

Developing Horizon Scanning Methods for the Discovery of Scientific Trends

Maja Karasalo

Department of Decision Support Systems
Swedish Defence Research Agency
SE-164 90 Stockholm, Sweden
maja.karasalo@foi.se

Johan Schubert

Department of Decision Support Systems
Swedish Defence Research Agency
SE-164 90 Stockholm, Sweden
johan.schubert@foi.se

Abstract—In this application-oriented paper, we develop a methodology and a system for horizon scanning of scientific literature to discover scientific trends. Literature within a broadly defined field is automatically clustered and ranked based on topic and scientific impact, respectively. A method for determining the optimal number of clusters for the established Gibbs sampling Dirichlet multinomial mixture model (GSDMM) algorithm is proposed along with a method for deriving descriptive and distinctive words for the discovered clusters. Furthermore, we propose a ranking methodology based on citation statistics to identify significant contributions within the discovered subject areas.

Keywords—horizon scanning; scientometrics; Gibbs sampling; Dirichlet multinomial mixture model; entropy; clustering; HSTOOL

I. INTRODUCTION

Horizon scanning methods aim to discover changes, disruptions, and trends with the potential to influence the development of a particular area of interest significantly. For scientific literature, the goal of horizon scanning is to discover emerging or rapidly growing research areas.

To scan broad scientific fields without making presumptions about specific topics worthy of further studies, large numbers of scientific articles must be included in the scanning process. This requirement motivates the need for a semiautomatic approach, where software tools provide some initial filtering and structuring of the data.

In this paper, we propose a method for semiautomatic horizon scanning of scientific literature and present the horizon scanning system HSTOOL that supports the proposed method. The goal of the method is to identify rapidly developing fields and their most significant contributions by first scanning the scientific literature using relatively general search criteria and then structuring and filtering the discovered articles. HSTOOL accesses the Thomson Reuters Web of Science¹ (WOS) Core Collection through a set of APIs that allow searches and retrieval of article data, as well as citation statistics.

The key steps of the method are clustering of the discovered literature to identify topics and ranking of articles in the resulting clusters based on scientific citation statistics to find the most significant contributions within the respective topic.

We use the Gibbs sampling Dirichlet multinomial mixture model (GSDMM) algorithm [1] for clustering and introduce a complementary method to determine the optimal number of clusters. We find the optimal clustering by evaluating the quality of placement of every article in each specific cluster using an

entropy measure [2, 3]. Furthermore, we develop a method for automatically presenting two sets of descriptive words for each cluster based on the cluster's contents. The first set consists of the words that most often occur in the cluster, while the second set consists of the most distinctive words in the sense that their occurrence throughout the entire set of articles is concentrated in the current cluster. In combination, the sets provide a description of the articles that are part of the cluster and an account of what primarily distinguishes these articles from articles in other clusters.

For scientific ranking, we propose a set of scientometric measures that identify articles that have made a significant impact in the respective fields. Influence is measured as either collecting many citations over a short period of time, or having a strong citation trend, or frequently being cited in prestigious journals. Finally, the measures are aggregated into a total ranking within each discovered cluster. The top-ranked articles can thus be selected for detailed study.

The paper is organized as follows. In Section II, we describe a workflow model that contains all process steps of searching, organizing and analyzing scientific articles. In Section III, we develop methods for performing horizon scanning to discover and analyze trends in scientific literature. In Section IV, we develop processes for scientific trend discovery and describe a literature scanning system. We apply the system to a case study of literature on military applications of artificial intelligence (Section V). Finally, conclusions are provided in Section VI.

II. WORKFLOW

Fig. 1 shows the proposed workflow of horizon scanning of scientific literature in five steps. The process is intended to facilitate scanning of broad areas defined by a general topic search string (step 1). Once a search has been performed and records downloaded (step 2), topics are automatically discovered using a clustering algorithm that groups the scientific articles based upon textual contents (step 3). Clusters of articles can then be selected for further studies. To find the key contributions from a cluster of interest, a ranking method is proposed that uses a set of statistical citation measurements to capture various aspects of scientific impact (step 4). Once top-rated contributions for a subject area have been identified, a manageable subset of articles can be selected for detailed studies (step 5).

¹ <http://www.webofknowledge.com> (March 2019).

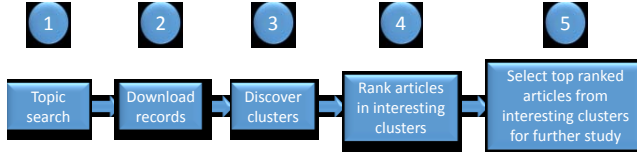


Fig. 1. Proposed workflow for horizon scanning of scientific literature.

III. METHODOLOGY

In this section, we describe methods for searching scientific literature, clustering articles in groups that correspond to subject areas and evaluating the scientific impact of all articles with citation statistics.

A. Searching Scientific Publications

All searches are performed using search terms provided by subject matter experts. These search terms should be tested before use in HSTOOL to ensure that they yield results within the area of interest.

We use HSTOOL to search for publications through an API that provides access to WOS. We limit the search to the Core Collection because that database has the citation statistics that we need for scientometric analysis and ranking of articles.

B. Clustering of Articles

Once a search result has been downloaded from WOS (Fig. 1, steps 1–2) we want to group all articles that deal with the same subject area into a cluster to be treated as a separate subproblem (Fig. 1, step 3) and then use scientometric information to determine which articles within each cluster are most important to that area (Fig. 1, step 4).

In the following two sections, we describe how to use a GSDMM algorithm to organize articles into clusters with common subject areas and how we determine the optimal number of clusters. It is important to point out that the search terms used in the previous step are not used in the clustering phase.

1) *Clustering with GSDMM*: To group articles within the same subarea, we use the abovementioned GSDMM [1, 4]. Simply described, this method starts from a large number of clusters and a random distribution of articles among clusters. Then, the method examines each article to determine if it fits better in any other cluster than where it is currently placed. This procedure is repeated iteratively for all the articles until there are no more changes.

The method proceeds by comparing for every article all words in the article’s title and abstract with the corresponding words in all other articles. If a word is missing or occurs a different number of times when comparing with another article, the probability that these articles belong together is assigned a lower value. These probabilities are combined for all articles within each cluster (and also for the cluster where the article is currently located). This results in an evaluation for all clusters of how well this article fits into all the different clusters. Then, the article is moved to a cluster where it fits well according to these probabilities.

The procedure is applied to all articles and repeated iteratively until all articles are placed in their best clusters.

During the process, the number of clusters will decrease dramatically, often by 80–85%. For example, if we start with

500 clusters and thousands of articles, we can finish with 75–100 clusters.

The clustering process is performed by a sequence of Gibbs sampling iterations. During each iteration, we calculate the probability of each article belonging to each cluster k , resulting in the probability that the article should be moved to that cluster.

We have [1]

$$p_{aki}(k_d = k | \vec{k}_{-d}, \vec{d}) \propto \frac{m_{k,-d} + \alpha}{D - 1 + K\alpha} \times \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (n_{k,-d}^w + \beta + j - 1)}{\prod_{i=1}^{N_d} (n_{k,-d} + V\beta + i - 1)} \quad (1)$$

where on the left-hand side, k_d is the cluster position of article d , k is the k th cluster, \vec{k}_{-d} is the set of cluster positions of all other articles excluding d , and \vec{d} is the set of all articles. In the first term on the right-hand side, $m_{k,-d}$ is the number of articles in cluster k not including d , α is a cluster parameter set to 0.1 in our test case, D is the total number of articles under consideration, and K is the initial number of clusters. In the second term on the right-hand side, w is the w th word of article d , N_d^w is the number of times word w appears in article d , $n_{k,-d}^w$ is the number of times word w appears in cluster k when article d has been removed, β is a cluster parameter that will determine the number of final clusters, N_d is the number of words in article d , $n_{k,-d}$ is the number of words in cluster k when article d has been removed, and V is the number of words in the vocabulary.

During the first iteration, a new cluster position is sampled for each article using (1). After each sampling, (1) is updated. When all articles in D have been reassigned to a new cluster position, the second iteration starts. The process continues for a fixed number of iterations. The final cluster positions of all articles at the last iteration is the result of the clustering process.

2) *Managing the number of clusters*: To select the best number of clusters, we need to evaluate various options. To this end, we evaluate various numbers of clusters based on the quality of clustering.

The GSDMM algorithm does not require a predetermined number of clusters to assign the articles of a given corpus. However, the number of clusters depends on parameter $\beta \in (0, 1)$ that appears in (1). A value of β near zero results in many clusters, while β near one produces fewer clusters.

Several standard internal clustering performance metrics [5] utilize some definition of distance between data points. However, since the GSDMM algorithm does not utilize any distance measure between documents to define clusters, these metrics are inapplicable.

Instead, we focus on the articles that have been clustered and study how well they fit in the clusters where they have been placed.

Each article has a probability distribution across all clusters that indicates the probability that each cluster is the optimal location for that article (1). This distribution is calculated and used in the clustering process for GSDMM and is recalculated in each step of the clustering process for all articles.

At the end of the clustering process, we use the final calculated probability distribution for each article. This is a distribution over all initial clusters, although most of the original

clusters are empty at the end of the clustering process and thus have a nearly zero probability.

We consider $\{p_{dki}\}$, where p_{dki} is the probability that article d belongs to cluster k at iteration i (1), with

$$\sum_{k=1}^K p_{dki} = 1 \quad (2)$$

for any constant d and i , and where K is the initial number of clusters.

If the placement of a particular article is almost certain, that article will have a probability near one for the respective cluster. Sometimes, an article may have more than one probability that is not near zero because the placement is uncertain. A clustering can be considered to be of high quality if as many articles as possible have as certain a placement as possible. Consequently, the entropy of the probability distribution is a good measure of the quality of placement of a particular article [2]. To study the convergence of the GSDMM algorithm, we calculate at each Gibbs sampling iteration i the entropy for each article d as

$$Ent_{di} = - \sum_{k=1}^K p_{dki}(k_d = k | \vec{k}_{-d}, \vec{d}) \log[p_{dki}(k_d = k | \vec{k}_{-d}, \vec{d})]. \quad (3)$$

To determine the quality of a specific clustering (i.e., the clustering at a specific iteration i for a specific value of β), we calculate its entropy as

$$Ent_i = \sum_{d=1}^D Ent_{di}. \quad (4)$$

Fig. 2 shows the convergence of Ent_i for $i \in [0, 14]$, averaging over 100 runs for a test case. The clustering quality is not significantly improved after 10 iterations. However, entropy convergence can vary between runs, which would motivate a dynamic choice of iterations, whereby the entropy reduction rate determines when the algorithm is complete. This would be an interesting future direction to investigate.

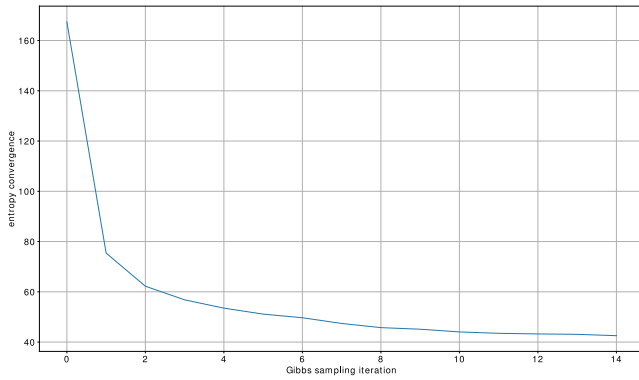


Fig. 2. Entropy convergence (4) over 15 Gibbs sampling iterations for all articles in a test case, averaging over 100 runs.

From the above, it follows that a good measure of quality of the entire partition of all articles for a particular clustering process is the sum of entropy over all articles after the final iteration, where Ent_{14} is the sought-after entropy to be minimized.

In Fig. 3, the average number of discovered clusters is shown for a test case for various values of parameter β . If β is small, we obtain a large number of remaining clusters at the end of the process. The number of clusters drops rapidly if β is increased. For values of β above 0.2, the decrease in the number of final clusters is more gradual. It is clear that choosing the right value of β is key to obtaining an appropriate number of clusters for the contents of the corpus.

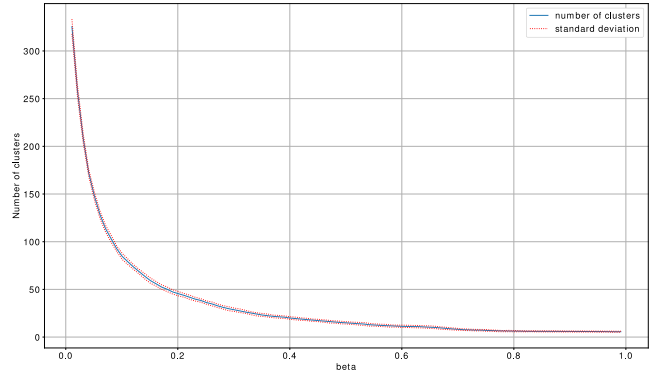


Fig. 3. Average number of clusters discovered by GSDMM as a function of β , averaged over 100 runs for each value of β for a test case.

To find the best number of final clusters, we study the entropy at the end of each clustering process for values of β between zero and one. The results are shown in Fig. 4. As β increases, there is a decline in the final entropy for each clustering process. Note that most of the decline in entropy occurs when β is increased to 0.1. However, the number of clusters keeps decreasing as β approaches 1, as shown in Fig. 3, without any improvement in entropy (Fig. 4), which ultimately results in a few large clusters, each containing multiple topics. Ideally, we want to find a partition that has well-defined clusters that correspond to subject areas yet has the lowest possible entropy.

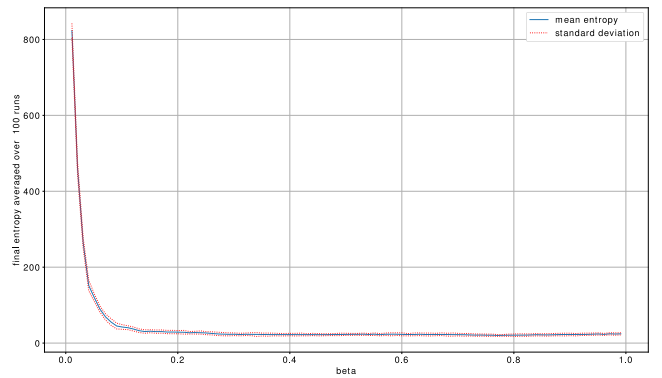


Fig. 4. Final entropy Ent_{14} (4) for a test case summed over all articles as a function of β , averaged over 100 runs.

To estimate the correct number of clusters, the final entropy derived from clusterings with various values of β is calculated. It is evident that there is a change of behavior of the entropy at a point that we consider to yield the best number of clusters; that point is determined as follows [3]. The concave lower envelope of entropy is determined by a convex hull algorithm. At any abscissa, the envelope function is bisected into left and right parts. The acute angle between the left and right line segments is minimized across all bisection values of abscissa, and the minimizing abscissa is selected as the best value of β .

C. Describing the Contents of Clusters

In this section, we outline a method for describing the contents of a cluster. A high-level description is given by the most representative and the most distinctive words.

The most representative words are those that most often occur in the cluster. For cluster k , we have $F_k^w = n_k^w$, where n_k^w is the number of times word w occurs in cluster k . We rank all words in cluster k according to F_k^w and present the highest-ranked words with the maximum F_k^w as representatives of cluster k .

Words that distinguish a cluster from other clusters are determined by calculating the entropy of each word in the corpus as

$$E_w = - \sum_{k=1}^K \frac{n_k^w}{\sum_{j=1}^K n_j^w} \log \left(\frac{n_k^w}{\sum_{j=1}^K n_j^w} \right) \quad (5)$$

where E_w is the entropy of word w , and K is the number of clusters. For each cluster k , the words in this cluster with the lowest entropy (i.e., the words that occur in the least number of clusters) are listed as distinctive words. We have

$$E_k^w = E_w | n_k^w > 0 \quad (6)$$

where E_k^w is the entropy of word w in cluster k . We rank all words in cluster k according to E_k^w and present the highest-ranked words with the minimum E_k^w as the most distinctive words.

Together, F_k^w and E_k^w identify the most representative and distinctive words for each cluster, describing the contents of that cluster.

D. Ranking of Articles within Clusters

The ranking of articles is done using citation statistics in several different ways [6]. We use the statistics provided by Thomson Reuters' WOS. With these statistics, we can rank all articles based on the interest that other scientists have expressed according to their citations.

Our focus is on finding the most important articles within the clusters. This is done independently for each cluster by ranking all its articles. The ranking results from several independent methods with measures that perform alternative assessments.

We start by calculating the number of citations for an article during each of the preceding six years. We then define four different impact measures based on citation impact and citation trends for all articles. Using the four measures, the articles

within each cluster are assigned four alternative impact rankings that are then aggregated into a total ranking. The aggregation of the four rankings is designed to maximize robustness such that no single method dominates the final ranking. The process is repeated independently for each cluster.

1) *Impact measures*: The first measure is called *Impact1*. With this measure, we can rank all articles within a cluster according to the number of times they have been cited in the WOS database over the past year (i.e., the preceding 365 days). This can be done by the operator *citingArticles* in the WOS API. We denote by s_{1j}^k the number of citations of A_j (i.e., the numerical value of the impact measure *Impact1*), where A_j is the j th article in the search, and k is the cluster position of A_j . The highest-ranked article A_j is that with the maximum value of s_{1j}^k for all $\{A_j\}$.

The second impact measure is called *Impact5*. This measure is similar to *Impact1*, except that it includes all citations over the past five years. We denote by s_{5j}^k the number of citations of A_j over the past five years. With *Impact5*, we rank all articles in the second ranking independently from the ranking made with *Impact1*.

The third impact measure is called *ImpactAIS*. Similarly, to *Impact5*, this measure uses citation statistics from the past five years. It is extended by weighting the source according to the source's importance with the Article Influence Score (AIS).

AIS is a measure developed to quantify the importance of a journal's articles during the first five years after publication. AIS is calculated for all publications covered by the Journal Citation Reports² (JCR).

For our purpose, Thomson Reuters provides AIS for all JCR journals. We have

$$s_{AISj}^k = \sum_{A_l \in Y_j} AIS(A_l) \quad (7)$$

where s_{AISj}^k is the number of citations in the preceding five years of A_j , where each citation is weighted by AIS of the citing source, and Y_j is the set of citing articles in the past five years.

The fourth impact measure is called *ImpactReg*. This method performs a least-squares fit of a line to data on five-year citations changes (based on six years of data) for each article. The method ranks all articles according to the average change in citations during these five years, as defined by the slope of the regression line.

The purpose of this method is to capture new articles with a strong trend that have not yet received enough citation coverage to receive high rankings by *Impact1*, *Impact5*, and *ImpactAIS*. We denote by s_{Regj}^k the slope of the regression line of six data points. We use s_{Regj}^k in cluster k as our fourth independent ranking of all articles A_j in the cluster.

2) *Combining all impact measures for aggregated ranking*: The measures derived in the previous section capture different aspects of scientific impact. The aggregated ranking should be able to reflect all these different aspects. A fairly good ranking

² <http://www.webofknowledge.com/JCR> (March 2019).

by all four measures should result in a fairly good aggregated ranking. Furthermore, to receive an acceptable aggregated ranking index, it should be sufficient for an article to have an excellent ranking by one measure, even if the rankings by the other measures are mediocre. Finally, to ensure robust sampling, we want to eliminate any skewness in the distribution for a particular measure. In what follows, we derive a method for aggregating the four impact measures into an overall ranking that meets these criteria.

When selecting r articles for further study from m articles ($r \leq m = |\{A_j\}|$) contained in cluster k , we use the four impact measures calculated for the articles in that cluster. For each impact measure, *Impact1*, *Impact5*, *ImpactAIS*, and *ImpactReg*, we sort all articles $\{A_j\}$ in cluster k in the decreasing order of impact according to $\{s_{ij}^k\}_j$ and renumber all articles within this cluster in the same decreasing order. Thus, the first article (A_1) has the highest impact according to s_{i1}^k within the current cluster k .

We assign a ranking score to r selected articles $\{A_j\}_{j=1}^r$. For article A_j in cluster k that received the j th highest $\{s_{ij}^k\}_j$, we calculate a ranking score with label P_{ij}^k determined according to

$$P_{ij}^k = \frac{r-j+1}{\sum_{l=1}^r l} = \frac{r-j+1}{\frac{1}{2}r(r+1)} \quad (8)$$

where $i = \{1,5,AIS,Reg\}$, and j is the index of article A_j in position j in the ranking of all articles. If $r < j \leq m$, then $P_{ij}^k = 0$ applies by definition.

Since

$$\sum_{j=1}^r P_{ij}^k = 1 \quad (9)$$

we can consider $\{P_{ij}^k\}_j$ as a probability distribution where P_{ij}^k is the probability that A_j is the most preferred article according to impact measure i . This approach turns out to be immediately useful: for rankings that take into account more than one measure, we will use the calculated $\{P_{ij}^k\}_j$ instead of $\{s_{ij}^k\}_j$ because the former is more robust, as some bias in the distribution of $\{s_{ij}^k\}_j$ is eliminated, since $\{P_{ij}^k\}_j$ decreases linearly for all $\{s_{ij}^k\}_j$.

Consequently, we substitute in place of scores $\{s_{ij}^k\}_j$ for each measure the corresponding ranking scores $\{P_{ij}^k\}_j$ and calculate the probabilistic sum of all ranking scores $\{P_{ij}^k\}_i$ for each article A_j . This will be the total measure we use for the final ranking of articles in each cluster.

Within each cluster, we have so far had four different numberings with an individual numbering for each impact measure, since we have sorted all articles separately according

to $\{s_{ij}^k\}_j$. We now number all articles within each cluster such that j always refers to the same article A_j for $\{P_{ij}^k\}_i$ and $\{s_{ij}^k\}_j$.

Finally, we calculate the total ranking score $\{P_{kj}^{Total}\}_j$ for each article A_j . We obtain

$$P_{kj}^{Total} = 1 - \prod_{i \in \{1,5,AIS,Reg\}} (1 - P_{ij}^k) \quad (10)$$

for each article A_j and cluster k , where P_{ij}^k is the ranking score of article A_j according to measure $i \in \{1,5,AIS,Reg\}$. This is the probabilistic sum of all $\{P_{ij}^k\}_i$ [7].

This is our final ranking of articles in cluster k . We can now select the highest-ranked articles within each cluster for further study, as shown in Fig. 1.

IV. SYSTEM DESCRIPTION

In this section, we provide an overview of the horizon scanning software HSTOOL. HSTOOL is a web application built mainly in Scala³. The functionalities of the software correspond to the proposed workflow and include

1. Topic search in WOS,
2. Downloading of article records to a local database,
3. Clustering of articles according to topics,
4. Ranking of articles within each cluster based on scientometric impact, and
5. Outputting the resulting ranked clusters for further study.

The user interface of HSTOOL is shown in Fig. 5. In the following sections, the functionality listed above is described in further detail.

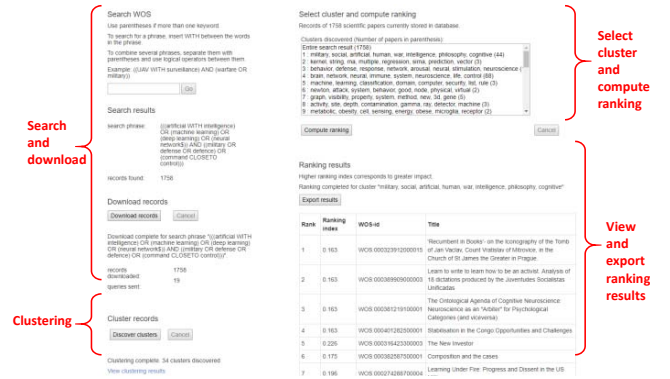


Fig. 5. HSTOOL user interface. Various functionalities are highlighted in red.

A. Searching and Downloading

Topic searches with HSTOOL are performed by combining search terms with logical operators. When a topic search is performed, the number of articles found is displayed in the HSTOOL user interface along with the search string used. A button is available for downloading the search result to a local Postgres database. Records of each article in the search result are

³ <https://www.scala-lang.org> (March 2019).

saved, including keywords, abstract, and scientific field data that will later be used for clustering.

B. Clustering

Pressing the button “Discover clusters” initiates the article clustering algorithm. For each article in the downloaded corpus, a representative text is constructed by combining title, abstract, keywords, subjects, headings, and subheadings provided in the records from WOS. These are the texts that are fed to the GSDMM algorithm.

Once the clustering has been completed, a list of the discovered clusters is displayed on a separate web page, represented by representative and distinctive words for the cluster and a list of the included articles. The discovered clusters are also displayed in a scrollable list in the HSTOOL main view, from which clusters can be selected for ranking.

C. Ranking

To rank articles within a cluster, it is necessary to select the cluster from the scrollable list and click “Compute ranking.” Once the ranking has been completed, a ranked list of articles, including the calculated ranking index of each article, is displayed in the HSTOOL main view.

D. Output

Once the ranking of a cluster has been completed, the results can be exported by clicking the button “Export results.” This action produces a CSV file, including the WOS ID, article title, abstract, keywords, subjects, headings, journal title, ISSN, AIS factor, the estimated impact factors, and the resulting ranking index.

V. CASE STUDY OF ARTIFICIAL INTELLIGENCE IN MILITARY APPLICATIONS

In this section, we report findings and results of a case study carried out to validate the proposed methodology and software tool. The topic of the case study was chosen within the authors’ field of expertise to facilitate the evaluation of the results of the horizon scanning process.

A. Topic Search

A search was performed using a combination of rather broad concepts, aiming to capture articles related to artificial intelligence in the context of defense applications. We compiled a topic search string of the form

$$[\text{AI terms}] \text{ AND } [\text{defense terms}].^4$$

The search resulted in 1358 hits in the WOS Core Collection, with publication years ranging from 1991 to 2019. Fig. 6 shows the number of articles per year for the search result. We will refer to the set of discovered articles as the AI corpus.

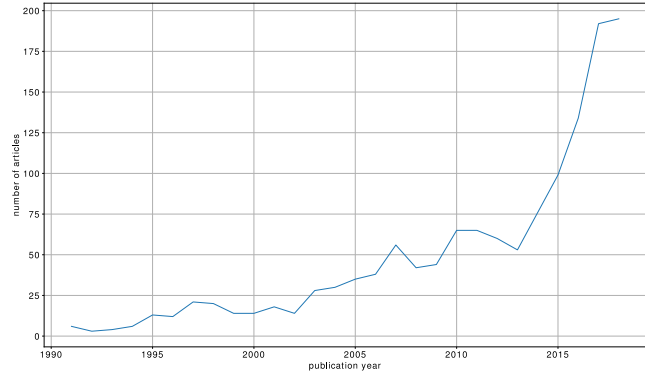


Fig. 6. Number of articles per year in the AI corpus.

B. Clustering Search Results

Clustering of the search results encompasses two steps: first, determining the optimal value of parameter $\beta \in (0, 1)$ (which in turn yields the optimal number of clusters), and, second, performing the actual clustering with the optimal settings. The GSDMM algorithm clusters articles in the course of a set of Gibbs sampling iterations, during which the articles converge to a subset of the initial clusters. The size of this subset is determined by parameter settings. To understand how the algorithm works, we will first study the Gibbs sampling iterations for a fixed value of β , after which we will determine the value of β that yields the best clustering of the AI corpus.

1) *Gibbs sampling iterations and convergence of entropy:* At each Gibbs sampling iteration, the conditional probability p_{dki} given by (1) is calculated for each article d and cluster k , yielding the probability that d is generated by k . Fig. 7 shows how p_{dki} varies over 15 Gibbs sampling iterations for a sample article. At first, the probability density function has spikes at a few different clusters, but for most articles, it converges quickly to a Dirac pulse at a certain cluster.

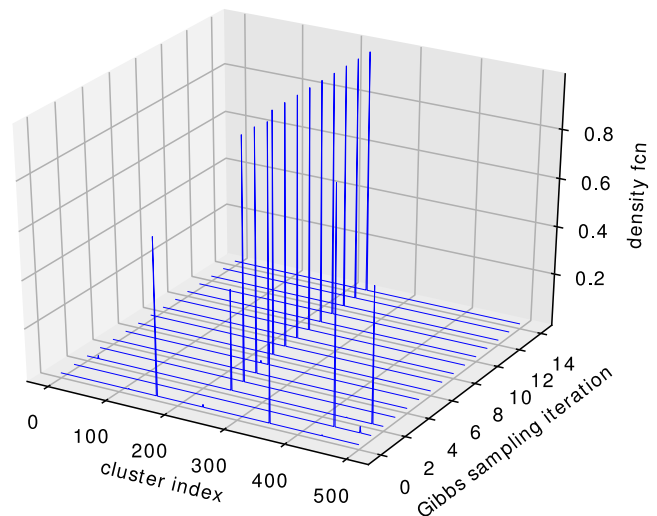


Fig. 7. Probability density function for an article in the AI corpus, as it varies over the 15 Gibbs sampling iterations.

⁴ (“artificial intelligence” OR “machine learning” OR “deep learning” OR “neural network\$”) AND (military OR defense OR defence OR (command

NEAR\1 control)) – the operator NEAR/n signifies that the words on either side of it must be at most n words apart, and \$ denotes the option of a plural s.

To determine the quality of a specific clustering (i.e., the clustering at a specific iteration i for some specified value of β), we calculate its entropy using (4) and the entire set of articles. Fig. 2 shows how the entropy for the entire AI corpus converges for a fixed value of β .

2) *Determining the optimal number of clusters:* Following the approach outlined in Section III.B.2, the final entropy of the entire AI corpus for values of $\beta \in (0, 1)$ is calculated, and the value of β that minimizes the angle of the lower envelope curve is determined to be $\beta = 0.101$, yielding an average of 84 clusters over 100 runs, as shown in Fig. 3.

C. Analysis of Clusters

We perform an individual clustering of the AI corpus using the optimal $\beta = 0.101$, this time yielding 90 clusters. The number of articles in the discovered clusters varies from 1 to 224. To select clusters for further study, we use three criteria:

1. The cluster should be of sufficient size to represent a significant topic for the corpus.
2. The cluster should be well defined, i.e., have as low an entropy as possible.
3. The topic of the cluster should be relevant with respect to the original intention of the search query. Once a subset of clusters has been selected according to steps 1 and 2, clusters with irrelevant topics are removed from this subset.

Fig. 8 shows the number of articles in each of 90 discovered clusters. We will select clusters of size 15 or larger for further study. Fig. 9 shows the mean entropies for all such clusters. Among these, we examine the nine clusters with the lowest entropy in detail.

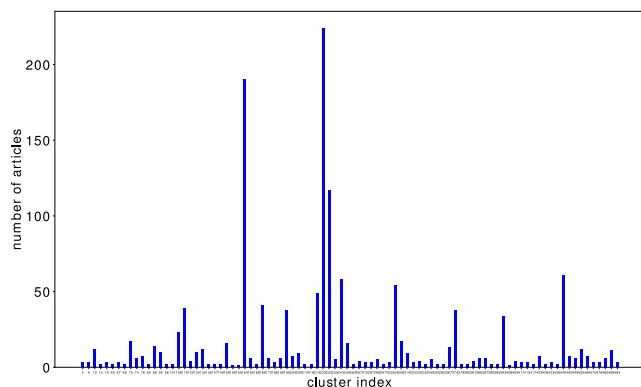


Fig. 8. Number of articles in each of 90 discovered clusters.

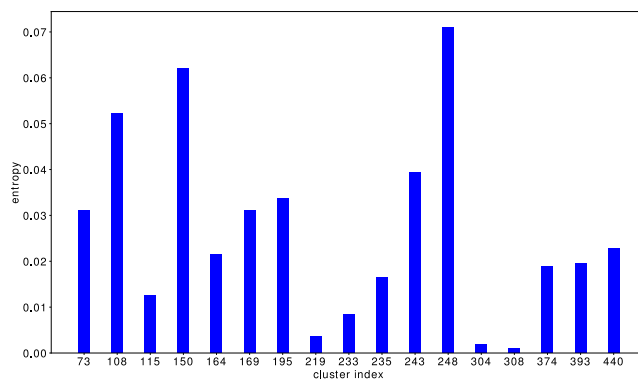


Fig. 9. Entropies for all clusters of size 15 or larger.

The search queries used for topic search in WOS aimed to find articles on AI applications in the military domain. Among the selected clusters, three clusters with irrelevant topics are detected and removed. The descriptive and distinctive words of the six remaining clusters are shown in TABLE I.

TABLE I. DISCOVERED TOPIC WORDS IN THE CLUSTERS CHOSEN FOR FURTHER STUDY

Id	Common words	Distinguishing words	Size
164	image, target, network, recognition, neural	correlator, mstar, foreground, dividing, eo	190
219	classification, signal, network, neural, feature	amc, instantaneous, cepstral, mel, warped	49
233	attack, system, network, detection, computer	multicore, bodyguard, port, nash, protocol	224
235	system, computer, agent, intelligence, decision	illustration, bdi, nec, ner, automates, succession	117
308	game, player, computer, artificial, defense	offense, beginner, dda, neuroevolution, warcraft	17
393	network, sensor, system, application, neural	fence, transceivers, steganography	34

Fig. 10 and Fig. 11 show, respectively, the number of articles in each of the studied clusters, and the number of citations of articles in each cluster over time. It can be noted that even though cluster 233 (attack, system, network) has the most articles and the strongest publications trend, the most cited cluster is cluster 164 (image, target, network), a trend that has been strong over the past 15 years and is still holding. It should be noted that the number of citing articles is based on all citing articles, whether part of the search result or not. A conclusion that can be drawn from Fig. 10 and Fig. 11 is that computer vision applications (cluster 164) remain a dominating topic within the “AI for the military” field, while defense against adversarial attacks for neural networks (cluster 233) has gained interest over the past few years. We will therefore choose the computer vision cluster 164 as an example for further study.

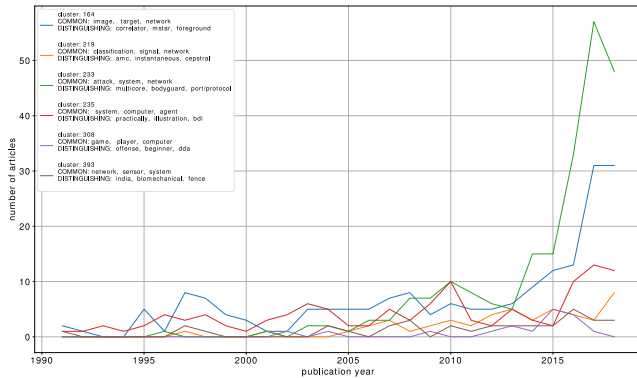


Fig. 10. Number of articles per year in the clusters of interest.

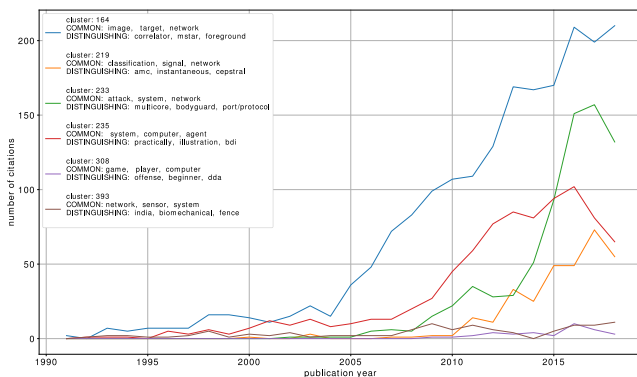


Fig. 11. Number of citations per year generated by the clusters of interest.

All articles in the computer vision cluster are ranked according to the four impact measures. The top results are shown in TABLE II. In the overall ranking, a set of representation learning and object recognition articles receives the highest scores. This is a reasonable result, given the impressive progress made recently with regard to these topics. As discussed, the different impact measures aim to capture different impact aspects. Examining the ranking results for the computer vision cluster, we note that *ImpactAIS* – that takes into account the AIS score of the journal an article is cited in – is the measure most in disagreement with the total ranking. As *ImpactAIS* does not in general reward recent articles, this would be an expected result, given the rapid development in the field of AI for computer vision (Fig. 10 and Fig. 11).

The case study indicates that the proposed horizon scanning methodology and tool are useful for finding trending and significant topics in the scientific literature. The top-ranked

articles within the studied cluster cover topics that have received significant attention in recent years, which validates the soundness of the impact measures and the aggregation method.

VI. CONCLUSIONS

We have developed new methods for horizon scanning, integrated them with an existing clustering method, and implemented all methods in a system for horizon scanning of scientific literature to discover scientific trends. In particular, we have developed methods for finding an optimal number of clusters by developing an entropy-based method that focuses on the clustered articles rather than on the clusters themselves. We conclude and show in a case study that with these methods, we can identify distinct clusters. These clusters can be categorized by automatically producing the most descriptive and distinctive words. Furthermore, we develop methods for a robust ranking of articles based on citation statistics and demonstrate in the case study how to produce an overall ranking of all articles in each category. Overall, these methods automatically discover previously unknown categories, describe such categories with their most important words, rank all articles within each category by importance and deliver categories of ranked articles as the system output.

ACKNOWLEDGMENT

This work was supported by the FOI research project “Horizon scanning”, which is funded by the R&D programme of the Swedish Armed Forces.

REFERENCES

- [1] J. Yin and J. Wang, “A Dirichlet multinomial mixture model-based approach for short text clustering,” in Proc. 20th ACM SIGKDD Int. Conf. Knowl. Discov. Data Min., Aug. 2014. New York: ACM, 2014, pp. 233–242.
- [2] C. E. Shannon, “A mathematical theory of communication,” Bell Syst. Tech. J., vol. 27, pp. 379–423, 623–656, Jul.–Oct. 1948.
- [3] S. Ahlberg, P. Hörling, K. Johansson, K. Jöred, H. Kjellström, C. Mårtensson, G. Neider, J. Schubert, P. Svenson, P. Svensson, and J. Walter, “An information fusion demonstrator for tactical intelligence processing in network-based defense,” Inf. Fusion, vol. 8, pp. 84–107, Jan. 2007.
- [4] K. Nigam, A. K. McCallum, S. Thrun, and T. Mitchell, “Text classification from labeled and unlabeled documents using EM,” Mach. Learn., vol. 39, pp. 103–134, May 2000.
- [5] E. Rendon, I. Abundez, A. Arizmendi, and E. M. Quiroz, “Internal versus external cluster validation indexes,” Int. J. Comput. Commun., vol. 5, pp. 27–34, Mar. 2011.
- [6] J. Mingers and L. Leydesdorff, “A review of theory and practice in scientometrics,” Eur. J. Oper. Res., vol. 246, pp. 1–19, Oct. 2015.
- [7] J. Schubert, “Constructing and evaluating alternative frames of discernment,” Int. J. Approx. Reason., vol. 53, pp. 176–189, Feb. 2012.

TABLE II. TOP 5 RANKED ARTICLES IN THE COMPUTER VISION CLUSTER, CONSIDERING ALL IMPACT MEASURES

Ranking	Title	Year	Impact 1	Impact 5	Impact AIS	Impact Reg
1	Spiking Deep Convolutional Neural Networks for Energy-Efficient Object Recognition	2015	2	2	3	1
2	Remote Sensing Scene Classification by Unsupervised Representation Learning	2017	3	3	4	2
3	Learning Race from Face: A Survey	2014	4	4	7	3
4	Neural networks for automatic target recognition	1995	7	7	6	8
5	Adaptive fusion method of visible light and infrared images based on non-subsampled shearlet transform and fast non-negative matrix factorization	2014	8	6	13	4