

METHODS FOR PATENT LANDSCAPE MODELLING AND PREDICTIVE ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

The article is dedicated to the problem of forecasting changes in the patent landscape based on the analysis of multimodal data. This article examines three main approaches to the analysis of patent information: 1) a clustering-based approach; 2) a resource-based approach; 3) a machine learning approach. Each approach is considered in terms of its advantages, flaws and potential for use in predicting the emergence of new technologies. The article proposes new methods for constructing a patent landscape and predictive models, as well as a method for assessing development directions in technological areas. The article also considers the problem of visualizing the results of analysis and forecasting of technological trends in order to provide a clear visual representation of the patent landscape. The article presents the results of experiments demonstrating the ability of the proposed methods to make correct predictions of technological trends.

Keywords—Clustering, machine learning, patent analysis, trend forecasting.

1 INTRODUCTION

Modern big data processing systems analyze multimodal data considering it as a combination of text, images, audio and video taken from different sources. This kind of analysis requires methods that effectively process data of heterogeneous formats and extract mutually complementary information from it.

Neural networks have proven their high effectiveness in complex and heterogeneous data processing tasks, solving them due to the ability to learn from large volumes of data and to identify the hidden patterns. Neural networks play a special role in the analysis of images (Anwar, 2018), texts (Tarwani & Edem, 2017) and graphs (Wu et al., 2020). There are many examples in the literature, and they include the following. The image analysis identifies diseases of human organs based on tomographic data (Anwar, 2018). The text data analysis identifies key topics and sentiments, as well as extracting data from unstructured texts (Tarwani & Edem, 2017). The graph data analysis identifies intricate relationships and dependencies, such as influential nodes, it detects communities, or finds the shortest paths to address real-world challenges in such areas as social network analysis, recommendation systems, fraud detection and so on (Wu et al., 2020).

The use of deep neural network architectures such as transformers increases the accuracy and depth of analysis considering the context and interrelations of the different elements (Gillioz, 2020). The use of neural networks for multimodal data analysis is of great importance in the development of intelligent systems for big data processing and decision support. This allows the opportunity to offer more informed solutions. The development of methods and models for multimodal data analysis leads to the creation of large language models (LLM) (Wang et al., 2024) and large multimodal models (LMM) (Li et al., 2024), which are used successfully when working with large sets of multimodal data.

In the modern world where innovation is so important, there is a significant increase in patent information associated with the intensive development of technologies and global digitalization. Thousands of new patents are published annually containing valuable data on technical innovations and industry trends. The patents themselves are an important source of information for companies, research organizations, and government agencies, because they help to track technological developments and to predict future directions for scientific research, industrial development and commercial

054 opportunity. However, a huge amount of this patent data is often analyzed and processed manually.
055 This is a complex task that requires significant time and resource, and hence cost.

056 To solve these problems, and to create an effective tool, it is necessary to create automated systems
057 for patent information analysis and processing. Automated systems that can process and analyze
058 large amounts of patent data efficiently have the potential to reduce data processing time signifi-
059 cantly and, in particular, improve the accuracy of forecasts regarding the emergence of new tech-
060 nological areas. The importance of developing such software systems is explained by the need for
061 strategic planning and decision-making tasks that are based on reliable and timely information about
062 technological trends (Berezkin et al., 2024). There are many approaches to patent data analysis and
063 processing, which include clustering, networks research and machine learning, some of these ap-
064 proaches are described in Section II of this article.

065 In recent years the Natural Language Processing (NLP) and Deep Learning (DL) technologies have
066 made significant progress, opening new opportunities for the creation of advanced and efficient
067 patent analysis systems. One of the main tasks in this context is to predict the emergence of new
068 technological areas, which requires the use of complex clustering and data classification algorithms.

069 This article describes a method for Predicting the Emergence of New Technological Areas based
070 on multimodal patent information (PENTA). The PENTA method allows loading and processing
071 data which is extracted from patent documents. PENTA converts the text information into a vector
072 representation of the clusters of patents and visualizes the patents' landscape. The main emphasis
073 of PENTA is on optimizing the clustering threshold and associating the obtained clusters with Inter-
074 national Patent Classification (IPC) classes (Hoshino et al., 2023), thus improving the accuracy of
075 forecasts and forming well-founded hypotheses about the emergence of new technologies.

077 2 RELATED WORKS

078 This section examines three main approaches to the patent information analysis from the point of
079 view of the ability to predict the emergence of new / promising technological areas.

083 2.1 CLUSTERING-BASED APPROACH TO PATENT INFORMATION ANALYSIS

084 The clustering-based approach to patent information analysis aims to identify and predict new tech-
085 nological areas by grouping patents based on similar characteristics.

086 The first step of such analysis is to transform the textual patent information into numerical vec-
087 tors. Traditional vectorization methods such as Term Frequency-Inverse Document Frequency (TF-
088 IDF) (Qaiser & Ali, 2018) consider the frequency of word occurrence in documents only. The
089 word embeddings models such as Word2Vec (Mikolov et al., 2013), GloVe (Ji et al., 2021), and
090 FastText (Joulin et al., 2016) transform words into vectors (embeddings) reflecting their semantic
091 meaning and context. The modern transformer models such as BERT (Devlin et al., 2019) or GPT-
092 3 (Brown et al., 2020) can create deep contextual representations of text, which can account for
093 more complex relationships between words and phrases. These models can be extended to the sen-
094 tence and document level using Sentence-BERT (Reimers & Gurevych, 2019) or Universal Sentence
095 Encoder (Sarkar et al., 2022).

096 The next step of such analysis involves the application of clustering algorithms for grouping patents:
097 the algorithm K-means (MacQueen, 1967) divides data into K clusters minimizing the intra-cluster
098 distance; hierarchical clustering (Nielsen, 2016) constructs a dendrogram, where clusters are formed
099 by sequentially combining or dividing groups of data; the algorithm Density-Based Spatial Cluster-
100 ing of Applications with Noise (DBSCAN) (Ester et al., 1996) selects clusters based on the density
101 of points in the vector space and allows an automatic determination of the number of clusters as well
102 as separates off the noise data; the algorithm Affinity Propagation (Manoj et al., 2015) determines
103 cluster centers by passing messages between data points, which allows for efficient grouping of large
104 volumes of data.

105 An important aspect of the approach is the analysis of changes in the structure of clusters over
106 time. The growth or shrinkage of clusters corresponds to increasing or decreasing interest in certain
107 technological areas. The merger or separation of clusters corresponds to the combination of techno-

108 logical areas or emergence of new, specialized areas. The emergence of new clusters indicates the
109 birth of new technological areas.

110 To help researchers and analysts to understand better and to interpret the complex relationships
111 between patents and technology areas, various visualization tools are used: dendrograms which
112 allow the researcher to see the structure and the connections between the clusters; heat maps which
113 show a proximity between patents and clusters; network graphs which display the connections and
114 interactions between the patents and the clusters.

115 It is important to mention that clustering-based methods can be used in combination in order to
116 achieve the best results in predicting the new technology trends and identifying gaps in patent data.
117

118 2.2 RESOURCE-BASED APPROACH TO PATENT INFORMATION ANALYSIS

119 A resource network is a directed graph in which nodes can store resources (for example, patents),
120 and edges have a certain capacity that limits the amount of resource transferred between nodes. At
121 each discrete time step, the resource is redistributed between nodes observing the conservation law.
122

123 The main methods of the resource-based approach are: hidden Markov models (HMM) in com-
124 bination with Bayesian networks (Jurafsky & Martin, 2023; Lee et al., 2017b); Markov processes
125 for systems analysis (Aristodemou & Tietze, 2018); random walks and Markov chains (Vassiliou &
126 Georgiou, 2021).
127

128 Each of these methods offers distinct advantages and limitations. Hidden Markov models and
129 Bayesian networks excel at modeling complex dependencies and are effective for forecasting, but
130 their primary drawbacks are the need for significant amounts of data and the complexity of their
131 setup and interpretation. Markov processes are valued for their simplicity and power in modeling
132 systems with discrete states; however, they are limited when dealing with more complex history
133 dependencies and require fine-tuning of transition probabilities. Finally, random walks and Markov
134 chains are highly adaptable for modeling network structures and accounting for temporal dynamics,
135 though they can be computationally expensive for large networks and require careful parameter cal-
136 ibration. These methods facilitate a deep understanding of the dynamics of technology development
137 and can serve as a basis for strategic planning in innovation activities and research.

138 2.3 MACHINE LEARNING APPROACH TO PATENT INFORMATION ANALYSIS

139 The general workflow for this approach has the following steps:

- 140 1. collecting the patent data, including description text, filing date, citations, and IPC classes;
- 141 2. cleaning data, processing text, vectorizing text information;
- 142 3. using historical data to train a classification model that predicts the emergence of new patent
143 classes;
- 144 4. evaluating the accuracy of the model on test data, adjusting hyperparameters to improve
145 the accuracy.
146
147
148

149 The main machine learning methods used for forecasting in the field of technological innovation are
150 the following.

151 Gradient Boosting Decision Trees (GBDT) (Friedman, 2002). This method can be applied to patent
152 data in order to identify trends and dependencies that indicate the possible emergence of new tech-
153 nological trends. For example, analyzing patent texts, their classifications, and timestamps can help
154 in predicting future technological innovations (Li et al., 2020).
155

156 Random Forests (RF) (Breiman, 2001). This method can be applied to patent data in order to identify
157 important features that influence the emergence of new technological areas. This information can
158 then be used for accurate prediction (Santiago et al., 2021).

159 Multilayer Perceptron (MLP) (Alzubaidi et al., 2021; Popescu et al., 2009). The study (Lee et al.,
160 2017a) used MLP to analyze patent data, which made it possible to identify new technologies at an
161 early stage. The technique includes the use of various patent indicators and association rules for
accurate and timely detection of developing technologies.

162 Convolutional Neural Networks (CNN) (Purwono et al., 2022). The article (Ji et al., 2024) discusses
163 a technique for predicting the emergence of new technologies using CNN and MLP. The study uses
164 18 input and 3 output indicators from the US Patent Office database, making it possible to identify
165 complex nonlinear relationships and analyze technology development trends at an early stage, which
166 helps in strategic planning and forecasting technological trends.

167 Recurrent Neural Networks (RNN) (Das et al., 2023). The article (Ji et al., 2024) uses RNN and
168 MLP to analyze patent data, which helps to accurately and promptly identify new technologies at
169 early stages of their development.

170 Long Short-Term Memory (LSTM) (Staudemeyer & Morris, 2019). The article (Hou et al., 2024)
171 uses LSTM to analyze the life cycle of technologies. The study proposes a new method for predicting
172 the emergence of new technologies by analyzing patent data. The use of LSTM allows for more
173 accurate and effective identification of technology development trends at different stages of their
174 life cycle, from emergence to growth, maturity, and saturation.

175 Deep neural networks with Transformer architecture (Transformer) (Vaswani et al., 2023). The
176 article (Zhao et al., 2024) focuses on identifying trends and selecting collaboration partners, taking in
177 to account the interdependence of knowledge and collaboration. The study uses transformer models
178 to analyze collaboration and knowledge networks, which allows organizations to strategically select
179 partners for innovation projects and predict future technology trends. The use of transformers helps
180 to analyze large amounts of data more accurately and improve the decision-making process, in turn
181 helping to increase the efficiency of innovation.

182 Association Rule Mining (ARM) (Jun et al., 2012). The study (Sunghae, 2013) uses ARM and
183 Box-Jenkins modeling methods to identify patterns in patent data, which allows the prediction of
184 technology trends and the identification of new promising technologies. The study focuses on ana-
185 lyzing data from the European Patent Office using cluster analysis, text mining and social network
186 methods, which allows for the formation of informative technology maps and improves strategic
187 innovation planning.

188 Bayesian Networks (Peal, 1985). The study (Jeong et al., 2021) describes a methodology for creating
189 technology roadmaps that not only adapt to changing circumstances, but also predict the emergence
190 of new technologies. The authors of (Jeong et al., 2021) use Bayesian networks and topic modeling
191 analysis to identify and assess risks, which allows for dynamically adjusting plans and predicting
192 new technology trends. This machine learning approach has been successfully applied in the field of
193 artificial intelligence, demonstrating its effectiveness and sustainability in technology management.

194 Each of the discussed methods has its own advantages and disadvantages, which makes them suit-
195 able for solving different types of problems: GBDT and Random Forest provide high accuracy and
196 robustness to noise, but require significant computational resources and careful tuning; MLP and
197 CNN are effective for analyzing data with nonlinear dependencies and features, but also require
198 large amounts of data and computational costs; RNN and LSTM are especially useful for sequential
199 data and time series, although they have problems with vanishing gradients; deep neural networks
200 with Transformer architecture are a modern approach that provides high efficiency in natural lan-
201 guage processing tasks, but also requires significant resources.

202 203 204 205 206 207 3 DATA USED FOR PATENT ANALYSIS

208
209
210
211 Patent documents are used as the input data for the development of the PENTA method. The docu-
212 ments were collected from the patent database of the United States Patent and Trademark Of-
213 fice (United States Patent and Trademark Office). When preparing the input data, the main goal was
214 to collect a representative data set that would cover various technological areas and time periods.
215 The data set includes patents published between 1995 and 2022. The patent information used for
the research includes patent title, text description, filing date and IPC classes.

4 THE PROPOSED METHODS FOR PATENT LANDSCAPE MODELING

4.1 METHOD FOR CONSTRUCTING A PATENT LANDSCAPE MODEL

The method for constructing a patent landscape model includes several stages. These ensure the transformation of text data into vector representations, clustering, and parameter optimization. This method allows the creation of highly efficient models that can identify hidden patterns in patent data and predict the emergence of new technological areas.

First, patent texts (Title and Abstract) are transformed into vector representations using a language model. These vector representations are stored in the database and can be reused at a later time without the need for recalculation. There are language models that have been specifically adapted for analyzing patent documents. Some of these models are universal, trained on patent data from different domains. Examples: PatentBERT (Lee & Hsiang, 2019) and Patent Scientific BERT augmented (PatentSBERTa) (Bekamiri et al., 2024).

The PatentSBERTa and PatentBERT models demonstrate similar quality indicators (such as Precision, Recall and F1) when solving classification problems (Bekamiri et al., 2024). The preference for one model over another is determined by the characteristics of a particular problem being solved. The authors chose the PatentBERT model as it is best suited for analyzing the title, text description and IPC class of patents. This model is an adapted version of the BERT model for processing patent texts.

After the vector representations have been made for each patent, the clustering process is launched using the DBSCAN algorithm available at the scikit-learn library. The DBSCAN is a density algorithm well suited for identifying clusters of arbitrary shape and it is robust to noise in the data. It identifies groups of patents that have similar thematic and technological characteristics.

The quality of clustering is assessed using the Adjusted Rand Index (ARI) metric. The ARI measures the degree of correspondence between the clusters and the well-known IPC class labels. The ARI considers the random coincidence and provides a more accurate estimation in comparison to other metrics. During the assessment of the clustering quality, the clusters obtained at the first stage are analyzed and the correspondence between them and the IPC classes can be performed.

The DBSCAN algorithm can be adjusted to improve the quality of clustering and maximize the ARI metric. For this the gradient descent method allows the researcher to find the optimal value of the parameter “neighborhood size of a point” (that is, the maximum distance between two samples for one to be considered as in the neighborhood of the other) that maximizes the ARI metric. Optimization is carried out iteratively, which allows achieving the best division of patents into clusters.

This process includes a pairwise comparison of clustering accuracy and IPC labels, which allows determining how well the clusters match the classification. With high correspondence, cluster labels can be associated with IPC class labels, which simplifies further analysis and interpretation of the data.

After all stages are completed, the model is saved in the database. The clustering threshold, ARI metric, cluster labels and associated IPC class labels are included. This allows the model to be reused for analyzing new data and assessing its quality.

4.2 METHOD FOR CONSTRUCTING A PREDICTIVE MODEL

The method for building a forecast model is based on similar principles to the patent landscape model building method, but it includes additional steps accounting the time aspects in order to generate forecasts. This method allows the researcher to identify possible directions of technology development and potential new technology areas based on the analysis of changes in patent data over specific time periods.

In order to construct a predictive model, one needs to create a new dataset similar to the one used to build the base model, but covering a more recent time range. This new dataset is filtered using the same criteria as the original one to ensure comparability of the results. It is important that the IPC codes and other filtering parameters remain identical for both datasets.

270 Similar to the base model building method, the patent text data from the new dataset is transformed
271 into vector representations using the PatentBERT model. The resulting vectors are stored in the
272 database for subsequent use in the clustering and analysis process.

273 The DBSCAN method is used to cluster the data in the new dataset with the optimized parameter
274 “neighborhood size of a point” found during the base model building step. This ensures consistency
275 and allows comparison of clustering results across time periods. The DBSCAN method identifies
276 groups of patents that have similar thematic and technological characteristics, which helps identify
277 changes in the patent landscape.

278 After clustering, the quality of clustering is assessed using the Adjusted Rand Index (ARI) metric,
279 similar to the baseline model. Next, the cluster labels are associated with the IPC class labels present
280 in the filter of the new dataset. This helps determine how well the clusters reflect the classification
281 of patents by IPC classes and identifies changes in thematic areas.

282 Based on the clustering results and the label association, a predictive patent landscape is created.
283 To visualize changes in patent data, the Uniform Manifold Approximation and Projection (UMAP)
284 is used, which allows building a two- or three-dimensional representation of clusters and their rela-
285 tionships. This landscape clearly shows how the technological areas have changed and where new
286 clusters may appear.

287 After completing all stages, the predictive model is saved in the database. The clustering threshold,
288 ARI metric, cluster labels and associated IPC class labels are included. This allows the model to be
289 reused to analyze new data and assess its quality.

292 4.3 METHOD FOR ASSESSING THE DEVELOPMENT DIRECTIONS OF TECHNOLOGICAL 293 AREAS

294
295 The method for assessing development directions is based on the analysis of changes in patent data
296 and the identification of significant patterns using an association matrix. This process allows for the
297 formation of hypotheses about possible directions of technological development and the assessment
298 of the probability of their occurrence.

299 At the initial stage an association matrix is created, in which the X axis represents the clusters of
300 the base model, and the Y axis represents the clusters of the forecast model. Each cell of the matrix
301 contains the percentage of correspondence between the clusters. This indicator is calculated as the
302 ratio of the number of patents present in both clusters to the total number of patents in the base
303 cluster. This allows for the assessment of changes in the patent landscape and the identification of
304 significant patterns.

305 The next step is the initialization of the list of hypotheses, which will be considered, based on the
306 analysis of the association matrix. This list will contain possible development scenarios.

307 For each row in the association matrix, a model class (X-cluster) is determined. Then, for this row,
308 cells with a match above a specified threshold are found. The number of such cells determines the
309 further course of the analysis:

- 311 1. if there is more than one cell, a check is made to see if they belong to different classes
312 (Y-clusters); if so, a hypothesis is created about the merging of classes and the probability
313 of the emergence of a new technological area; otherwise, a hypothesis is created about
314 diversification and a change in focus within the class;
- 315 2. if there is one cell, a hypothesis is created about maintaining the focus of the technological
316 area;
- 317 3. if there are no cells and the overall compliance is less than the threshold, a hypothesis is
318 created about a decrease in activity in this technological area.

319 After the hypotheses are formed, their probability is assessed. This assessment is based on an
320 analysis of the degree of change in the association matrix, the number of patents and other indicators.
321 Each type of hypothesis uses its own formula for calculating the probability: the hypothesis of class
322 merging – the probability P_s is calculated based on the sum of the percentage correspondence of all
323

324 cells participating in the merger and the average growth rate of patents in these clusters:
 325

$$326 \quad P_s = \left(\frac{\sum_{i=1}^n C_i}{n} \right) \times G \quad (1)$$

327
 328 where C_i is the percentage match for each cell, n is the number of cells, G is the average growth
 329 rate of patents; the hypothesis of diversification and change of focus – the probability P_d is based on
 330 the percentage of patents that have moved to a new cluster and the change in the topic of keywords
 331 identified using LDA:

$$332 \quad P_d = \frac{P_m}{P_t} \times K \quad (2)$$

333 where P_m is the number of patents that have moved to the new cluster, P_t is the total number of
 334 patents in the original cluster, K is the coefficient of change of keywords; the hypothesis of a decrease
 335 in activity – the probability P_a is estimated based on the general decrease in activity in the cluster,
 336 calculated as the difference in the number of patents between the base and forecast periods:
 337

$$338 \quad P_a = 1 - \frac{P_p}{P_b} \quad (3)$$

339 where P_p is the number of patents in the forecast period, P_b is the number of patents in the base
 340 period.
 341

342 At the final stage, the hypotheses and their probabilities are stored in the database. This allows
 343 decision makers to access the analyzed data and make informed decisions based on the forecasts
 344 received.

345 The method for assessing the development directions of technological innovations is the main ele-
 346 ment of the proposed approach, ensuring the identification of significant patterns and the formation
 347 of hypotheses about possible directions of technology development. The use of this method allows
 348 for the effective analysis of large volumes of data and provides users with valuable insights for
 349 making strategic decisions in the field of innovation and technology.
 350

351 5 EXPERIMENTS

352 5.1 EXPERIMENTAL VERIFICATION OF THE PROPOSED METHODS

353 The patent data dated 2000 and 2010 has been selected for testing. The aim of the experiment was
 354 to compare the patent landscapes of 2000 and 2010 and to forecast technology trends based on this
 355 comparison.
 356

357 To simplify the testing procedure, it was decided to narrow the analysis to several subclasses. The
 358 selected subclasses are large enough to provide sufficient data for analysis. They are also diverse
 359 enough to demonstrate the efficiency of the method in different areas. The authors decided to select
 360 the following subclasses: G05B, H01L, H04L and A61B. The selection process was performed
 361 before transferring the data to the classifier. After training the classifier on patent data for 2000, its
 362 accuracy was tested using cross-validation as shown in the figure. (Fig. 1). The cross-validation
 363 results show that the classifier was well trained.
 364

365 Using the methods described in Section 4, patent landscape models were constructed for the selected
 366 subclasses (a base model based on 2000 data and a predictive model based on 2010 data). The visual
 367 representation of the models created using UMAP, is shown in Fig. 2 and Fig. 3 respectively.
 368

369 Comparison of Fig. 2 and Fig. 3 suggests shifts in emphasis in patent documents (movement of
 370 cluster centroids and changes in keywords). Patent subclasses A61B and H04L have less in common,
 371 and patent class G05B has shifted to the side, which may also indicate the development of this cluster
 372 into an independent area. The proposed visualisation facilitated monitoring of centroid movements
 373 and changes in cluster boundaries, thereby allowing a comprehensive assessment of the dynamics
 374 of the patent landscape.

375 5.2 EVALUATING THE PERFORMANCE OF THE PROPOSED METHODS

376 In order to test the proposed approach for predicting the emergence of new technologies based
 377 on patent information, a series of 15 experiments were conducted on different datasets. Expert

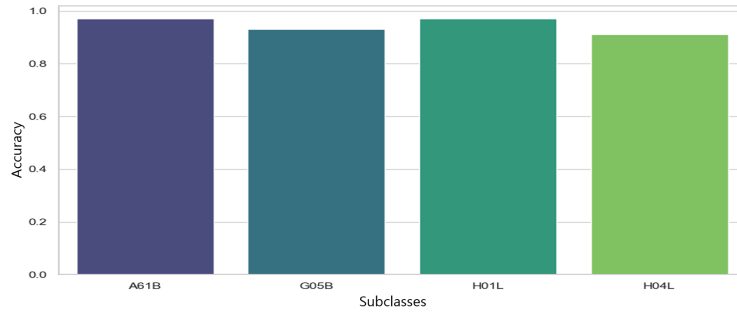


Figure 1: Classifier cross-validation results

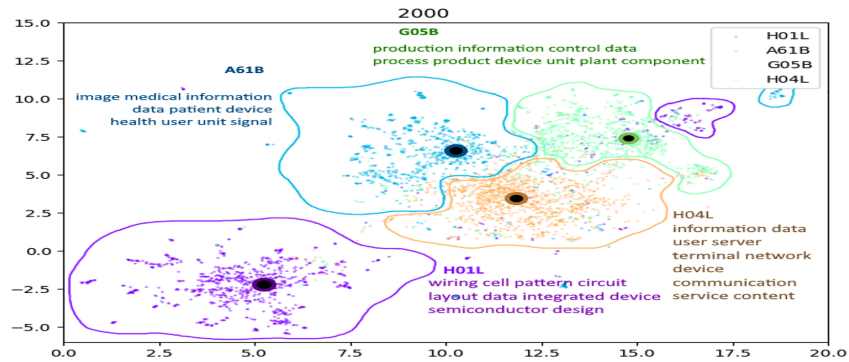


Figure 2: Distribution of patents in vector space for 2000

assessments were used to verify the accuracy and relevance of the predictions made via the proposed approach. In particular, one of the experiments analyzed the behavior of subclasses H01L, A61B, G05B, H04L in the period from 2010 to 2020. Comparison of data between 2010 and 2020 showed that the content of patents in G05B (regulatory and control systems) has become significantly closer to patents in the areas of medical solutions in diagnostics and surgery (A61B) and infrastructure systems for transmitting digital information (H04L). This indicates an increase in the number of solutions at the intersection of these areas. This trend is expected to intensify in the coming years due to the synergetic effect of the introduction of artificial intelligence and neuromorphic computing in the medical domain. Thus, it is likely that patent classes G05B and H04L will continue to converge, and at their intersection a new class dedicated to intelligent medical solutions may emerge.

To evaluate the accuracy of the proposed methods, an accuracy indicator was used, calculated using the following formula:

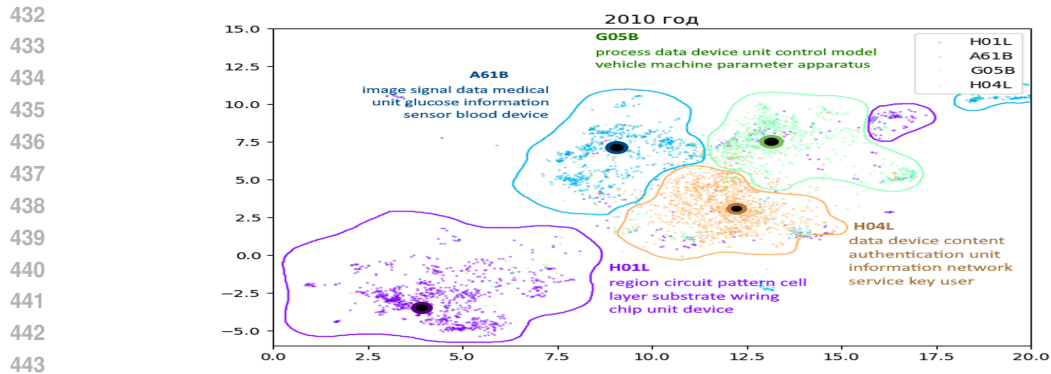
$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} \quad (4)$$

In this case, correct predictions are those where the outcome produced via the proposed approach matches the expectations of experts.

Of the 15 scenarios, 11 were marked as correct predictions and 4 as incorrect predictions:

$$\text{Accuracy} = \frac{11}{15} \approx 73.3\% \quad (5)$$

The testing results showed that the forecasting method is able to correctly identify the main directions of technology development in most cases. Expert assessment confirmed the high accuracy and relevance of the predictions.



445 Figure 3: Distribution of patents in vector space for 2010

446 6 CONCLUSION AND FUTURE WORKS

447

448

449

450 The article considers the challenge of developing a new approach to analyzing and forecasting new

451 technological areas based on patent information (PENTA). The proposed approach allows for effi-

452 cient analysis of large arrays of patent data and identification of potential directions of technological

453 development, which is critical for strategic planning and innovation.

454 The article considered various approaches and algorithms used for patent data analysis, including

455 machine learning and natural language processing (NLP) methods. The capabilities and limitations

456 of each method were studied, which made it possible to select the optimal solutions to develop a

457 novel approach.

458 The authors developed methods for constructing a patent landscape model and a forecast model, as

459 well as a method for assessing development directions in technological areas. Particular attention

460 was paid to algorithms for optimizing the clustering threshold and forming associations between

461 patents and technological directions.

462 The proposed methods uses the PatentBERT model to transform patent texts into a vector representa-

463 tion, the DBSCAN method for patent clustering, and the Adjusted Rand Index metric to optimize the

464 clustering threshold. UMAP and LDA methods are used to visualize the patent landscape, providing

465 a clear observation of the structure and relationships of the data.

466

467 The next steps will be to develop an understanding of model and forecasting accuracy more deeply

468 by applying the tools to a wider range of technology areas. In a world where the cost of technological

469 innovation is high, opportunities to better understand where to focus, capitalize on advancements

470 and ‘break new ground’ can be of value to research, government and business organizations in many

471 ways. An informed approach to research and innovation is a fundamental element of successful

472 foresighting.

473 REFERENCES

- 474
- 475
- 476 Laith Alzubaidi, Jinglan Zhang, Amjad J. Humaidi, Ayad Al-Dujaili, Ye Duan, Omar Al-Shamma,
- 477 J. Santamaría, Mohammed A. Fadhel, M. Al-Amidie, and Laith Farhan. Review of deep learning:
- 478 concepts, CNN architectures, challenges, applications, future directions. *Journal of Big Data*, 8
- 479 (1), 2021.
- 480
- 481 S. M. Anwar. Medical image analysis using convolutional neural networks: a review. *Journal of*
- 482 *Medical Systems*, 42:1–13, 2018.
- 483
- 484 L. Aristodemou and F. Tietze. The state-of-the-art on intellectual property analytics (IPA): A liter-
- 485 ature review on artificial intelligence, machine learning and deep learning methods for analysing
- intellectual property (IP) data. *World Patent Information*, 55:37–51, 2018.

- 486 H. Bekamiri, D. S. Hain, and R. Jurowetzki. Patentsberta: A deep nlp based hybrid model for patent
487 distance and classification using augmented sbert. *Technological Forecasting and Social Change*,
488 206:123536, 2024.
- 489 D. Berezkin, I. Kozlov, and P. Martynyuk. Predictive analytics of scientific and technological trends
490 for decision making in university management. *Procedia Computer Science*, 234:270–277, 2024.
- 491 L. Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
- 492 Tom Brown et al. Language models are few-shot learners. In *Advances in neural information*
493 *processing systems*, volume 33, pp. 1877–1901, 2020.
- 494 S. Das, A. Tariq, T. Santos, S. S. Kantareddy, and I. Banerjee. Recurrent neural networks (RNNs):
495 Architectures, training tricks, and introduction to influential research. In O. Colliot (ed.), *Machine*
496 *Learning for Brain Disorders*, volume 197 of *Neuromethods*, pp. 117–138. Humana, 2023.
- 497 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
498 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*
499 *the North American Chapter of the Association for Computational Linguistics: Human Language*
500 *Technologies*, volume 1, pp. 4171–4186, 2019.
- 501 Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for
502 discovering clusters in large spatial databases with noise. In *Proceedings of the 2nd International*
503 *Conference on Knowledge Discovery and Data Mining (KDD-96)*, pp. 226–231, 1996.
- 504 J. H. Friedman. Stochastic gradient boosting. *Computational Statistics and Data Analysis*, 38(4):
505 367–378, 2002.
- 506 A. Gillioz. Overview of the transformer-based models for NLP tasks. In *Proceedings of the 15th*
507 *Conference on Computer Science and Information Systems (FedCSIS)*, pp. 179–183. IEEE, 2020.
- 508 Y. Hoshino, Y. Utsumi, Y. Matsuda, Y. Tanaka, and K. Nakata. IPC prediction of patent documents
509 using neural network with attention for hierarchical structure. *PLoS One*, 18(3), 2023.
- 510 J. Hou, S. Tang, and Y. Zhang. A novel technology life cycle analysis method based on LSTM and
511 CRF. *Scientometrics*, 129:1173–1196, 2024.
- 512 Y. Jeong, H. Jang, and B. Yoon. Developing a risk-adaptive technology roadmap using a bayesian
513 network and topic modeling under deep uncertainty. *Scientometrics*, 126:3697–3722, 2021.
- 514 T. Ji, N. Self, K. Fu, Z. Chen, N. Ramakrishnan, and C.-T. Lu. Citation forecasting with multi-
515 context attention-aided dependency modeling. *ACM Transactions on Knowledge Discovery from*
516 *Data*, 18(6):1–23, 2024.
- 517 Tong Ji, S. Nguyen, and X. Bai. Automated feature engineering for machine learning: A critical
518 survey. In *Proceedings of the 14th ACM International Conference on Web Search and Data*
519 *Mining (WSDM '21)*. Association for Computing Machinery, 2021.
- 520 Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hervé Jégou, and Tomas
521 Mikolov. Fasttext.zip: Compressing text classification models. arXiv:1612.03651, 2016.
- 522 S. Jun, S. S. Park, and D. S. Jang. Patent management for technology forecasting: A case study of
523 the bio-industry. *Journal of Intellectual Property Rights*, 17(6):539–546, 2012.
- 524 D. Jurafsky and J. H. Martin. Hidden markov models in speech and language processing. *TLTB*, 2:
525 27–68, 2023.
- 526 C. Lee, O. Kwon, M. Kim, and D. Kwon. Early identification of emerging technologies: A machine
527 learning approach using multiple patent indicators. *Technological Forecasting and Social Change*,
528 127:1–13, 2017a.
- 529 J. S. Lee and J. Hsiang. Patentbert: Patent classification with fine-tuning a pre-trained BERT model.
530 arXiv:1906.02124, 2019.

- 540 K. Lee, D. Go, I. Park, and B. Yoon. Exploring suitable technology for small and medium-sized
541 enterprises (SMEs) based on a hidden markov model using patent information and value chain
542 analysis. *Sustainability*, 9(7):1100, 2017b.
- 543 Y. Li, X. Wang, C. Chen, C. Jing, and T. Wu. Exploring firms’ innovation capabilities through
544 learning systems. *Neurocomputing*, 409:27–34, 2020.
- 546 Zhen Li, Jinyu Yu, Jiejun Yang, Weihong Wang, Lei Yang, and Rui Xia. Generative multimodal
547 data augmentation for low-resource multimodal named entity recognition. In *Proceedings of the*
548 *32nd ACM International Conference on Multimedia (MM ’24)*, pp. 7336—7345. Association for
549 Computing Machinery, 2024.
- 550 J. B. MacQueen. Some methods for classification and analysis of multivariate observations. In
551 *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*, volume 1,
552 pp. 281–297, 1967.
- 554 P. K. Manoj, Z. J. Mohamed, and Sh. Subramarian. Method and system for natural language pro-
555 cessing, 2015. US Patent No. US-20150263925-A1.
- 556 Tomas Mikolov, Greg Corrado, Kai Chen, and Jeffrey Dean. Efficient estimation of word represen-
557 tations in vector space. In *Proceedings of the Workshop at ICLR*, pp. 1–12, 2013.
- 558 F. Nielsen. *Introduction to HPC with MPI for Data Science*, chapter 8. Springer, 2016.
- 560 J. Peal. Bayesian networks: A model of self-activated memory for evidential reasoning. In *Proceed-*
561 *ings of the Annual Meeting of the Cognitive Science Society*, volume 7, pp. 329–334, 1985.
- 563 M.-C. Popescu, V. Balas, L. Perescu-Popescu, and N. Mastorakis. Multilayer perceptron and neural
564 networks. *WSEAS Transactions on Circuits and Systems*, 8(7):579–588, 2009.
- 565 Purwono, A. Ma’arif, W. Rahmiani, H. Imam, H. I. K. Fathurrahman, A. Frisky, and Q. M. U. Haq.
566 Understanding of convolutional neural network (CNN): A review. *IJRCS*, 2(4):739–748, 2022.
- 567 Shumaila Qaiser and Ramsha Ali. Text mining: Use of TF-IDF to examine the relevance of words
568 to documents. *International Journal of Computer Applications*, 181(1):25–29, 2018.
- 570 Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using siamese BERT-
571 networks. arXiv:1908.10084, 2019.
- 572 M. V. Santiago, T. Frosch, E. Weber, E. Csernai, E. Gerencser, B. Lantos, A. Kerekes, A. Lanczky,
573 and S. J. Ingles. Workflow predictive analytics engine, 2021. European Patent Application. EP 3
574 826 029 A1.
- 576 Sourav Sarkar, Dading Feng, and Swapna S. Karmaker. Exploring universal sentence encoders for
577 zero-shot text classification. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter*
578 *of the Association for Computational Linguistics and the 12th International Joint Conference on*
579 *Natural Language Processing*, volume 2, pp. 135–147, 2022.
- 580 R. C. Staudemeyer and E. R. Morris. Understanding LSTM – a tutorial into long short-term memory
581 recurrent neural networks. arXiv:1909.09586, 2019.
- 582 J. Sunghae. A new patent analysis using association rule mining and boxjenkins modeling for
583 technology forecasting. *Information (Japan)*, 16:555–562, 2013.
- 585 K. M. Tarwani and S. Edem. Survey on recurrent neural network in natural language processing.
586 *International Journal of Engineering Trends and Technology*, 48(6):301–304, 2017.
- 587 United States Patent and Trademark Office. Bulk data storage system (BDSS). [https://](https://bulkdata.uspto.gov/)
588 bulkdata.uspto.gov/. Accessed: 2025-02-27.
- 590 P.-C. G. Vassiliou and A. C. Georgiou. Markov and semi-markov chains, processes, systems and
591 emerging related fields. *Mathematics*, 9(19):2490, 2021.
- 592 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
593 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. arXiv:1706.03762, 2023.

594 Zhaopan Wang, Liming Wang, Zhipin Zhao, Muxin Wu, Chen-Chieh Lyu, Huasen Li, Deli Cai,
595 Luyu Zhou, Shuai Shi, and Zhaoxiang Tu. GPT4Video: A unified multimodal large language
596 model for instruction-followed understanding and safety-aware generation. In *Proceedings of the*
597 *32nd ACM International Conference on Multimedia (MM '24)*, pp. 3907—3916. Association for
598 Computing Machinery, 2024.

599 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A
600 comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and*
601 *Learning Systems*, 32(1):4–24, 2020.

602
603 J. Zhao, Z. Dong, X. Yao, and Xi Xi. Optimizing collaboration decisions in technological innovation
604 through machine learning: identify trend and partners in collaboration-knowledge interdependent
605 networks. *Annals of Operations Research*, pp. 1–42, 2024.

606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647