# A SERENDIPITOUS RESEARCH PAPER RECOMMENDER SYSTEM

**Benard Magara Maake**
Computer Systems Engineering
Tshwane University of Technology
E-mail: 214832445@tut4life.ac.za

**Sunday O. Ojo**
Computer Systems Engineering
Tshwane University of Technology
E-mail: OjoSO@tut.ac.za

**Tranos Zuva**
Department of ICT
Vaal University of Technology
E-mail: tranosz@vut.ac.za

─Abstract─

In recent times, the rate at which research papers are being processed and shared all over the internet has tremendously increased leading to information overload. Tools such as academic search engines and recommender systems have lately been adopted to help the overwhelmed researchers make right decisions regarding using, downloading and managing these millions of available research paper articles. The aim of this research is to model a spontaneous research paper recommender system that recommends serendipitous research papers from two large and normally mismatched information spaces using Bisociative Information Networks (BisoNets). Set and graph theory methods were employed to model the problem, whereas text mining methodologies were used to process textual data which was used in developing nodes and links of the BisoNets graph. Nodes were constructed from weighty keywords while links between these nodes were established through weightings determined from the co-occurrence of corresponding keywords originating from both domains. Final results from our experiments ascertain the presence of latent relationships between the two habitually incompatible domains of magnesium and migraine. Word clouds indicated that there was no obvious relationship between the two domains, but statistical significance investigations on the terms indicated the presence of very strong associations that formed information networks. The strongest links in the established information networks were further exploited to show bisociations

between the two habitually incompatible matrices. BisoNets were consequently constructed, exposing terms and concepts from two discordant domains that were bisociated. These terms and concepts were utilised in querying the one domain for recommendations in another domain. Hence, serendipitous recommendations were made since our bisociative knowledge discovery methodologies revealed hidden relationships between research papers from diverse domains. Finally, it was postulated that latent relationships exist between two incompatible domains, and when well exploited, it leads to the discovery of new information and knowledge that is useful to researchers in various fields, especially those engaged in multi-disciplinary research. Further research is being conducted to identify outlier linkers and connectors between domains of diverse subjects.

**Key Words:** *Bisociation, Creative knowledge discovery, Information networks, Recommender systems, Serendipity*

**JEL Classification:** L86 Information and Internet Services • Computer Software

## 1. INTRODUCTION

Information and Communication Technologies (ICT) recent advancement has led to the development of Scientific Social Networks (SSNs), which are convenient platforms through which researchers communicate. The quantity of information that these researchers encounter daily as they try to make research decisions keeps growing exponentially (Kötter, Thiel, & Berthold, 2010; Wang, He, & Ishuga, 2018), and the larger it grows, the harder it gets in accessing important information and knowledge (Ferrara, Pudota, & Tasso, 2011). In the year 2015, it was estimated that an average American was consuming 15.5 hours of media content and approximately 74 gigabytes of information daily. This is a staggeringly large amount of information to be consumed by an individual every day. Consequently, this calls for powerful filtering systems that are going to help these online users manage the big data, and receive only what is relevant and useful (Nodus Labs, 2015).

The exponential growth in the number of online research paper publications is rendering current internet-based search systems inadequate in providing better services to researchers. This has posed a big challenge to scholars who are searching for relevant published work (Bollacker, Lawrence, & Giles, 1998, 2000). Unfortunately, researchers do not have the ability and knowledge of processing this large information space and get what they need. On the other hand, search engines and recommender systems that are currently being adopted

to meet these problems are falling short of addressing the quest by scholars for relevant, accurate and serendipitous research papers. Regrettably, it has now become almost impossible to track all existing research paper publications on the internet.

Researchers are increasingly working with interdisciplinary environments where information needs are required from not only one research domain but from a variety of fields. The task of finding such relevant information is a tough and difficult task since there is no one person who is an expert in all fields (Kötter et al., 2010). This has given rise to the need for recommender systems in the field of research papers publications with the goal of helping researchers access relevant, useful and novel research papers. However, it is unfortunate that in practice, good research paper recommender systems are very few, and the current ones have limited functionalities like limiting their focus on citations and textual analysis when processing recommendations (Gipp, Beel, & Hentschel, 2009). Nevertheless, both generic and academic search engines are trying to recommend relevant research papers, but at the same time worsening the problem of information overload since they are constantly flooding the researchers with more journals and research papers articles.

The challenges of an information-rich environment that most people are getting themselves into nowadays is: offering the right type of content to the right person and in the most appropriate way (Benlamri & Zhang, 2014). And in our case, offering serendipitous research papers from domains that are normally regarded as unrelated. Part of the solution necessitates adopting information filters that remove unwanted information, presenting only what is useful and relevant information to researchers, of which recommender systems are part of them (Sridharan, 2014). Hence, recommender systems will be used to identify a set of items that are likely to interest a user (Bogers & Van den Bosch, 2008), further eliminating and minimising the efforts used by users in accessing highly relevant and useful items.

A research paper recommender system needs to have enough information about the user in order to better recommend relevant articles (Beel, Langer, Genzmehr, & Nürnberger, 2013). To increase a researcher's ability to discover new knowledge from a digital library, multi-dimensional and hybridization methodologies in recommender systems are to be utilised (Vellino & Zeber, 2007). On the other hand, recommender systems do come with pitfalls like: not building user confidence (trust failure), not generating any recommendation (knowledge failure), generating incorrect recommendation (personalisation

failure), generating recommendations to meet a wrong need (context failure) (McNee, Kapoor, & Konstan, 2006). That is why (Zarrinkalam & Kahani, 2012) suggested that every recommender system should have background data, input data and a recommendation algorithm for the sake of improving the quality of recommendations.

Though recommender systems are perceived pejoratively, their use has brought out more benefits to online users, for instance, reducing the amount of time and energy customers use in navigating e-commerce sites (Lémdani, Polaillon, Bennacer, & Bourda, 2011). Then again, they can discover and recommend surprisingly new and unexpected information that is fortuitous to the user without necessarily searching for it (André, Teevan, & Dumais, 2009; Sridharan, 2014). Therefore, this research is aimed at modelling and developing a serendipitous research paper recommender system that will recommend research papers from domains that are distant and unrelated.

## 2. LITERATURE REVIEW

To recommend relevant research papers, we need to understand the needs of the researcher, and this can be done through understanding his profile. Several methods can be used to analyse and measure the relationships existing between a researcher and other personalised perspectives surrounding him, such as email communication, social media interaction, project cooperation, etc. Other supplementary components that can be exploited together with user preference when processing recommendations include domain knowledge, user background, learning target and cognitive patterns (Zhao, Wu, Dai, & Dai, 2015). In this research, we will utilise the user's research preference represented by the title of the article that is currently being read by the researcher, which again is used to construct the user's profile and query the database. See Figure-1.

In the work done by (Paraschiv, Dascalu, Dessus, Trausan-Matu, & McNamara, 2016), a paper recommender system was developed that graphically visualised all the related concepts in research papers. Their system automatically tagged paper abstracts with semantic words that summarised the paper content, and then the Latent Semantic Indexing (LSI) algorithm was utilised to create a semantic representation of associated words that co-occurred in similar contexts. Our research will also graphically visualise important terms in research paper titles and then different statistical measures will be used to evaluate the strength of relationships between co-occurring terms, within a single domain as well as between different discordant domains.

To enable cross-domain associations, (Kötter et al., 2010) integrated information units from two dissimilar domains which resulted in the formation of a single network that enabled cross-domain associations. The network from this kind of association whereby latent relationships are discovered between two discordant domains is called Bisociation, and the network formed is termed as a Bisociation Network (BisoNet). We extend on their work by modelling a research paper recommender system that will search, find and recommend articles between the two unrelated domains. We also seek to create BisoNets and further use the bisociation links to make recommendations of unexpected, relevant and surprisingly interesting research papers from those distant domains.

Traditional recommender systems recommend items that are similar to user profiles or what users have indicated as interesting, a problem known as overspecialisation, which leads to obvious and uninteresting recommendations. The overspecialisation problem can be reduced or solved by introducing serendipity (Kotkov, Wang, & Veijalainen, 2016) or randomness (Vellino & Zeber, 2007) into a recommender system, and we will adopt the former method in our research. Knowledge discovery focusses on finding patterns in reasonably well-connected domains, and lately, research is moving into the discovery of connections between ordinarily unconnected domains (Berthold, 2012).

Majority of recommender systems are developed using traditional paradigms that are based on accuracy and prediction metrics. It should be noted that other aspects, such as user satisfaction and experience should also be considered when making recommendations. It is now widely agreed that accurate predictions are very important but insufficient to deploy a good recommender system (Shani & Gunawardana, 2011). Therefore, recommending beyond accuracy and the term-frequency inverse-document-frequency (TF-IDF) metrics will be regarded as an advanced approach to content-based filtering. Finally, different goals or needs of researchers require different approaches and algorithms to solve, hence affirming that there is no specific algorithm which will solve the various aspects of research paper recommendation (Nascimento, Laender, da Silva, & Gonçalves, 2011). Therefore, this has motivated us to use creative knowledge discovery methodologies (bisociation) in solving the problem of overspecialisation in research paper recommender systems while providing serendipitous recommendations (Adamopoulos, 2013).

**Figure-1: Bisociative discovery model for research paper recommender systems**



## 3. METHODOLOGY

### 3.1 Set theory

In this research, we take two known domains that are not related and then we will attempt to find concepts that interconnect them (bridging concepts/ terms), a process known as the closed discovery process (Juršič, Sluban, Cestnik, Grčar, & Lavrač, 2012). Tables 1 and 2 contain the notations that are used to formulate our problem and construct the BisoNet respectively. Domains not related in any way are represented as matrices $M_1$ and $M_2$ (see Figure-2) that intersect at right angles implying their incompatibility or un-relatedness. There is a possibility of having similar themes within one domain, but more importantly, we can have these themes extending and appearing on other domains/ matrices. A shared idea or problem $\pi$ within the two matrices is denoted by a line that runs from one matrix to another (illustrated by the line traversing the orthogonal planes in Figure-2). The concepts $c_1, c_2,$ and $c_3$ on matrix $M_1$ are connected to one another, denoting that they have some sort of relation. Similarly, on matrix $M_2$, the concepts $c_1$ and $c_6$ are connected. This kind of relationship that exists on each separate matrix is known as association. Similar concepts and ideas perceived on these two domains act as links between those two incompatible matrices, forming an association that is known as Bisociation.

**Table 1: Bisociation and Graph Theory Notations**

| Bisociation | | Graph Theory | |
|---|---|---|---|
| Notation | Definition | Notation | Definition |
| $M_1$ | First matrix/ knowledge base/ domain/ field | $V$ | Vertices |
| $M_2$ | Second matrix/ knowledge base/ domain/ field | $E$ | Edges |
| $c_1, \ldots c_6$ | Concepts/ ideas/ themes that run through both matrices | $B$ | BisoNet |
| $D_i$ | A set of concepts $i$ associated with domain $D$ | $\lambda$ | Function to assign a unique label to a vertex |
| $K_i$ | A knowledge base that is a subset of the domain $D$ | $\omega$ | Function weights to the edges |
| $U$ | All the concepts/ ideas/ themes in both matrices $M_1$ and $M_2$ | $k$ | Number of partitions |
| | | $\beta$ | Relevance of concepts |
| $X$ | Concepts that are associated with the problem | $\delta$ | Unexpectedness of a concept in a domain |
| $K_i^R$ | Reference system for knowledge base $i$ | $\mu$ | Novelty of a concept(s) |
| $K_j^R$ | Reference system for knowledge base $j$ | | |
| $\pi$ | A problem or concept that links two incompatible domains | | |
| $t$ | time | | |

Let domain $D_i$ represent a set of concepts associated with a particular domain $i$, and a knowledge base $K_i$ be represented as a subset of domain $D_i$ such that $K_i \subseteq D_i$. Let $U$ be the space that contains all the concepts and ideas in both matrices $M_1$ and $M_2$. Let $c \in U$ represent all the identified concept belong to the universal space $U$. Let $X \subset U$ denote a collection of concepts which belong to the problem $\pi$ that originates from both matrices. Let $R$ denote a reference system that owns exactly one knowledge base for each domain $D_i$, hence $K^R \subset D_i$ will

denote a knowledge base as per the domain $D_i$. A union of the two reference systems can be denoted as $K^R = U_i K_i^R$. Two reference systems $K_i^R$ and $K_j^R$ are said to be unrelated where $(i \neq j)$, meaning that there are no concepts that can be perceived simultaneously in both the matrices. Therefore, bisociation can be denoted using the following expression. Let $\pi$ be utilised to represent a problem, and let $X \subset U$ be used to represent all the concepts that are associated with the problem $\pi$. Additionally, let the two unrelated, habitually incompatible matrices be represented as $K_i^R$ and $K_j^R$ such that $(i \neq j)$. Let domain $D_i$ represent a set of concepts associated with a particular domain $i$, and a knowledge base $K_i$ be represented as a subset of domain $D_i$ such that $K_i \subseteq D_i$. Let $U$ be the space that contains all the concepts and ideas in both matrices $M_1$ and $M_2$. Let $c \in U$ represent all the identified concepts belong to the universal space $U$. Let $X \subset U$ denote a collection of concepts which belong to the problem $\pi$ that originates from both matrices. Let $R$ denote a reference system that owns exactly one knowle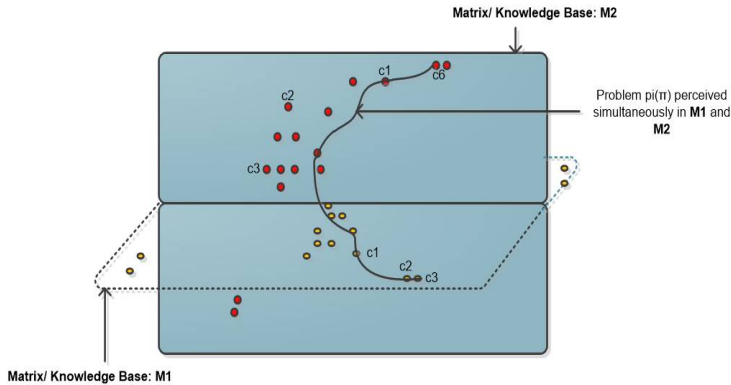dge base for each domain $D_i$, hence $K^R \subset D_i$ will denote a knowledge base as per the domain $D_i$. A union of the two reference systems can be denoted as $K^R = U_i K_i^R$. Two reference systems $K_i^R$ and $K_j^R$ are said to be unrelated where $(i \neq j)$, meaning that there are no concepts that can be perceived simultaneously in both the matrices. Therefore, bisociation can be denoted using the following expression. Let $\pi$ be utilised to represent a problem, and let $X \subset U$ be used to represent all the concepts that are associated with the problem $\pi$.

## 3.2 Graph theory

Additionally, let the two unrelated, habitually incompatible matrices be represented as $K_i^R$ and $K_j^R$ where $(i \neq j)$. Bisociation will occur if and only if elements of $X$ are perceived simultaneously in both matrices $K_i^R$ and $K_j^R$ and at a time $t$.

**Figure-2: Two habitually incompatible domains intersecting at right angles**



Concepts and ideas in a domain can be represented as a graph, where the vertices, $V$ can represent information units, and the edges, $E$ represent relationships between these information units. BisoNets are networks of information units whereby the information units have associations amongst themselves. A BisoNet $B = (V_1, \ldots, V_k, E, \lambda, \omega)$ is a graph data structure that has got attributes with at least a minimum of two partitions $k \geq 2$. Information units are represented as vertices, $v \in V$ and relations between the vertices are the edges of the graph, where a set of edges $E = \{\{u, v\}: u \in V_i; v \in V_j; j \neq i\}$ whereby the two vertices $\{u, v\}$ from two different partitions are connected. For visualisation, a function $\lambda: V \to \Sigma^*$ is used to assign a unique label to a vertex, while another function $\omega: E \to [0, 1]$ assigns the edges with weights.

## 3.3 Text mining

Text mining is a knowledge-intensive process used to assist a person searching for information to interact with document collections over time using a suite of analysis tools. It derives its techniques from data mining, machine learning, natural language processing, information retrieval and knowledge management; all which try to solve the information overload problem. Textual data can be decomposed from long sentences into single terms and words, which in essence can represent the vertices discussed in the previous section. A research paper recommender system will often utilise textual data from research papers to form themes and concepts. Using the graph data structure, information units from a

research paper will be utilised to form vertices, while the relationships between the terms and concepts will be represented with edges. Weights applied to the edges will denote the existence of a relationship between the terms, whereby higher weights will signify the certainty of an existing relationship, whereas lower weights will infer the presence of a weak or no relationship. Bisociation will only occur if BisoNets are formed from one matrix into another in that graph.

## 3.4 Evaluation

The dataset used in this research is a well-researched migraine-magnesium domain pair which was introduced by Swanson (Swanson, 1988) which contains research paper titles retrieved from the PubMed database using two queries/ keywords for the separate domains. The first keyword used was "migraine", for the migraine domain, while the second keyword used was "magnesium", for the magnesium domain. An additional condition to the query that we used was the restriction on the range in terms of date of publishing the articles. We wanted to get only titles that were published before the year 1988, for a simple reason of reproducing bridging terms that were discovered by Swanson (Swanson, 1988). A total of 3360 titles were returned when the migraine keyword was used, and a total of 8843 titles were returned when the magnesium title was used as a query. We, therefore, used 12203 titles from both magnesium and migraine in our analysis to try to discover whether two domains are related in any way, and which are those terms that bridge the two domains. Bridging terms were revealed through co-occurrence of terms, and they assisted in linking the two domains. This was a different approach compared to the Swanson's approach.

## 3.5 Statistical methods to measure relatedness of text in discordant domains

In this work, different statistical significance methods were used to extract semantic links existing between terms of both matrices. Statistical significance in our case was the likelihood that a relationship between two or more discordant domains was caused by something other than chance. The titles retrieved were decomposed into short sentences using a probabilistic model before the traditional pre-processing steps, and this is because the sentence-segmentation-model required intact words and punctuations to form the new-sentences (quasi-documents). The pre-processing stage then followed to resolve issues, such as errors in data, incompleteness, inconsistency and lacking in certain behaviour and trends. Further cleansing was done by removing all the punctuation marks, numbers, white spaces, stop words and transforming all terms into lower case.

A document-term-matrix (DTM) was created with words that occurred at least 10 times. Words that co-occurred in the DTM were counted by multiplying a transposed matrix with its original matrix to get an encoded term-term matrix. See Table 3 below, which shows the number of joint co-occurrences and the frequency of a single term. A sample interpretation of the table is as follows: excretion occurs 239 times, and it co-occurs 69 times with the term renal. It occurs zero times with the terms clin and nephrol.

**Table 3: Sample terms co-occurrence table**

|          | excretion | renal | trace | clin | nephrol | acute |
|----------|-----------|-------|-------|------|---------|-------|
| excretion | 239 | 69 | 1 | 0 | 0 | 8 |
| renal | 69 | 308 | 1 | 0 | 0 | 17 |
| trace | 1 | 1 | 41 | 0 | 0 | 0 |
| clin | 0 | 0 | 0 | 820 | 2 | 0 |
| nephrol | 0 | 0 | 0 | 2 | 50 | 0 |
| acute | 8 | 17 | 0 | 0 | 0 | 180 |

It is not enough to only count the joint occurrence of terms, hence we endeavoured to determine their statistical significance. We utilised a co-occurrence (query) term to calculate the significance of all the joint occurrences that occurred. We utilised the mutual information, dice and log-likelihood measures to implement this work. The results were then sorted out in decreasing mode whereby the most significant co-occurrences appeared top of the list, while the less significant appeared towards the bottom of the list. See Table 4 below.

**Figure-3: Word clouds observations of un-relatedness of two domains**



Dataset Source: PubMed: http://www.ncbi.nlm.nih.gov/pubmed

**Table 4: Statistical significance of terms in both domains**

| | Freq-terms | Freq | MI-terms | MI | Dice-Terms | Dice | LL-Terms | LL |
|---|---|---|---|---|---|---|---|---|
| 1 | vitamin | 87 | vitamin | 5.636471 | vitamin | 1.000000 | magnesium | 142.51101 |
| 2 | magnesium | 81 | rickets | 4.680960 | rickets | 0.100000 | calcium | 45.56182 |
| 3 | calcium | 29 | resistant | 4.027033 | combined | 0.077519 | rickets | 39.38735 |
| 4 | effect | 17 | hypopara | 4.027033 | admin | 0.071065 | admin | 28.05691 |
| 5 | deficiency | 7 | dose | 3.621568 | term | 0.043795 | women | 18.56123 |

**Table 5. Three-column data-frame for the BisoNet construction**

| | from | to | sig |
|---|---|---|---|
| 314 | ergotamine | administration | 45.040669 |
| 613 | blood | levels | 39.536234 |
| 122 | med | state | 48.195894 |
| 1215 | childhood | complicated | 4.732060 |

## 3.6 BisoNet creation (Graph construction)

To visualise the latent relationship that exists within individual domains and between the two distant domains, we construct a BisoNet. The graph was created from a three-column data-frame (see Table 5), the source, the sink and lastly the weight attached to the edge. A triple data frame was constructed encoding the source, edge and weight information of this graph. Each triple consisted of a target term, a term that co-occurred and the significance of their joint co-occurrence. To fully complete the BisoNet, the following packages from the R-programming environment were utilised; 'tm package' for text processing, 'igraph package' for constructing the BisoNet graph. Preliminary results from the word clouds (see Figure-3) at the onset indicated that the two domains were not related. However, we got different results when we calculated the statistical significance of all terms in all domains (see Table 4). Latent relationships between different domains were also revealed necessitating the action of recommendation.

## 4. RESULTS

From Table 4, it is clear that different statistical measures give varying statistical significance results. The log-likelihood outperformed the dice and the mutual information measures by giving the best values for the relationships that exist between terms. To further see how this relationship extends to the other domain, we constructed a BisoNet graph (See Figure-4). The graph showed that there is a strong relationship between Magnesium and Migraine, i.e. the orange link between migraine and magnesium indicates a very strong relationship. Terms like

serum, calcium, potassium, phosphorus and effect(s) on the magnesium domain have very strong relationships and this may infer that there is a high probability of having a serendipitous recommendation through those strong relationships. Titles containing these terms were later recommended from the magnesium domain to the migraine domain.

## 5.  CONCLUSIONS AND RECOMMENDATIONS

This research attempts to address the overspecialisation problem in research paper recommender systems. Constantly receiving recommendations of articles that are similar to one's profile, or what a user has been working on, always leads to uninteresting and obvious recommendations. We envisioned an ideal system to be one that can recommend relevant, unexpected and novel research papers from domains that are least expected. We utilised bisociative knowledge discovery concepts to mine for new knowledge and information. Word clouds at the beginning of our experiments indicated that the domains were unrelated, however, BisoNets were created to visualise various statistical significance measures that existed between various terms. These measures proved that various terms and concepts in the two domains were highly related. It was further shown that serendipitous recommendations were able to be made through exploiting latent

relationships that existed between the two domains. These hidden and underlying relationships led to the discovery of new information that could be of use by a research paper recommender system. All these concepts when put together revealed the possibility of using technology in recommending serendipitous research papers to researchers who are currently engaged in multi-disciplinary research. Further research is being conducted to find bridging terms within and between outliers of distinct domains. These research findings will benefit multi-disciplinary researchers and entities seeking relevant information from seemingly unrelated domains.

**Figure-4: BisoNet generated from statistical significance values from all the terms > 10.**

## REFERENCES

Adamopoulos, P. (2013). *Beyond rating prediction accuracy: on new perspectives in recommender systems.* Paper presented at the Proceedings of the 7th ACM conference on Recommender systems.

André, P., Teevan, J., & Dumais, S. T. (2009). *From x-rays to silly putty via Uranus: serendipity and its role in web search.* Paper presented at the Proceedings of the SIGCHI Conference on Human Factors in Computing Systems.

Beel, J., Langer, S., Genzmehr, M., & Nürnberger, A. (2013). *Introducing Docear's research paper recommender system.* Paper presented at the Proceedings of the 13th ACM/IEEE-CS joint conference on Digital libraries.

Benlamri, R., & Zhang, X. (2014). Context-aware recommender for mobile learners. *Human-centric Computing and Information Sciences, 4*(1), 1-34.

Berthold, M. R. (2012). Towards bisociative knowledge discovery. In R. B. Michael (Ed.), *Bisociative Knowledge Discovery* (pp. 1-10): Springer-Verlag.

Bogers, T., & Van den Bosch, A. (2008). *Recommending scientific articles using citeulike.* Paper presented at the Proceedings of the 2008 ACM conference on Recommender systems.

Bollacker, K. D., Lawrence, S., & Giles, C. L. (1998). *CiteSeer: An autonomous web agent for automatic retrieval and identification of interesting publications.* Paper presented at the Proceedings of the second international conference on Autonomous agents.

Bollacker, K. D., Lawrence, S., & Giles, C. L. (2000). Discovering relevant scientific literature on the web. *IEEE Intelligent Systems and their Applications, 15*(2), 42-47.

Ferrara, F., Pudota, N., & Tasso, C. (2011). A keyphrase-based paper recommender system *Digital Libraries and Archives* (pp. 14-25): Springer.

Gipp, B., Beel, J., & Hentschel, C. (2009). *Scienstein: A research paper recommender system.* Paper presented at the Proceedings of the international conference on Emerging trends in computing (ICETiC'09).

Juršič, M., Sluban, B., Cestnik, B., Grčar, M., & Lavrač, N. (2012). Bridging concept identification for constructing information networks from text documents *Bisociative Knowledge Discovery* (pp. 66-90): Springer.

Kotkov, D., Wang, S., & Veijalainen, J. (2016). A survey of serendipity in recommender systems. *Knowledge-Based Systems, 111*, 180-192.

Kötter, T., Thiel, K., & Berthold, M. R. (2010). *Domain bridging associations support creativity*.

Lémdani, R., Polaillon, G., Bennacer, N., & Bourda, Y. (2011). *A semantic similarity measure for recommender systems.* Paper presented at the Proceedings of the 7th International Conference on Semantic Systems.

McNee, S. M., Kapoor, N., & Konstan, J. A. (2006). *Don't look stupid: avoiding pitfalls when recommending research papers.* Paper presented at the Proceedings of the 2006 20th anniversary conference on Computer supported cooperative work.

Nascimento, C., Laender, A. H., da Silva, A. S., & Gonçalves, M. A. (2011). *A source independent framework for research paper recommendation.* Paper presented at the Proceedings of the 11th annual international ACM/IEEE joint conference on Digital libraries.

Nodus Labs. (2015). Divinatory Recommender Systems: between Similarity and Serendipity.   Retrieved from https://noduslabs.com/research/divinatory-recommender-systems-similarity-serendipity/

Paraschiv, I. C., Dascalu, M., Dessus, P., Trausan-Matu, S., & McNamara, D. S. (2016). A Paper Recommendation System with ReaderBench: The Graphical Visualization of Semantically Related Papers and Concepts *State-of-the-Art and Future Directions of Smart Learning* (pp. 445-451): Springer.

Shani, G., & Gunawardana, A. (2011). Evaluating recommendation systems *Recommender systems handbook* (pp. 257-297): Springer.

Sridharan, S. (2014). Introducing serendipity in recommender systems through collaborative methods.

Swanson, D. R. (1988). Migraine and magnesium: eleven neglected connections. *Perspectives in biology and medicine, 31*(4), 526-557.

Vellino, A., & Zeber, D. (2007). *A hybrid, multi-dimensional recommender for journal articles in a scientific digital library.* Paper presented at the Proceedings of the 2007 IEEE/WIC/ACM international conference on web intelligence and international conference on intelligent agent technology.

Wang, G., He, X., & Ishuga, C. I. (2018). HAR-SI: A novel hybrid article recommendation approach integrating with social information in scientific social network. *Knowledge-Based Systems, 148*, 85-99. doi:https://doi.org/10.1016/j.knosys.2018.02.024

Zarrinkalam, F., & Kahani, M. (2012). *A multi-criteria hybrid citation recommendation system based on linked data.* Paper presented at the

Computer and Knowledge Engineering (ICCKE), 2012 2nd International eConference on.

Zhao, W., Wu, R., Dai, W., & Dai, Y. (2015). *Research Paper Recommendation Based on the Knowledge Gap.* Paper presented at the 2015 IEEE International Conference on Data Mining Workshop (ICDMW).