (Fall et al., 2003a; Fall et al., 2003b). Such expositions provide  
175 automatic and efficient methods for inventors and examiners to label the new patents with existing  
176 classes, but do not provide opportunities to advance the understanding of the real world in the  
177 target field. Although using deep learning models may lead to better performance than traditional  
178 machine learning models, rare studies employ deep learning models in automatic patent  
179 classification tasks.

### 180 **3. The proposed approach integrating deep** 181 **learning and NLP techniques**

182 To screen ICTC patents from a patent collection, a binary classifier is developed to classify pieces  
183 of patents into two classes: ICTC-related or not. Figure 1 shows the overall procedure of the  
184 approach to achieve the classifier. The first step is to collect a database for training, incorporating

185 the full texts of the instances annotated with target labels (the two classes). Then, NLP tools are  
 186 used to process textual data to achieve clean texts. Based on the processed texts, N-gram and Tf-  
 187 Idf algorithms are employed for the vectorization to represent each of the patents as a numerical  
 188 vector that could be fed into the MLP that would be trained by gradient descent in which the  
 189 hyperparameters are tuned. At last, a validation experiment is conducted by means of k fold cross-  
 190 validation in two datasets. The succeeding sections discuss these steps.

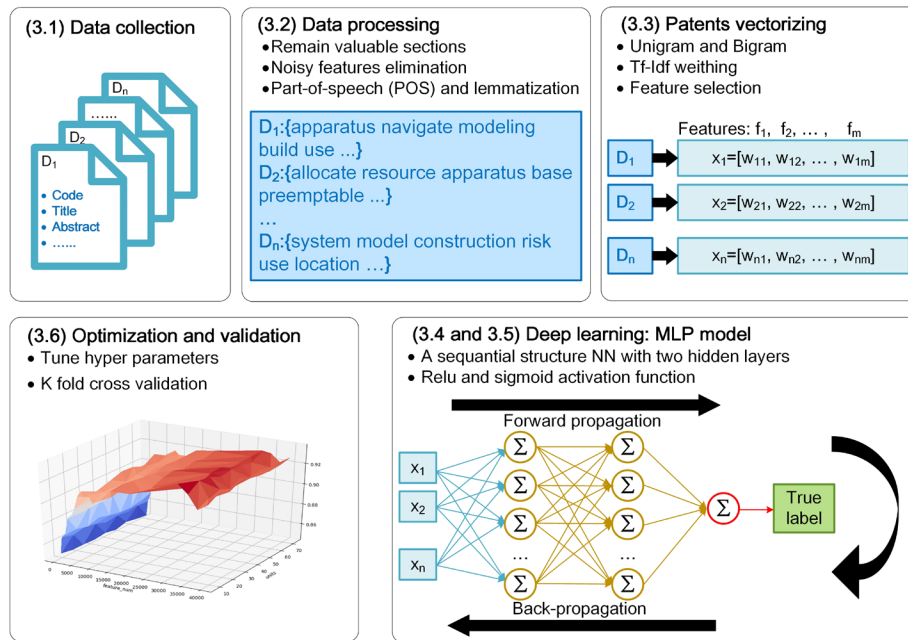


Figure 1 Overall procedure of the approach

192

### 193 3.1. Data collection and annotation

194 The target of this step is to obtain training data - the full texts of the patents that are manually  
 195 labeled as *ICTC* or *non-ICTC*. All the required patent text were crawled from USPTO, because (1)  
 196 USPTO is the largest international patent grant office, and (2) USPTO is recognized as the most  
 197 representative database to analyze the technological knowledge, providing patents that are well  
 198 written and structured according to its requirement (Wang, 2018). The authors retrieved the patents  
 199 on July 30, 2018. Totally, we have collected and annotated 348 patents as the dataset for further  
 200 training and testing. The detailed processes are described in the following two paragraphs.

201  
202  
203  
204  
205  
206  
207

Figure 2 depicts the data collection and annotation process, whereby patents were collected and annotated as *ICTC* or *non-ICTC*. As for the *ICTC class*, patents were gathered in the following steps: (1) By querying search strategy 1 in Table 1 (ICT classes and AEC domain terms), 5311 patents were obtained in collection 1. (2) Totally 1500 patents were randomly selected from the 5311 patents. (3) Through the process of manually checking<sup>1</sup>, 174 patents were obtained as *ICTC class* from the 1500 patents.

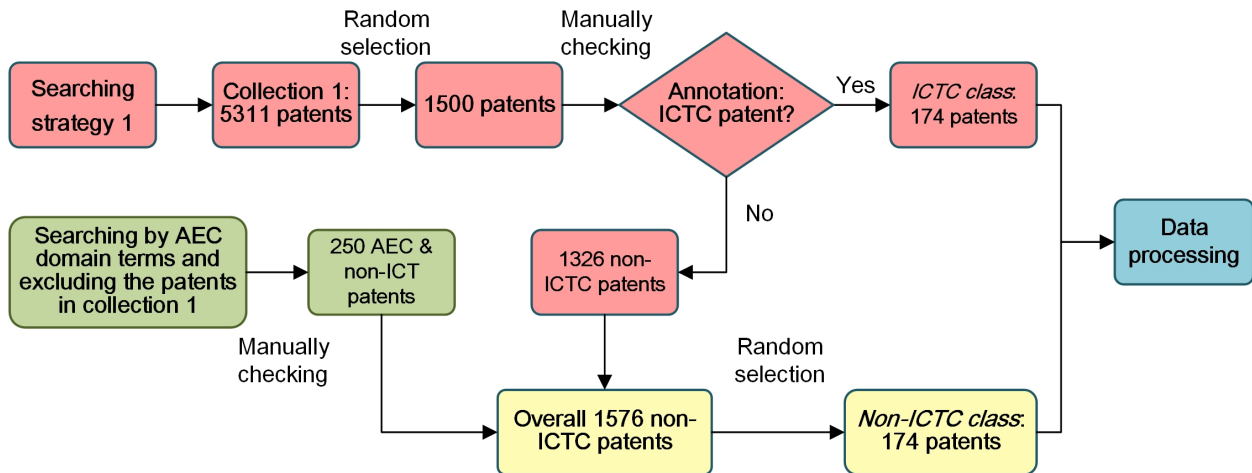


Figure 2 The process of data collection

210  
211  
212  
213  
214  
215

As for the *non-ICTC class*, the patents were collected from two different sources. One was through the annotation process mentioned above, in which 1326 patents were identified as *non-ICTC class*. The other was obtained from the patents retrieved by searching AEC domain terms and excluding patents in collection 1. This results in a combined collection of 1576 non-ICTC patents and 174 of

---

<sup>1</sup> The process of manually checking labels a patent as either *ICTC* or *non-ICTC*, performed by three Ph.D. students (their research directions are related to the AEC area) through in-depth reviewing of the title, abstract, claim and description. A patent can be labeled as *ICTC class* if the content expresses that the essence of the technology application is under the ICT scope and the AEC industry is a major embodiment in which the technology application can be implemented. To prevent mistakes as much as possible, two students annotated the patents independently, and the third student would make a judgment when the labels of a patent are inconsistent.

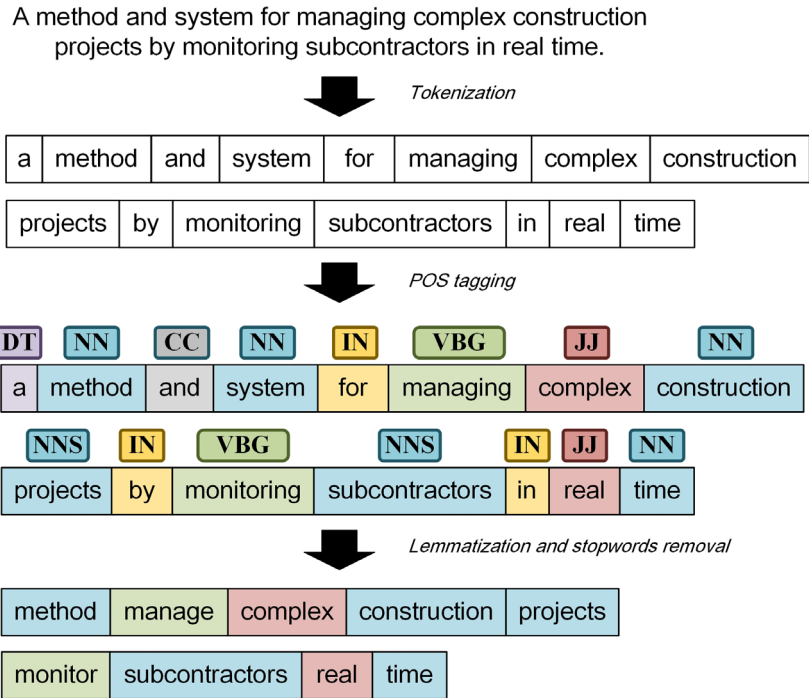
216 them were randomly selected as training instances for the *non-ICTC class*. The complex collection  
217 process has two advantages: (1) The *non-ICTC* class contains not only common ICTs that exclude  
218 ICTC patents, but also the technologies of the AEC industry that exclude ICTC patents. This can  
219 prevent the data over-fitting, thus generating a more generalized model that is able to distinguish  
220 ICTC patents from ICT, as well as from AEC patents; (2) This study uses the negative sampling to  
221 make the two classes have the same size, because the balanced size for each class is proven as a  
222 key factor to achieve high accuracy in training (Brown and Mues, 2012; Zhao et al., 2015).

### 223 **3.2. Data processing by NLP techniques**

224 The raw text of each patent contains several sections (i.e., code, title, abstract, CPC classes,  
225 inventors and countries and description). Among them, *title*, *abstract* and *claim* are frequently  
226 utilized and remained for further analysis in this study because they were recognized as useful  
227 items providing basic technological information (Niemann et al., 2017; Venugopalan and Rai,  
228 2015). *Title* and *abstract* convey the essence about the technology, which were always written in  
229 a restricted pattern within short content (Lee et al., 2013). In addition, *claim* defines the protection  
230 right of the invention, always providing articulated expressions about the technical boundaries and  
231 specifications (Niemann et al., 2017).

232  
233 The selected text of patents is raw data, which is pre-processed by NLP techniques for further  
234 analysis. Without pre-processing, the texts would contain a lot of noisy features (in a typical case,  
235 the number of features can be close to the number of words in the dictionary of the training  
236 instances), thus creating higher-dimension vectors. To process the selected raw text, this study  
237 employs three NLP techniques (Figure 3 plots the pre-processing procedure using these  
238 techniques): (1) Tokenization. For each raw sentence in the texts, tokenization is utilized to split

239 the sentence into words. In addition, all the words are converted to lowercase and punctuations are  
240 removed. Through this step, all the raw sentences would be replaced with sequenced and lowercase  
241 words. (2) POS tagging. In this step, each word is tagged with POS tag indicating its syntactic role  
242 (i.e. noun, adverb) according to the surrounding words. POS tagging plays a central role in text  
243 processing, which could increase the accuracy for lemmatization and stemming (Habash et al.,  
244 2009). (3) Lemmatization and stop-words removal. The purpose of this step is to correctly match  
245 the words with different forms, such as plural forms for nouns and presenting and past forms for  
246 verbs. Lemmatization transforms the different forms into the stem forms (root words). However,  
247 lemmatization may generate a number of mistakes without POS tagging. For example, “modeling”  
248 may be a present participle of a verb (with lemma “model”) or a noun (with lemma “modeling”)  
249 according to the context, and the lemma of noun “modeling” would be wrongly identified as  
250 “model” without the POS information (Vlachidis and Tudhope, 2016). This study utilizes NLTK  
251 toolkits to perform POS tagging and lemmatization (Bird and Loper, 2004). Moreover, stop-words  
252 (i.e., a, an, of, one, two, three and so on) are removed, because they are non-descriptive and do not  
253 convey any semantic meanings.  
254



**Figure 3 The processing procedure for textual data in patents**

### 3.3. Vectorizing patents

The processed patent texts have to be converted into numerical vectors that can be fed into MLP. This study adopts N-gram model with Tf-Idf weighting algorithm to vectorize the patents. N-gram considers the N words in a sequence as a feature, which has been proposed in the 1940s (Shannon, 1948) and has been employed in a large and growing body of literature (Bengio et al., 2003; Benson and Magee, 2013). In this case, two typical N-gram models, N = 1 (unigram) and 2 (bigram) are used to extract unigrams and bigrams as features from the patent texts, constituting the vocabulary with size v (overall v unigrams and bigrams are identified from the patent texts). A vector with v-dimension in which each position is the Tf-Idf (term frequency & inverse document frequency, see Sparck Jones (1972) for details) value of the feature in the vocabulary is generated to represent a patent. Another necessary step is to filter useful features because many of the features do not contribute to the training and prediction. This study, according to the Tf-Idf vectors, uses ANOVA F-value to select top features (number = K). In this study, K is set as a hyperparameter that would

270 be tuned in the optimization step.

271

272 This study adopts N-gram and Tf-Idf as the vectorizing approach but not the topic models or SAO  
273 structures, because (1) BOW and Tf-Idf have been widely used in NLP studies and have been  
274 proven as the prominent vectorizing algorithm due to the simplicity and effectiveness (García  
275 Adeva et al., 2014; Mirończuk and Protasiewicz, 2018; Pavlinek and Podgorelec, 2017); (2) Topic  
276 models and SAO structures are suitable in clustering or classification tasks that have more than  
277 two classes to be distinguished (Blei et al., 2003; Choi et al., 2012); (3) Topic models and SAO  
278 structures replace the N-grams with latent topics or subject-active-objective structures. This would  
279 generate vectors with much lower dimensions which is not necessary in this case, because the  
280 proposed MLP model may get better performance when the number of input features is large (Cakir  
281 and Yilmaz, 2014).

### 282 **3.4. MLP architecture**

283 To improve the performance and take the non-linear relations into consideration, this study  
284 proposes MLP to learn and train the complex relations between inputs and outputs. MLP is  
285 typically designed with a feed-forward-based architecture and back-propagation learning process  
286 (Rosenblatt, 1961). There is a number of neurons in the MLP, and each of them receives signals  
287 from the former layer and pass transformed signals by an activation function to the subsequent  
288 layer (Riedmiller, 1994). Although it is a general wisdom that deep learning models are better than  
289 machine learning models, neural network design and hyperparameter choice are more important  
290 than deep learning models themselves (Levy et al., 2015). This section describes the MLP  
291 architecture.

292

293 After the vectorizing, the input matrix in this study is  $X \in \mathbb{R}^{N \times F}$ , in which N and F represent the  
 294 number of instances and features respectively. Features are set as columns, and thus each patent is  
 295 reflected as a row vector  $x_i \in \mathbb{R}^{1 \times F}$ . The output is a column vector  $Y \in \mathbb{R}^{N \times 1}$ . The main target of  
 296 the MLP model is to obtain learned neurons in layers of neural networks that could predict from  
 297 X to Y. Figure 4 illustrates the architecture of the MLP consisting of four layers: one input layer,  
 298 two hidden layers and an output layer, labeled from layer 0 to layer 3. The weigh matrixes connect  
 299 the layers in sequence, and the neurons in the hidden and output layers are processing units,  
 300 embodied with activation functions to transform the inputs to outputs. The number of the neurons  
 301 in the input and output layers are set as N and 1, which are subject to the dimensions of the input  
 302 and output vectors. The numbers of the neurons in the hidden layers are set as hyperparameters.

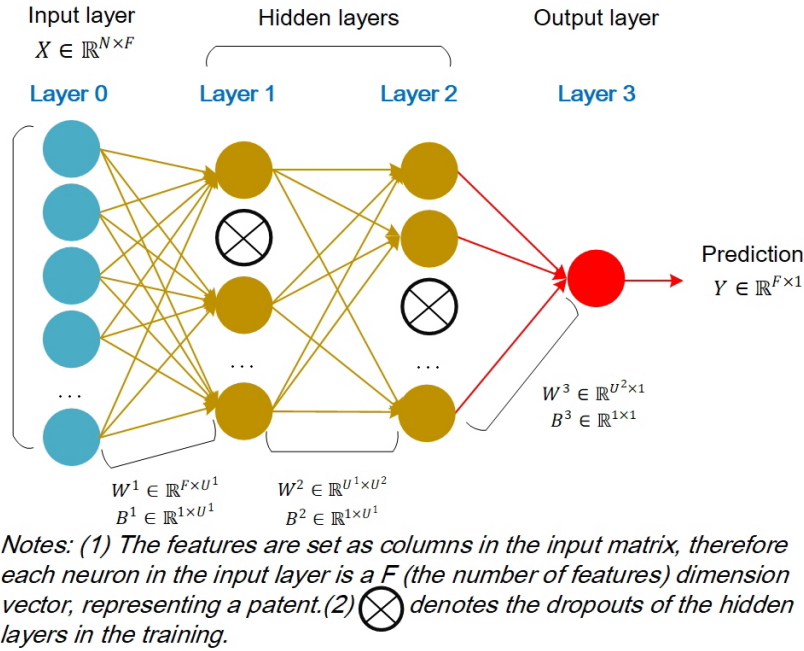


Figure 4 The MLP architecture

304  
 305  
 306 The MLP predicts the outputs based on the connection weights and the activation functions. In  
 307 specific, the  $j$ -th neuron in  $l$ -th layer transforms an output based on the following equations:

$$\begin{cases} h_i^l = f^l \left( \sum_{i=1}^{U^{l-1}} h_i^{l-1} w_{ij}^l + b^l \right), l = 1, 2 \\ h^l = h^3 = f^3 \left( \sum_{i=1}^{U^2} (h_i^2 w_i^3 + b_i^3) \right), l = 3 \end{cases} \quad (1)$$



309

310 where  $l$  represents the layer sequence,  $U^{l-1}$  indicates the number of neurons in the  $(l - 1)$ -th  
311 layer,  $x_i^{l-1}$  denotes the output of  $i$ -th neuron it receives,  $w_{ij}$  is the weight connecting  $x_i^{l-1}$  and  $j$ -  
312 th neuron in  $l$ -th layer, and  $b$  is the bias function for this neuron.  $f^l$  is the activation function in  $l$ -  
313 th layer. In this case, the two hidden layers (layer 1 and layer 2) and the output layer (layer 3) use  
314 *Rectified Linear Unit* (ReLU) and *Sigmoid* functions as the activation functions respectively. With  
315 the back-propagation process, the neurons in hidden and output layers can be trained with unique  
316 weight matrix and bias, producing different outputs according to the tasks (Garcia-Laencina et al.,  
317 2013). Moreover, the “Dropouts” is adopted in the hidden layers to prevent the overfitting  
318 (Srivastava et al., 2014).

### 319 **3.5. Model training by gradients and dropouts**

320 As mentioned above, the main task of MLP is to make the neurons to be learned, which could  
321 predict  $Y$  from  $X$ . The learning process is achieved by certain iterations, each of which is a loop  
322 consisting of a feed-forward and a back-propagation process (Haykin, 1999; Riedmiller, 1994). In  
323 the feed-forward process, the weights and bias in the hidden and output layers are randomly  
324 generated and propelled forward, calculating the output value  $h^3$  from input  $X$ . Since a sigmoid  
325 function is selected as the activation function in the output layer, the errors follow a logistic  
326 distribution between the predictions (with values between 0 and 1) and true labels (with values are  
327 only 0 or 1). The loss function is the following:

$$328 \quad J = -\sum_{n=1}^N y_n \log(h_n^3) + (1 - y_n) \log(1 - h_n^3) \quad (2)$$

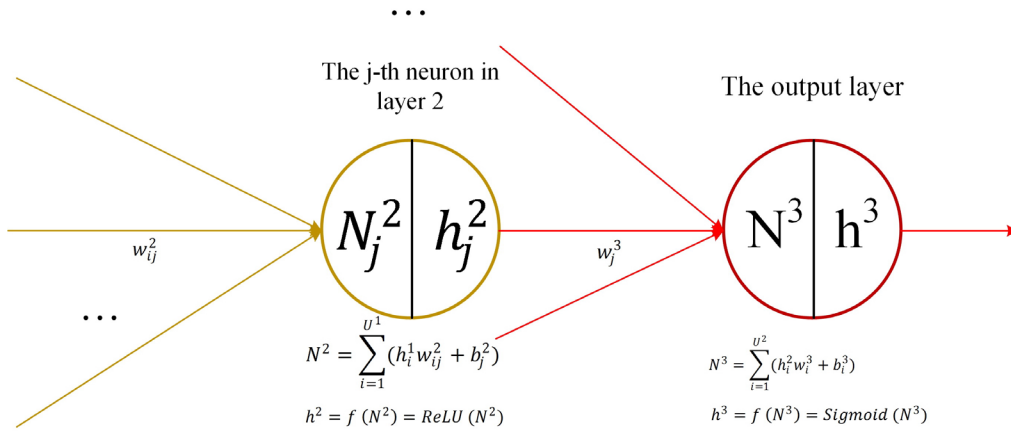
329 In the back-propagation, the parameters  $\theta$  (including all the weights and bias in hidden and output  
330 layers) would be updated by stochastic gradient descent. Two types of signals constitute the  
331 gradients: (1) global signals that can be computed from the derivatives, which transform the errors

332 from the loss function; (2) local signals that are the inputs from the former layer. The  $\theta$  would be  
 333 updated from back to front, as the gradients are computed from the loss value to the former layers,  
 334 one by one. For the clarity of the back-propagation process, this study illustrates the updating  
 335 process of  $w_{ij}^2$  and  $w_j^3$  in layer 2 and 3. Figure 5 shows the functions of the neurons in layer 2 and  
 336 3. The gradient of  $w_j^3$  ( $\nabla w_j^3$ ) is defined as the derivative from  $J$  to  $w_j^3$ , which could be computed  
 337 by the chain rule of derivatives:

$$338 \quad \nabla w_j^3 = \frac{dJ}{dw_j^3} = \left( \frac{dJ}{dh^3} \times \frac{dh^3}{dN^3} \right) \times \frac{dN^3}{dw_j^3} = f'(N^3) \times h_j^2 \quad (3)$$

$$339 \quad w_j^3_{new} = f(w_j^3_{old}, \nabla w_j^3) \quad (4)$$

340 where the  $f'(N^3)$  is the global signal that could be computed by the derivative with loss value,  $h_j^2$   
 341 is the local signal (the output of the j-th neuron in layer 2), and  $a$  is the learning rate that is pre-  
 342 defined.



344 **Figure 5 The neurons with input and activation functions in the last two layers**  
 345

346 Similar to layer 3,  $\nabla w_{ij}^2$  could be computed by the following:

$$347 \quad \nabla w_{ij}^2 = \frac{dJ}{dw_{ij}^2} = \frac{dJ}{dh^2} \times \frac{dh^2}{dN^2} \times \frac{dN^2}{dw_{ij}^2} = f'(N^2) w_j^3 f'(N^3) \times h_i^1 \quad (5)$$

348 
$$w_{ij}^{2new} = f(w_{ij}^{2old}, \nabla w_{ij}^2) \quad (6)$$

349 where  $f'(N^2)w_j^3 f'(N^3)$  is the global signal that is propagated from the loss value, and  $h_i^1$  is the  
350 local signal. The computations of other parameters, such as  $w$  and  $b$  are similar to equation (3) and  
351 (5). According to the gradients, the parameters could be updated by algorithms (equation (4) and  
352 (6)). Typical optimization algorithms include Stochastic Gradient descent (Robbins and Monro,  
353 1985), AdaGrad (Duchi et al., 2011), RMSProp (Tieleman and Hinton, 2012), and Adam (Kingma  
354 and Ba, 2014). This study applies the Adam algorithm as the optimizer, as it has been recognized  
355 as the most effective in most cases with less computation time.

356  
357 Dropouts is applied in the training process. “Dropouts” refers to temporarily eliminating some  
358 neurons and their incoming and outgoing connections in the neural networks. The dropped neurons  
359 are selected randomly based on a pre-defined ratio  $a$  ( $a=0.2$  in this case). In the back-propagation  
360 of a training loop, a new thinned neural network is achieved with the proportion of  $1 - a$  neurons  
361 remained. The parameters updating process would be implemented within the thinned neural  
362 network. In the feed-forward process of the subsequent loop, the removed neurons would turn on,  
363 and their parameters are obtained from the remaining neurons by a scale of  $1/a$ . which parameters  
364 are obtained from the remaining neurons by a scale of  $1/a$ . Therefore, training MLP with Dropouts  
365 can be regarded as training a larger number of thinned neural networks that share the same  
366 parameters. Such a training fashion effectively prevents neurons from co-adapting, and thus  
367 preventing the overfitting issues (Al Rahhal et al., 2018). As for the details of Dropouts, please see  
368 Al Rahhal et al. (2018).

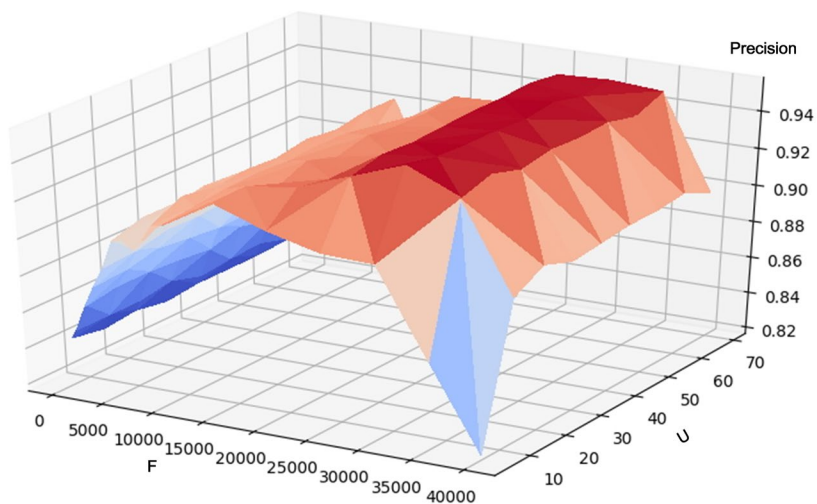
369  
370 After updating  $\theta$ , a loop with a feed-forward and back-propagation finishes, which would be

371 iterated in training. In this case, the maximum of epochs is set as 1000, and the consecutive tries  
372 of loss value without decrease is set two. The training process would iterate the loops until any of  
373 the above stop conditions is met. A small number of self-developed python programs are used to  
374 build, train, optimize and validate the model.

## 375 4. Results

### 376 4.1. Results of hyperparameters tuning

377 The purpose of hyperparameters tuning is to achieve an MLP model with the best performance by  
378 tuned hyperparameters. The range of the features is from 1000 to 40000, with steps of 1000 and  
379 2500 for F in (1000, 10000) and (10000, 40000) respectively. With regard to the number of units,  
380 this study adopts the measurement proposed by Fan et al. (2015), which proposed a range around  
381  $\sqrt{N + 1}$  (N denotes the number of neurons in the input layer). The resulting range of the number  
382 of units is from 5 to 69 and the step is set as 8. Figure 6 shows the hyperparameters tuning process.  
383 The MLP model reaches the highest accuracy when F is 30000 and U is 13.



384

385

386

**Figure 6 The hyperparameters tuning process**

## 387 4.2. Validation

### 388 4.2.1. Validation methods

389 This study validates the proposed approach not only over the dataset in which 348 patents (labeled  
390 as *ICTC* or non-*ICTC*) were collected from USPTO, but also the patents from Derwent Innovations  
391 Index (DIX). The additional validation over DIX patents can evaluate the performance of the  
392 proposed model in processing texts that were written in different genres. Following prior machine  
393 learning studies (Sokolova and Lapalme, 2009), this study utilizes precision, recall, F-score to  
394 validate the deep learning model based on true positives (TP), false positives (FP) and false  
395 negatives (FN). Generally, TP is the number of instances the model correctly predicted. FP denotes  
396 the instances the model incorrectly predicted. FN reflects the number of instances the model failed  
397 to predict. The precision, recall, and F-score can be computed by

$$398 \quad P = \frac{TP}{TP+FP}, R = \frac{TP}{TP+FN}, F1 = \frac{2 \times P \times R}{P + R} \quad (7)$$

399  
400 Specifically, as for the 348 USPTO patents in the dataset, we used the k-fold cross-validation along  
401 with the training process. In the training process, all the dumping data would be randomly split  
402 into k folds with the same size, and one of them is set as test instances and others are used for  
403 training. Such a training process is performed in k times, each of which has a different fold for  
404 testing and a different composition of k-1 folds for training. The final validation value is the  
405 average of k validation values. In this way, k-fold cross-validation prevents the bias in data  
406 selection and ensures the measures of the performances with objections (Friedman et al., 2001). In  
407 this study, k is set as 5, and all the annotated data (totally 348 USPTO patents annotated as *ICTC*  
408 or non-*ICTC* class were obtained in Section 3.1) were randomly divided into 5 folds. For each  
409 training round, four folds consist of the training collection and the other fold consists of the testing

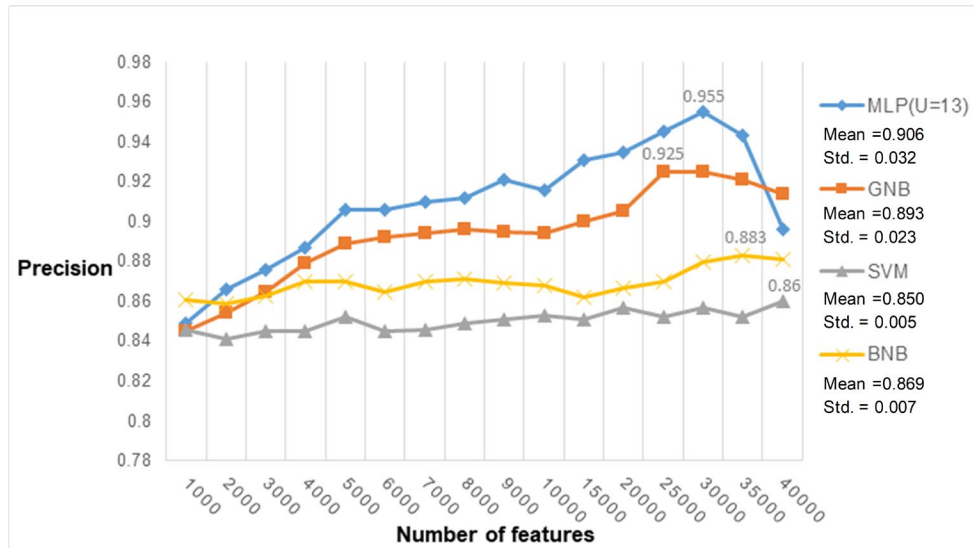
410 collection.

411

412 Besides the annotated data, this study collected and randomly selected 200 patents from Derwent  
413 Innovations Index (DIX) as an additional testing dataset, where the patents were written by  
414 inventors from a much wider range of countries and agencies. The search strategy is the same as  
415 the retrieval from USPTO: topic = *construction project, project management, infrastructure*  
416 *project, civil engineering* and *transportation project*. As the DIX does not provide fields of claim  
417 and description, title and abstract were collected as raw data. Before validation, the raw data of the  
418 200 patents have to be annotated, processed and vectorized by the same processes mentioned in  
419 *Section 3.1, 3.2 and 3.3*.

#### 420 **4.2.2. Validation results**

421 In this validation, the goal is to verify if the MLP has better screening accuracy than the traditional  
422 machine learning models. The performance of the MLP is compared against existing machine  
423 learning models, including Gaussian Naive Bayes (GNB), SVM and Bernoulli Naive Bayes (BNB).  
424 Figure 7 shows the accuracy over the different feature numbers. The highest precision value for  
425 each model are marked above the lines. By examining the figure, we can verify that the MLP  
426 model is superior to those machine learning models over all the features except  $K = 1000$  and  $K =$   
427  $40000$ . We can also observe that the MLP model is more sensitive to the number features, with the  
428 highest standard deviation value (0.032) over the models. This is consistent with one of the major  
429 differences between deep learning and traditional machine learning models: the traditional  
430 machine learning models are not capable of adjusting the model complexity according to the inputs,  
431 whereas the deep learning could tune the structure (number of layers and neurons) that is most  
432 suitable for input features (Moraes et al., 2013).



**Figure 7 The precision values for MLP and machine learning models over the features**

435  
 436  
 437 Table 3 illustrates the cross-validation results over the optimized MLP (K= 30000, U = 13), GNB  
 438 SVM, and BNB. As was mentioned above, 5-fold cross-validation is used to verify the  
 439 performance of the trained model. All the annotated instances were shuffled and randomly divided  
 440 into 5 folds with same size, and 4 of them were used for training and the rest were testing instances.  
 441 The training and testing would be performed in 5 times, and each time has a different fold for  
 442 testing and a different combination of 4 folds for training. The performance value is obtained by  
 443 the mean of the 5 testing results. It can be observed that MLP (K= 30000, U = 13) has the best  
 444 performance in all the three indexes (precision, recall and F1 score).

**Table 3 Cross-validation results over MLP, GNB, SVM and BNB in the initial dataset**

	<b>Precision</b>	<b>Recall</b>	<b>F1 score</b>
MLP (K=30000,U=13)	<b>0.955</b>	<b>0.954</b>	<b>0.954</b>
GNB (K=25000)	0.925	0.919	0.918
SVM (K=40000)	0.86	0.852	0.848
BNB (K=35000)	0.883	0.86	0.853

445  
 446  
 447 Table 4 shows the validation results over another database, which is used to evaluate the generality  
 448 of the proposed model. The validation results, described by Table 4, indicate that the learned

449 classifier based on MLP could be precisely implemented in the database from DIX, in which  
450 patents are written in different levels by a variety of inventors from different countries. In addition,  
451 the MLP also outperforms the machine learning methods.

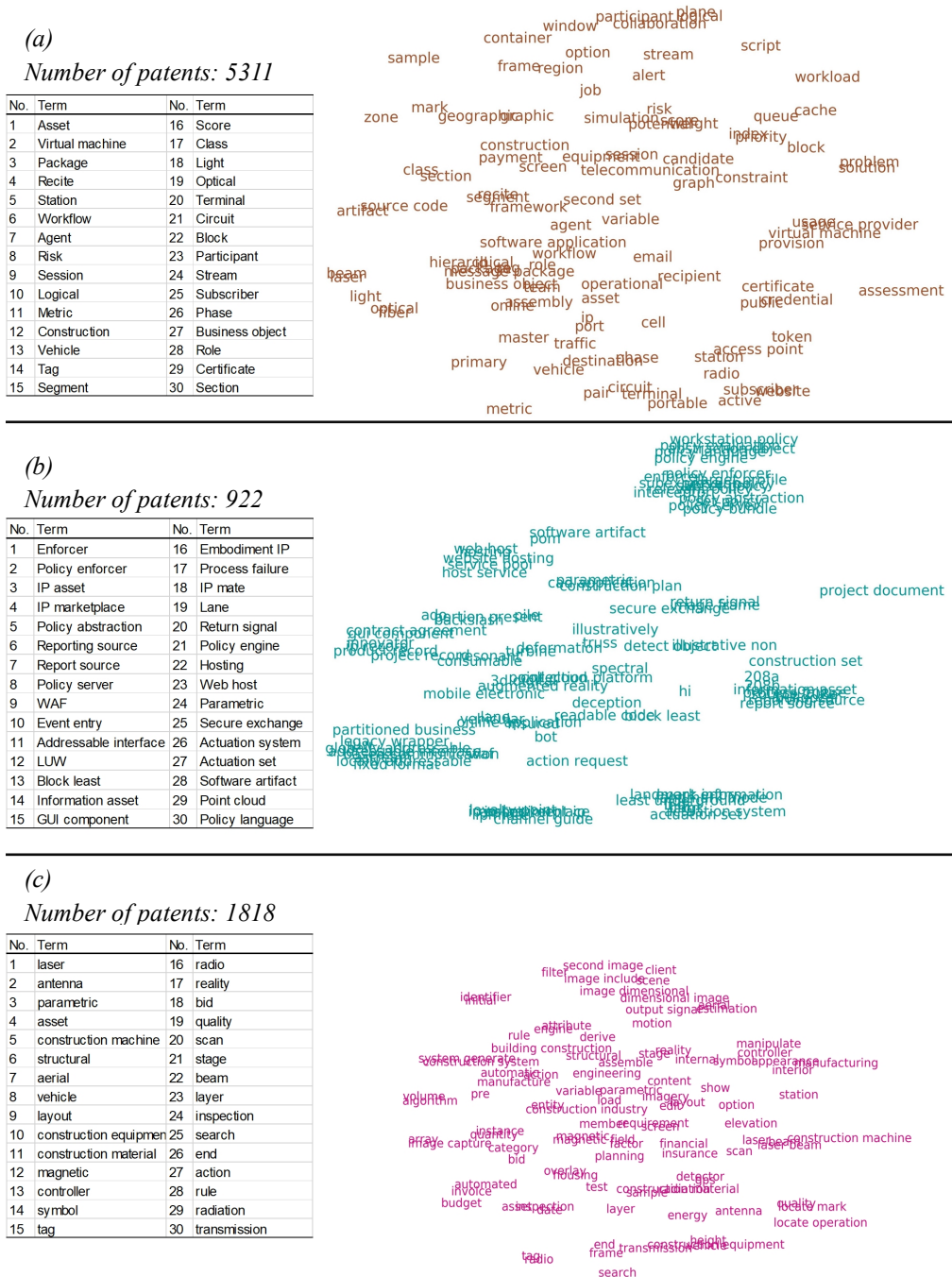
452 **Table 4 Cross-validation results over MLP, GNB, SVM and BNB in the dataset from DIX**

	<b>Precision</b>	<b>Recall</b>	<b>F1 score</b>
MLP (K=30000,U=13)	<b>0.897</b>	<b>0.897</b>	<b>0.897</b>
GNB (K=25000)	0.849	0.852	0.849
SVM (K=40000)	0.85	0.852	0.848
BNB (K=35000)	0.444	0.667	0.533

### 453 **4.3. The screened ICTC patents**

454 Besides the validation results, some important implications should be further discussed. The  
455 authors use the proposed approach to automatically screen a corpus of ICTC patents. To compare  
456 the topic distribution of the patents in the corpus, as well as the patents in collection 1 and 2 (Table  
457 1), this study plots the figures of feature space for each of the collections (Figure 8). According to  
458 the processes in section 3.1 - 3.3, this study vectorizes each of the patents in the three collections.  
459 The t-Distributed Stochastic Neighbor Embedding (TSNE) algorithm (Czerniawski et al., 2018;  
460 Maaten and Hinton, 2008) is adopted to project the high dimensional feature vectors into a 2D plot,  
461 in which the physical distance between two features roughly represents the degree of association  
462 of them in the corresponding collection.





Notes: (2) In each sub-figure, top 100 features with highest average Tf-Idf value are plotted, and top 30 are list at left for clarity. (2) As for strategy 1 and 2, please see Table 6.2 for details.

463  
464 **Figure 8 TSNE plots of feature spaces (a) The plot of patents screened by strategy 1; (b) The plot of**  
465 **patents screened by strategy 2; (c) The plot of patents screened by the proposed approach**  
466

467  
468 As explained in the introduction, the searching engines for patents has two major flaws: (1) the  
469 searching engines can only perform “match” logic based on structured data; and (2) searching by

470 keywords cannot avoid personal preference, and thus the results highly depend on users’  
471 knowledge. Figure 8 (a) depicts the feature space of collection 1, in which patents were searched  
472 by ICT classes and AEC domain keywords. The features in this figure are averagely distributed,  
473 incorporating a large number of ICT-related features, but some typical ICTC terms do not appear.  
474 Such feature distribution indicates that the patents in this collection are mainly relevant to ICT, but  
475 not ICTC. A possible explanation for this might be that the AEC domain keywords are not capable  
476 of discerning ICTC patents from ICT patents using query-based methods. For example, the  
477 keyword “construction project” may match patents related to construction projects, but it can also  
478 match patents of “software project” containing sayings about “construction project” which means  
479 construct a project. Despite the miss matching problem, strategy 2 can lead to short coverage of  
480 ICT techniques. As Figure 8 (b) shown, the features are agglomerated into clusters, indicating an  
481 unbalanced distribution of topics. The features, not surprised, are mainly related to the searching  
482 keywords, such as wireless and mobile. The features in Figure 8 (C) are distributed averagely,  
483 incorporating a wide range of ICTC related terminologies, such as laser, construction machine,  
484 and radio. This indicates that the proposed approach is more suitable in gathering ICTC patents  
485 than traditional searching engines.  
486

## 487 **5. Discussion and conclusion**

488 ICT applications are a key determinant to improve the level of coordination and collaboration in  
489 the AEC industry. Even though patents have been recognized as a valuable resource to provide  
490 technological knowledge, the patent offices have not provided a specific classification of ICTC  
491 patents. Acknowledging this research opportunity, the presenting study accurately and widely  
492 screens a corpus of ICTC patents, by proposing an approach based on deep learning and NLP

493 techniques.

494 Specifically, this study has made the following contributions: (1) This study contributes an  
495 approach to widely and accurately retrieve and collect a corpus of patents for domains like ICTC  
496 that does not exist as a specific classification in patents and hardly being represented by queries.  
497 Although patent offices provide elaborate classification schemes, it cannot satisfy all the  
498 requirements in the real world. Therefore, when a collection of patents does not exist in the  
499 classification schemes, query-based methods become the only possible way to search these patents.  
500 However, query-based methods were developed for retrieving relevant documents for a specific  
501 patent application rather than a set of patents. For the collections like ICTC that incorporates a  
502 variety of technologies, it is an extremely challenging task for the query-based methods to retrieve  
503 the patents simply by a query. (2) The proposed approach takes advantage of deep learning and  
504 NLP techniques. Although deep learning has become prominent in processing textual data, the  
505 previous studies in the AEC area mainly utilized machine learning methods to perform  
506 classification tasks, which performance highly depends on feature selection because the traditional  
507 machine learning models could only learn the linear relations. Compared to traditional machine  
508 learning methods, deep learning models are more advanced by using the layers of neural networks  
509 to learn non-linear relations and more suitable for complex tasks with specific objectives. The  
510 validation results indicate that the MLP model outperforms the traditional machine models in  
511 classifying the ICTC patents. In addition, NLP techniques were employed to pre-process the raw  
512 data. In the AEC area, most previous studies only utilized the N-gram, tokenization and stop-words  
513 removal, ignoring the advanced NLP tools such as lemmatization and POS. This study utilizes  
514 lemmatization and POS to convert the words into stems for generating more accurate N-grams  
515 from the textual data. (3) In practice, this study contributes a specific collection for ICTC patents,

516 which is not provided by the patent offices. The collection widely and accurately covers the ICT  
517 applications in the construction, not only constituting a dictionary for searching ICTC, but also  
518 identifying problems to be solved by the state of art ICTC inventions and recognizing all possible  
519 specific embodiments of ICTC.

520

521 The presenting study is not without limitations. The feature extraction process is based on  
522 traditional BOW models, which does not take the semantic meanings in contexts into consideration.  
523 This limitation, however, has been largely offset by the proposed supervised MLP model, which  
524 learns complex relations between the inputs and outputs by training the deep layers of neurons.  
525 This could perform the prediction task with good performance without considering the semantic  
526 meanings. This study focusses on classifying ICTC and non-ICTC. Future research is needed to  
527 concentrate on AI-aided approaches that could automatically categorize the ICTC patents  
528 according to the technological components or the management issues in practice.

## 529 **References**

530

- 531 Adriaanse, A., Voordijk, H., Dewulf, G., 2010. The use of interorganisational ICT in United States  
532 construction projects. *Automation in Construction*, 19(1), 73-83.
- 533 Aggarwal, C.C., Reddy, C.K., 2013. *Data clustering: algorithms and applications*. CRC press.
- 534 Agrawal, A., Henderson, R., 2002. Putting patents in context: Exploring knowledge transfer from  
535 MIT. *Management Science*, 48(1), 44-60.
- 536 Ahuja, V., Yang, J., Shankar, R., 2009. Study of ICT adoption for building project management in  
537 the Indian construction industry. *Automation in Construction*, 18(4), 415-423.
- 538 Ahuja, V., Yang, J., Skitmore, M., Shankar, R., 2010. An empirical test of causal relationships of  
539 factors affecting ICT adoption for building project management: An Indian SME case study.  
540 *Construction innovation*, 10(2), 164-180.
- 541 Al Rahhal, M.M., Bazi, Y., Al Zuair, M., Othman, E., BenJdira, B., 2018. Convolutional Neural  
542 Networks for Electrocardiogram Classification. *Journal of Medical and Biological  
543 Engineering*, 38(6), 1014-1025.
- 544 Alberts, D., Yang, C.B., Fobare-DePonio, D., Koubek, K., Robins, S., Rodgers, M., Simmons, E.,  
545 DeMarco, D., 2017. Introduction to patent searching, *Current Challenges in Patent*

546 Information Retrieval. Springer, pp. 3-45.

547 Alsafouri, S., Ayer, S.K., 2018. Review of ICT Implementations for Facilitating Information Flow  
548 between Virtual Models and Construction Project Sites. *Automation in Construction*,  
549 86(August 2016), 176-189.

550 Azad, H.K., Deepak, A., 2019. Query expansion techniques for information retrieval: A survey.  
551 *Information Processing & Management*, 56(5), 1698-1735.

552 Bell, G., Hey, T., Szalay, A., 2009. Beyond the data deluge. *Science*, 323(5919), 1297-1298.

553 Bengio, Y., Ducharme, R., Vincent, P., Jauvin, C., 2003. A neural probabilistic language model.  
554 *Journal of Machine Learning Research*, 3(Feb), 1137-1155.

555 Benson, C.L., Magee, C.L., 2013. A hybrid keyword and patent class methodology for selecting  
556 relevant sets of patents for a technological field. *Scientometrics*, 96(1), 69-82.

557 Bird, S., Loper, E., 2004. NLTK: the natural language toolkit, *Proceedings of the ACL 2004 on*  
558 *Interactive poster and demonstration sessions. Association for Computational Linguistics*,  
559 p. 31.

560 Blei, D.M., Ng, A.Y., Jordan, M.I., 2003. Latent dirichlet allocation. *Journal of Machine Learning*  
561 *Research*, 3, 993-1022.

562 Bouadjenek, M.R., Sanner, S., Ferraro, G., 2015. A study of query reformulation for patent prior  
563 art search with partial patent applications, *Proceedings of the 15th International Conference*  
564 *on Artificial Intelligence and Law. ACM*, pp. 23-32.

565 Brown, I., Mues, C., 2012. An experimental comparison of classification algorithms for  
566 imbalanced credit scoring data sets. *Expert Systems with Applications*, 39(3), 3446-3453.

567 Cakir, L., Yilmaz, N., 2014. Polynomials, radial basis functions and multilayer perceptron neural  
568 network methods in local geoid determination with GPS/levelling. *Measurement*, 57, 148-  
569 153.

570 Cambria, E., White, B., 2014. Jumping NLP curves: A review of natural language processing  
571 research. *IEEE Computational Intelligence Magazine*, 9(2), 48-57.

572 Cassetta, E., Marra, A., Pozzi, C., Antonelli, P., 2017. Emerging technological trajectories and new  
573 mobility solutions. A large-scale investigation on transport-related innovative start-ups and  
574 implications for policy. *Transportation Research Part A: Policy and Practice*, 106(March),  
575 1-11.

576 Chakrabarti, S., Dom, B., Indyk, P., 1998. Enhanced hypertext categorization using hyperlinks.  
577 *SIGMOD Rec.*, 27(2), 307-318.

578 Chiarello, F., Cimino, A., Fantoni, G., Dell'Orletta, F., 2018. Automatic users extraction from  
579 patents. *World Patent Information*, 54, 28-38.

580 Choi, S., Kang, D., Lim, J., Kim, K., 2012. A fact-oriented ontological approach to SAO-based  
581 function modeling of patents for implementing Function-based Technology Database.  
582 *Expert Systems with Applications*, 39(10), 9129-9140.

583 Cocarascu, O., Toni, F., 2018. Combining Deep Learning and Argumentative Reasoning for the  
584 Analysis of Social Media Textual Content Using Small Data Sets. *Computational*  
585 *Linguistics*, 44(4), 833-858.

586 Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., Kuksa, P., 2011. Natural  
587 language processing (almost) from scratch. *Journal of Machine Learning Research*,  
588 12(Aug), 2493-2537.

589 Czerniawski, T., Sankaran, B., Nahangi, M., Haas, C., Leite, F., 2018. 6D DBSCAN-based  
590 segmentation of building point clouds for planar object classification. *Automation in*  
591 *Construction*, 88, 44-58.

592 Davies, R., Harty, C., 2013. Implementing 'Site BIM': A case study of ICT innovation on a large  
593 hospital project. *Automation in Construction*, 30, 15-24.

594 Dehlin, S., Olofsson, T., 2008. An evaluation model for ICT investments in construction projects.  
595 *Electronic journal of information technology in construction*, 13, 343-361.

596 Duchi, J., Hazan, E., Singer, Y., 2011. Adaptive subgradient methods for online learning and  
597 stochastic optimization. *Journal of Machine Learning Research*, 12(Jul), 2121-2159.

598 El-Ghandour, W., Al-Hussein, M., 2004. Survey of information technology applications in  
599 construction. *Construction innovation*, 4(2), 83-98.

600 Enesi, F.A., Oyefolahan, I.O., Abdullahi, M.B., Salaudeen, M.T., 2018. Enhanced Query  
601 Expansion Algorithm: Framework for Effective Ontology Based Information Retrieval  
602 System. *i-Manager's Journal on Computer Science*, 6(4), 1.

603 Fall, C., Benzineb, K., Guyot, J., Törösvári, A., Fiévet, P., 2003a. Computer-assisted categorization  
604 of patent documents in the international patent classification, *Proceedings of the*  
605 *International Chemical Information Conference (ICIC'03)*, Nîmes, France.

606 Fall, C.J., T, A., #246, rcsv, #225, ri, Benzineb, K., Karetka, G., 2003b. Automated categorization  
607 in the international patent classification. *SIGIR Forum*, 37(1), 10-25.

608 Fan, X., Li, S., Tian, L., 2015. Chaotic characteristic identification for carbon price and an multi-  
609 layer perceptron network prediction model. *Expert Systems with Applications*, 42(8),  
610 3945-3952.

611 Flyvbjerg, B., 2014. What You Should Know About Megaprojects and Why: An Overview. *Project*  
612 *Management Journal*, 45(2), 6-19.

613 Forman, G., 2002. Choose Your Words Carefully: An Empirical Study of Feature Selection Metrics  
614 for Text Classification. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 150-162.

615 Friedman, J., Hastie, T., Tibshirani, R., 2001. The elements of statistical learning. Springer series  
616 in statistics New York, NY, USA:.

617 Frits, S., 2007. Strategy to enhance use of ICT in construction, *Proceedings of CIB World Building*  
618 *Congress*, pp. 2527-2535.

619 Garcia-Laencina, P.J., Sancho-Gomez, J.L., Figueiras-Vidal, A.R., 2013. Classifying patterns with  
620 missing values using Multi-Task Learning perceptrons. *Expert Systems with Applications*,  
621 40(4), 1333-1341.

622 García Adeva, J.J., Pikatza Atxa, J.M., Ubeda Carrillo, M., Ansuategi Zengotitabengoa, E., 2014.  
623 Automatic text classification to support systematic reviews in medicine. *Expert Systems*  
624 *with Applications*, 41(4 PART 1), 1498-1508.

625 Gerken, J.M., 2012. A new instrument for technology monitoring: novelty in patents measured by  
626 semantic patent analysis. *Scientometrics*, 91(3), 645-670.

627 Giachanou, A., Salamasis, M., Paltoglou, G., 2015. Multilayer source selection as a tool for  
628 supporting patent search and classification. *Information Retrieval Journal*, 18(6), 559-585.

629 Gimpel, K., Schneider, N., O'Connor, B., Das, D., Mills, D., Eisenstein, J., Heilman, M., Yogatama,  
630 D., Flanigan, J., Smith, N.A., 2011. Part-of-Speech Tagging for Twitter: Annotation,  
631 Features, and Experiments, Meeting of the Association for Computational Linguistics:  
632 Human Language Technologies: Short Papers, pp. 42-47.

633 Girthana, K., Swamynathan, S., 2018. Semantic Query-Based Patent Summarization System  
634 (SQPSS), *International Conference on Intelligent Information Technologies*. Springer, pp.  
635 169-179.

636 Greiman, V.A., 2013. Megaproject management: Lessons on risk and project management from  
637 the Big Dig. John Wiley & Sons.

- 638 Gwak, J.H., Sohn, S.Y., 2018. A novel approach to explore patent development paths for subfield  
639 technologies. *Journal of the Association for Information Science and Technology*, 69(3),  
640 410-419.
- 641 Habash, N., Rambow, O., Roth, R., 2009. MADA+ TOKAN: A toolkit for Arabic tokenization,  
642 diacritization, morphological disambiguation, POS tagging, stemming and lemmatization,  
643 Proceedings of the 2nd international conference on Arabic language resources and tools  
644 (MEDAR), Cairo, Egypt, p. 62.
- 645 Haddi, E., Liu, X., Shi, Y., 2013. The role of text pre-processing in sentiment analysis. *Procedia*  
646 *Computer Science*, 17, 26-32.
- 647 Haykin, S., 1999. *Neural networks a comprehensive introduction*. Prentice Hall, New Jersey.
- 648 Hoetker, G., Agarwal, R., 2007. Death Hurts, But It Isn't Fatal: The Postexit Diffusion of  
649 Knowledge Created by Innovative Companies. *Academy of Management Journal*, 50(2),  
650 446-467.
- 651 Intarakumnerd, P., Charoenporn, P., 2015. Impact of stronger patent regimes on technology transfer:  
652 The case study of Thai automotive industry. *Research Policy*, 44(7), 1314-1326.
- 653 Ittoo, A., Nguyen, L.M., van den Bosch, A., 2016. Text analytics in industry: Challenges,  
654 desiderata and trends. *Computers in Industry*, 78, 96-107.
- 655 Jain, A., Kulkarni, G., Shah, V., 2018. Natural Language Processing. *International Journal of*  
656 *Computer Sciences and Engineering*, 6(1), 161-167.
- 657 Kaplan, S., Vakili, K., 2013. Novelty vs. usefulness in innovative breakthroughs: A test using topic  
658 modeling of nanotechnology patents, Technical Report, Working Paper.
- 659 Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. arXiv preprint  
660 arXiv:1412.6980.
- 661 Kurdi, M.Z., 2017. *Natural language processing and computational linguistics 2: semantics,*  
662 *discourse and applications*. John Wiley & Sons.
- 663 Lee, C., Song, B., Park, Y., 2013. How to assess patent infringement risks: a semantic patent claim  
664 analysis using dependency relationships. *Technology Analysis and Strategic Management*,  
665 25(1), 23-38.
- 666 Levitt, R.E., 2007. CEM research for the next 50 years: Maximizing economic, environmental,  
667 and societal value of the built environment. *Journal of Construction Engineering and*  
668 *Management-Asce*, 133(9), 619-628.
- 669 Levy, O., Goldberg, Y., Dagan, I., 2015. Improving distributional similarity with lessons learned  
670 from word embeddings. *Transactions of the Association for Computational Linguistics*, 3,  
671 211-225.
- 672 Li, H., Chan, G., Wong, J.K.W., Skitmore, M., 2016. Real-time locating systems applications in  
673 construction. *Automation in Construction*, 63, 37-47.
- 674 Li, S., Hu, J., Cui, Y., Hu, J., 2018. DeepPatent: patent classification with convolutional neural  
675 networks and word embedding. *Scientometrics*, 117(2), 721-744.
- 676 Li, X., Shen, G.Q., Wu, P., Yue, T., 2019. Integrating Building Information Modeling and  
677 Prefabrication Housing Production. *Automation in Construction*, 100, 46-60.
- 678 Li, Z., Tate, D., Lane, C., Adams, C., 2012. A framework for automatic TRIZ level of invention  
679 estimation of patents using natural language processing, knowledge-transfer and patent  
680 citation metrics. *Computer-Aided Design*, 44(10), 987-1010.
- 681 Liu, S.-H., Liao, H.-L., Pi, S.-M., Hu, J.-W., 2011. Development of a Patent Retrieval and Analysis  
682 Platform – A hybrid approach. *Expert Systems with Applications*, 38(6), 7864-7868.
- 683 Maaten, L., Hinton, G., 2008. Visualizing data using t-SNE. *Journal of Machine Learning Research*,

684 9(2605), 2579-2605.

685 Mahdabi, P., Crestani, F., 2014. The effect of citation analysis on query expansion for patent  
686 retrieval. *Information Retrieval*, 17(5-6), 412-429.

687 Mahdabi, P., Keikha, M., Gerani, S., Landoni, M., Crestani, F., 2011. Building Queries for Prior-  
688 Art Search, *Information Retrieval Facility Conference*.

689 Michel, J., Bettels, B., 2001. Patent citation analysis. A closer look at the basic input data from  
690 patent search reports. *Scientometrics*, 51(1), 185-201.

691 Mirończuk, M.M., Protasiewicz, J., 2018. A recent overview of the state-of-the-art elements of text  
692 classification. *Expert Systems with Applications*, 106, 36-54.

693 Mok, K.Y., Shen, G.Q., Yang, J., 2015. Stakeholder management studies in mega construction  
694 projects: A review and future directions. *International Journal of Project Management*,  
695 33(2), 446-457.

696 Moraes, R., Valiati, J.F., Neto, W.P.G., 2013. Document-level sentiment classification: An  
697 empirical comparison between SVM and ANN. *Expert Systems with Applications*, 40(2),  
698 621-633.

699 Munková, D., Munk, M., Vozár, M., 2013. Data pre-processing evaluation for text mining:  
700 transaction/sequence model. *Procedia Computer Science*, 18, 1198-1207.

701 Nadeau, D., Sekine, S., 2007. A survey of named entity recognition and classification. *Linguisticae*  
702 *Investigationes*, 30(1), 3-26.

703 Niemann, H., Moehrle, M.G., Frischkorn, J., 2017. Use of a new patent text-mining and  
704 visualization method for identifying patenting patterns over time: Concept, method and test  
705 application. *Technological Forecasting and Social Change*, 115, 210-220.

706 Onan, A., Korukoğlu, S., Bulut, H., 2016. Ensemble of keyword extraction methods and classifiers  
707 in text classification. *Expert Systems with Applications*, 57, 232-247.

708 Pavlinek, M., Podgorelec, V., 2017. Text classification method based on self-training and LDA  
709 topic models. *Expert Systems with Applications*, 80, 83-93.

710 Perez-Molina, E., 2018. The Role of Patent Citations as a Footprint of Technology. *Journal of the*  
711 *Association for Information Science and Technology*, 69(4), 610-618.

712 Riedmiller, M., 1994. Advanced supervised learning in multi-layer perceptrons-from  
713 backpropagation to adaptive learning algorithms. *Computer standards and interfaces*, 16(3),  
714 265-278.

715 Robbins, H., Monro, S., 1985. A stochastic approximation method, *Herbert Robbins Selected*  
716 *Papers*. Springer, pp. 102-109.

717 Rosenblatt, F., 1961. Principles of neurodynamics. perceptrons and the theory of brain mechanisms.  
718 CORNELL AERONAUTICAL LAB INC BUFFALO NY.

719 Saiki, T., Akano, Y., Watanabe, C., Tou, Y., 2006. A new dimension of potential resources in  
720 innovation: A wider scope of patent claims can lead to new functionality development.  
721 *Technovation*, 26(7), 796-806.

722 Sardroud, J.M., 2015. Perceptions of automated data collection technology use in the construction  
723 industry. *Journal of Civil Engineering and Management*, 21(1), 54-66.

724 Shalaby, W., Zadrozny, W., 2018. Toward an interactive patent retrieval framework based on  
725 distributed representations, *The 41st International ACM SIGIR Conference on Research &*  
726 *Development in Information Retrieval*. ACM, pp. 957-960.

727 Shalaby, W., Zadrozny, W., 2019. Patent retrieval: a literature review. *Knowledge and Information*  
728 *systems*, 1-30.

729 Shannon, C.E., 1948. A mathematical theory of communication. *Bell system technical journal*,



730 27(3), 379-423.

731 Shekarpour, S., Marx, E., Ngonga Ngomo, A.C., Auer, S., 2015. SINA: Semantic interpretation of  
732 user queries for question answering on interlinked data. *Journal of Web Semantics*, 30, 39-  
733 51.

734 Silva, F.N., Amancio, D.R., Bardosova, M., Costa, L.D., Oliveira, O.N., 2016. Using network  
735 science and text analytics to produce surveys in a scientific topic. *Journal of Informetrics*,  
736 10(2), 487-502.

737 Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser,  
738 J., Antonoglou, I., Panneershelvam, V., Lanctot, M., 2016. Mastering the game of Go with  
739 deep neural networks and tree search. *Nature*, 529(7587), 484.

740 Smith, H., 2002. Automation of patent classification. *World Patent Information*, 24(4), 269-271.

741 Sokolova, M., Lapalme, G., 2009. A systematic analysis of performance measures for  
742 classification tasks. *Information Processing & Management*, 45(4), 427-437.

743 Sparck Jones, K., 1972. A statistical interpretation of term specificity and its application in retrieval.  
744 *Journal of documentation*, 28(1), 11-21.

745 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: a  
746 simple way to prevent neural networks from overfitting. *The Journal of Machine Learning*  
747 *Research*, 15(1), 1929-1958.

748 Tannebaum, W., Rauber, A., 2014. Using query logs of USPTO patent examiners for automatic  
749 query expansion in patent searching. *Information Retrieval*, 17(5-6), 452-470.

750 Terragno, P.J., 1979. Patents as technical literature. *IEEE Transactions on Professional*  
751 *Communication*, PC-22(2), 101-104.

752 Tieleman, T., Hinton, G., 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of  
753 its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2), 26-31.

754 USPTO, 2007. *Manual of patent examining procedure* (8th ed.), Alexandria.

755 Venugopalan, S., Rai, V., 2015. Topic based classification and pattern identification in patents.  
756 *Technological Forecasting and Social Change*, 94, 236-250.

757 Vlachidis, A., Tudhope, D., 2016. A knowledge - based approach to I nformation E xtraction for  
758 semantic interoperability in the archaeology domain. *Journal of the Association for*  
759 *Information Science and Technology*, 67(5), 1138-1152.

760 Wang, J., 2018. Innovation and government intervention: A comparison of Singapore and Hong  
761 Kong. *Research Policy*, 47(2), 399-412.

762 Wang, X., Jiang, W., Luo, Z., 2016. Combination of convolutional and recurrent neural network  
763 for sentiment analysis of short texts, *Proceedings of COLING 2016, the 26th international*  
764 *conference on computational linguistics: Technical papers*, pp. 2428-2437.

765 Wu, C.H., Yun, K., Huang, T., 2010. Patent classification system using a new hybrid genetic  
766 algorithm support vector machine. *Applied Soft Computing Journal*, 10(4), 1164-1177.

767 Wu, P., Song, Y., Shou, W., Chi, H., Chong, H.-Y., Sutrisna, M., 2017. A comprehensive analysis  
768 of the credits obtained by LEED 2009 certified green buildings. *Renewable and Sustainable*  
769 *Energy Reviews*, 68, 370-379.

770 Zhang, L., Liu, Z., Li, L., Shen, C., Li, T., 2018. PatSearch: an integrated framework for  
771 patentability retrieval. *Knowledge and Information systems*, 57(1), 135-158.

772 Zhao, Z., Xu, S., Kang, B.H., Kabir, M.M.J., Liu, Y., Wasinger, R., 2015. Investigation and  
773 improvement of multi-layer perceptron neural networks for credit scoring. *Expert Systems*  
774 *with Applications*, 42(7), 3508-3516.

775 Zidane, Y.J.T., Johansen, A., Ekambaram, A., 2013. *Megaprojects - Challenges and Lessons*

776  
777

Learned, in: Pantouvakis, J.P. (Ed.), Selected Papers from the 26th Ipma, pp. 349-357.