
F-Droid Application: Arxiv Papers

Yiqiao Yin
W.Y.N. Associates, LLC

Abstract

1 This is the capstone project of this course. This report examines an
2 F-Droid application called ArXiv Explorer that search, archive, and
3 download papers at ease and with a move of finger. The report has
4 two parts. First, the report investigates the ArXiv Explorer applica-
5 tion. Next, the report discusses the procedure of design and launch
6 of this project. The report also explores interesting facts that can
7 be provided to the readers including benefits and area of improve-
8 ments. Last, the report provides several novel solutions to tackle the
9 areas for improvements. For the second part, the project has also
10 finished a complimentary presentation slides where a video recording
11 is prepared.

12 **1 Introduction**

13 This study reviewed more than 200 research papers and investigated the research-paper
14 recommendation system (Beel et al., 2016). Amongst these 200 studies, approximately
15 55% of the papers are content-based filtering, 18% of the papers are collaborative
16 filtering, and 16% of the papers are graph-based recommendations. There are other
17 recommendation concepts and they proposed stereotyping, item-centric recommender,
18 and hybrid recommender systems. Based on citation contexts, many evaluations have
19 attempted to develop tasks and models to be used on paper recommender systems Teufel
20 et al. (2006); Ghosh et al. (2016); He et al. (2010); Ebesu and Fang (2017); Mohammad
21 et al. (2009); Chen and Zhuge (2019)

22 **2 Examination of ArXiv Explorer Repository**

23 This section of the report focus on the examination of ArXiv Explorer Repository. The
24 section consists of the following subsections.

- 25 • “README”: This subsection summarizes the benefits of a README file for
26 the repository of ArXiv Explorer.
27
- 28 • “Features”: This subsection covers the important features of the package at
29 interest.
30

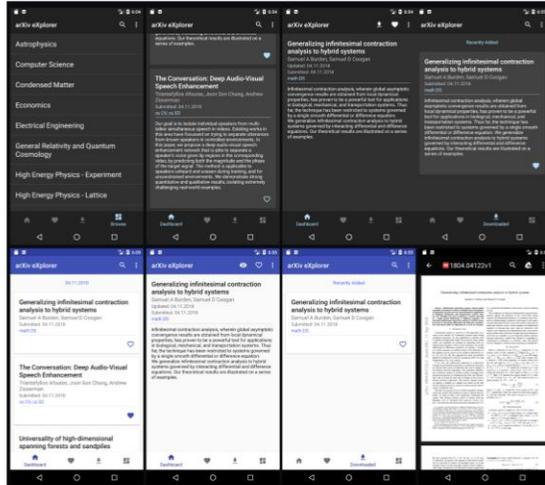


Figure 1: Screenshots of ArXiv Explorer Application UI.

- 31 • “Download Instructions”: This subsection introduces the instructions of
32 downloading such package and some user manual of navigating the package.
33
- 34 • “Repository Management”: This subsection covers the important pipeline of
35 the management of data-flow of this repository.
36

37 2.1 README

38 The repository of ArXiv Explorer has an instruction manual to lead readers through
39 steps of installation and deploying the application. This instruction manual is known as
40 a “READ ME” file. This READ ME file covers a list of features used in the application.
41 In addition, it also provides links to deploy on Google Play (mobile application store of
42 android devices). The F-Droid link is also made available for collaboration purposes.

43 The most useful information and import features from the User Interface (UI) is presented
44 in Figure 1. These are the screenshots of the application UI. The application has
45 navigation menu in the bottom of the screen. For simplicity, the application has four
46 buttons. The “Home” button is placed first to navigate to the main dashboard. The
47 “Favorite” button (the second one that looks like a heart shape) is a location when the user
48 saved articles that are to be read later. The “Download” button (the third one that looks
49 like an arrow pointing down) is a function key to allow offline reading. The “Browse”
50 button (the fourth one that looks like a grid of squares) is to view a list of articles in
51 ArXiv Explorer.

52 2.2 Features

53 The entire main function is written in java. The most fundamental feature of the main
54 function is the “browse” functionality. To establish the connection, an ArXiv Application
55 Programming Interface (API) is required instantiate the connection from local server to
56 the cloud server of the ArXiv website. The “browse” function has an access key which is

57 used to extract the information of the paper. There is a “getPapers” function coded inside
58 the “browse” script and related attributes are also built in such as sort by order, sort by
59 group, and maximum result to ensure the UI is comfortable to the end-users. There is
60 also a “callback” function in the “browse” script to respond to a user interaction. Upon a
61 click of the button, a query request is sent and the information of the paper is extracted.

62 The next fundamental feature of the main function is “dashboard”. The “dashboard”
63 functionality is wrapped around the “browse” function, because the “dashboard” will
64 generate a list of papers of which the information is extracted using “browse” functions.
65 To ensure the query list is comfortable and easy to read, there is a “toggleCategory-
66 Names” function written in the script. This implementation is to ensure categories and
67 subcategories are up to date and limited to certain length so that readers do not have to
68 scroll through pages of papers before they see the next section. In this script, information
69 of the paper can be concatenated and added to the list for review.

70 The “category” is introduced above as a function. However, the behind-the-scene action
71 of “category” is actually written as an independent script. The script starts off with a
72 “viewHolder” where information of papers can be hosted and stored as a placeholder.
73 Then functionalities such as “binding”, “recycling”, “rowCount” are implemented to
74 allow further categorization of list of queries. Afterwards, there is also a “category
75 presenter” script.

76 An attribute built upon “category” is the “category presenter” feature. This feature
77 allow the article in designated category to be reviewed upon a click motion, and it is
78 programmed as an reactive feature with the backend.

79 A download attribute is also made available for the UI. This feature allows users to read
80 offline and without the access of the internet after the file is successfully downloaded.
81 The feature currently is not supported with particular VPN setting.

82 The main model feature in the script is a function called “model”. The model is defined
83 using a “Paper” function where the function is a public object defined globally in
84 the application. This object carries attributes such as title, author, summary, publish
85 date, paper ID, URL, category, and so on. The object, once defined and fed into the
86 information, should present the end-user a full list of information from its attributes. The
87 API is presented in the following.

```
88  
89 public Paper(String title,  
90             String author,  
91             String summary,  
92             String updatedDate,  
93             String publishedDate,  
94             String paperID,  
95             String pdfURL,  
96             String categories,  
97             String paperURL) {  
98     this.title = title;  
99     this.author = author;  
100    this.summary = summary;  
101    this.updatedDate = updatedDate;  
102    this.publishedDate = publishedDate;  
103    this.paperID = paperID;  
104    this.pdfURL = pdfURL;  
105    this.categories = categories;  
106    this.paperURL = paperURL;  
107    this.downloaded = false;
```

```

108     this.favorited = false;
109     this.read = false;
110 }
111
112 public String getTitle() {
113     return title;
114 }
115
116 public void setTitle(String title) {
117     this.title = title;
118 }
119
120 public String getPublishedDate() {
121     return publishedDate;
122 }
123
124 public String getPaperID() {
125     return paperID;
126 }
127
128 public String getSummary() {
129     return summary;
130 }
131
132 public String getPDFURL() {
133     return pdfURL;
134 }
135
136 public String getAuthor() {
137     return author;
138 }
139
140 public String getUpdatedDate() {
141     return updatedDate;
142 }
143
144 public void setDownloaded(boolean downloaded) {
145     this.downloaded = downloaded;
146 }
147
148 public void setFavorited(boolean favorited) {
149     this.favorited = favorited;
150 }
151
152 public void setRead(boolean read){this.read = read;}
153
154 public String getCategories() {
155     return categories;
156 }
157
158 public String getPaperURL() {
159     return paperURL;
160 }

```

162 The “paper” object is fed in the query of information from an API developed called
163 “network”. The network is an API that feeds access into the paper object. It requires
164 a sufficient URL as “baseUrl” which is a parameter in the ArxivAPI object. This
165 ArxivAPI is a class object in java and it is built to ensure the client information is
166 communicated without error. A high level brief of the code structure is provided in the
167 following.

```
168 public class ArxivAPI {  
169     private static final String baseUrl = "http://export.arxiv.org  
170         ↪ /api/";  
171     private static final String querySegment = "query";  
172     private static OkHttpClient client;  
173  
174     public static OkHttpClient getClient() {  
175         if (client == null) {  
176             client = new OkHttpClient();  
177         }  
178         return client;  
179     }  
180     ...  
181 }  
182 }  
183 }  
184 }
```

185 The “ArxivAPI” class object is supported by a builder function of which arguments
186 are parsed into the object to be able to carry on the attributes throughout the rest of
187 the script. There is also a private static object that is called “addANDParamToQuery”,
188 which allows the app to concatenate the query, key, and parameters together to form a
189 more complete query names.

190 There is a public search function that can allow users to enter search queries such as
191 name, category, author, title, and so on. This search query, upon request of the user, then
192 sends the information request into the app. The app can use the information to retrieve
193 the corresponding articles, assuming that there is valid internet connection. If there is
194 not a completely match of those key words, a close or a similar one will be provided.
195 Afterwards, the same “download” feature is introduced to allow user to download the
196 paper if desired.

197 Another important feature is the search algorithm. The current edition of the app provides
198 UI of the search function. However, the query is fed into ArXiv using URL and the list
199 of results are directly retrieved using this URL provided sufficient internet connection.
200 However, this is sub-optimal, because the app itself does not carry any search algorithm
201 and can be improved based on the users’ needs.

202 **2.3 Download Instructions**

203 The ArXiv Explorer application or package can be downloaded using the following code.

```
204 git clone https://github.com/GarrettBeatty/arXiv-eXplorer.git  
205 bundle install  
206 }
```

208 **2.4 Repository Management**

209 The ArXiv Explorer application is publically available on F-Droid website. In addition,
210 the source code of the application can be found on the Github with the following URL.

211 It is made very convenient for software developers. If desired, one can simply go to the
212 following URL to download the source code and install the package.

213 `https://github.com/GarrettBeatty/arXiv-eXplorer`
214

216 **3 Areas of Improvement**

217 This section discusses some areas of improvement.

218 **3.1 AI-based Solution**

219 Despite many of its accomplishments, it is still unclear which recommendation system
220 is in effect and which can lead to promising search results during literature review.

221 **3.2 Category Refinement**

222 **3.2.1 User Recommendation**

223 The category has attributes such as category presenter mode. With a click of a button,
224 the presenter mode present summary information such as the title, author, and abstract.
225 However, to further tailor the need of the software package to researcher or PhD student,
226 it is often times beneficial to have the paper linked to a certain author or topic.

227 This is an area where user recommendation system can be implemented heavily while the
228 current version of the package is not equipped. In research period of graduate program,
229 the literature review is heavily built upon and tailored by each individual research topic.
230 Due to this nature, the review process could be drastically different and varying student
231 by student. In this case, a conventional system without user recommender could easily
232 provide insufficient or unrelated articles that may be a huge waste of time, because it
233 is up to the end-user to screen and filter whether the article can be used. However, if
234 features about the article can be extracted (which according to the current design of the
235 package it does have), then it is a small leap of a way to introduce machine learning or
236 even AI-based solutions.

237 **3.2.2 Build Machine Learning Model for User Recommendation**

238 To further investigate this direction, let us denote the interest of an article to be the target
239 which is what the model should be producing as a first screening process. Denote this
240 target Y . The features such as author, topic, keywords, abstract, can then be extracted as
241 features, which can be denoted as X . The task is to learn from X and try to produce
242 Y . In other words, the desired task is a supervised learning model that attempts to
243 understand and potentially guess the real relationship between X and Y . In other words,
244 there exist in nature a map, $f(\cdot)$ such that $Y = f(X)$. The goal is to train a machine
245 to understand and make the best guess that it can about the map, hence, the goal is to
246 search for the best $\hat{f}(\cdot)$ such that $\hat{Y} = \hat{f}(X)$ with the least error possible. In this learning
247 model, the error is measured between Y and \hat{Y} where Y can be the real target in the
248 training data and \hat{Y} is the educated guess from the model based on the features. The
249 proposed idea can be illustrated in Figure 2.

250 In this proposal, it is important to write the application to create a user-centric training
251 data set. The idea is that the application is built to learn from each and every user. There
252 could be an environment built and to collect the data only about the user at interest. This



Figure 2: **User Recommendation System.** How does it work?

253 way the target variable Y is tailored to the interest and the passion of that user only,
 254 which can be much more effective at recommending articles.

255 The downside of this approach is the potential lack of the data in the beginning of
 256 the user experience, which is the most crucial component of the user period. During
 257 this time, if the user does not use the application often enough for the machine to be
 258 well trained, it is possible to create similar dataset based on the user's demographics
 259 information. For example, a form can be created to collect this information at the sign-up
 260 stage. Information such as name, school, education level, researcher status, and so on
 261 can be extracted to build the model.

262 The following models can be available candidates for this task. For simple linear
 263 regression, it takes the form of

$$Y = \beta_0 + \beta_1 X + \epsilon \quad (1)$$

264 where Y is the target which states the interest a user has on an article and X being a
 265 feature such as education level. Under general scenario, the model can have multiple
 266 features or explanatory variables, and hence can take the form of

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon = \beta_0 + \sum_{j=1}^p X_j + \epsilon \quad (2)$$

267 where the explanatory data set is a p dimensional feature. In practice, the assumption of
 268 continuous variable for Y rarely holds true. An alternative is to use logistic regression
 269 model. To transform into logistic regression model, a non-linear link function such as
 270 sigmoid is used. The linear model is fed into the link function and hence produce the
 271 educated guess. The logistic regression model takes the following form

$$Y = \frac{1}{1 + \exp(-(\beta_0 + \sum_{j=1}^p \beta_j X_j))} \quad (3)$$

272 where the sigmoid function $1/(1 + \exp(-z))$ maps the output onto a range from 0 and
 273 1. Hence, this is also called a probabilistic model. The logistic model is actually a very
 274 simple neural network with one neuron. Instead of using ordinary least square to search
 275 for the parameters, we can use optimization techniques to search for the parameters. This
 276 is the fundamental idea for neural networks. A simple neural network is demonstrated in

Figure 3: **Diagram of an Artificial Neural Network.** Artificial Neural Network or ANN has input layers, hidden layers, and output layers. These layers are fully connected. An ANN regressor has one neuron in the output layer and an ANN classifier has an output layer with a number of neurons the same with the number of classes.

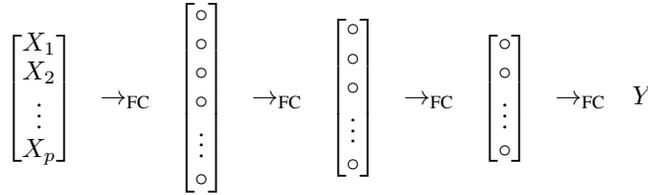
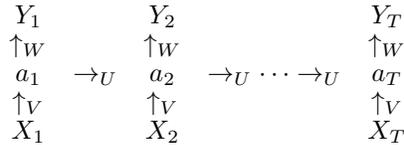


Figure 4: **Diagram of a Recurrent Neural Network.** A Recurrent Neural Network or RNN has layers designed as a Markov Chain. The weights are shared across different time stamps. Unlike ANN, RNN is not fully connected, because the time stamp is an important feature to capture in the network architecture.



277 the following diagram, Figure 3. The Artificial Neural Network or ANN has an input
 278 layer which are the explanatory features in the data, i.e. denoted as X_1, X_2, \dots, X_p . The
 279 input layer is sent into the ANN architecture with 3 hidden layers. Each hidden layer has
 280 certain number of neurons where the actual number is a tuning parameter. The layers are
 281 fully connected, denoted as an arrow with “FC”. This means that each neuron from the
 282 prior layer is connected with each of the neuron in the next layer. The final layer outputs
 283 a prediction which is the target.

284 The data serves as the key component of the research question in user recommendation
 285 system. In the beginning, it is discussed above that the system might be lack of data
 286 because the user is in the beginning stage of using the package. However, as time moves
 287 on there might be a change of situation where the user starts to use the application more
 288 dependent on their research demands. The alternative can also be true as well where a
 289 user might find the application not helpful and stops using it. In this situation, information
 290 changes across time and the features also change across time. It is somewhat beneficial
 291 to develop a Recurrent Neural Network to tackle and address some of the sequential
 292 pattern in the dataset. A Recurrent Neural Network (RNN), just like Artificial Neural
 293 Network (ANN), is constructed using forward propagation and backward propagation.
 294 However, the difference is that in RNN the features have time stamps, which is not the
 295 case for ANN. The diagram of a RNN is presented below.

296 3.2.3 Interpret Linear Regression Models

297 Another area to focus on is the interpretability of the machine learning component
 298 in category refinement. The purpose of category refinement is to be able to provide
 299 students and other end-users better and improved recommended articles in increase the

300 online traffic of the application. In addition to build more complicated models for better
 301 performance, the implication of models need to be analyzed as well. This functionality, if
 302 implemented well, will be able to give the organizers fundamental ideas what additional
 303 articles to recommend and why. In other words, a key component to introduce is the
 304 model interpretability and explainability of the machines used to predict category ratings
 305 of a user. In many cases, this can be treated as a model interpretability problem or even
 306 an inference problem.

307 In simple terms, a linear model is fundamentally interpretable. There is no additional
 308 work needs to be done to try to explain the model. A linear model such as one that
 309 follows ordinary least square (OLS) assumption searches for the linear coefficient of the
 310 model by using least square assumption. Denote the target to be Y and the educated
 311 guess to be \hat{Y} . The least square formula is written as

$$\sum_i (Y - \hat{Y})^2 \quad (4)$$

312 where the summation has a running index i to track the number of observations. If it
 313 is assumed that the model is a simple linear regression model, then it takes the form of
 314 equation 1. Under the OLS assumption, the least square is used (equation 4) to find the
 315 linear coefficient. The educated guess from the model is \hat{Y} , so the model formula can be
 316 used to replace it in order to update the least square formula.

$$\sum_i (Y - (\beta_0 + \beta_1 X))^2 \quad (5)$$

317 and the goal is to search for the linear coefficients β_0 and β_1 . This can be done by
 318 treating it as an optimization problem. At the optimal point, the gradient of least square
 319 is valued at 0. Hence, the set of equations can be set up as the following

$$\begin{aligned} \frac{\partial}{\partial \beta_0} \sum_i (Y - (\beta_0 + \beta_1 X))^2 &= 0 \\ \frac{\partial}{\partial \beta_1} \sum_i (Y - (\beta_0 + \beta_1 X))^2 &= 0 \end{aligned} \quad (6)$$

320 where the system is a two-equation-two-unknown. Hence, it can be solved and the OLS
 321 solution can be obtained

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}, \hat{\beta}_1 = \frac{\sum_i (X_i - \bar{X})(y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} \quad (7)$$

322 One important assumption for simple linear regression model is that the variance of the
 323 model is constant, i.e. σ^2 is constant. Though it is difficult to defend this assumption in
 324 practice, it does provide some convenient interpretation for us to understand the model
 325 once it is fitted. using σ^2 , it is mathematically sufficient to derive the following results

$$\text{SE}(\beta_0) = \sqrt{\frac{\sigma^2 \sum_i X_i^2}{n \sum_i (X_i - \bar{X})^2}}, \text{SE}(\beta_1) = \sqrt{\frac{\sigma^2}{\sum_i (X_i - \bar{X})^2}} \quad (8)$$

326 These estimators allow the end-users to make inference on the linear coefficients sought
 327 under the OLS assumption. For example, a simple model can be built to predict the
 328 housing price using square footage. The housing price is continuous and so is the
 329 square footage. One practical question useful for us is: “does square footage affect the
 330 housing price?” This question can be answered using linear regression model. A linear
 331 regression model can be built using square footage as X to explain the housing price Y .
 332 The housing price, under OLS assumption, is a linear additive form of the square footage.
 333 The linear coefficients along with their standard error are all known to the end-users

334 once the model is fitted. To answer the above question, a hypothesis test can be set up in
335 the following form

$$\begin{aligned} H_0 &: \text{Square footage has no impact on housing price} \\ H_1 &: \text{Square footage has a big impact on housing price} \end{aligned} \quad (9)$$

336 where the alternate hypothesis points to the obvious answer that end-user has in mind
337 and the null hypothesis is set up as the opposite end of the stick. To translate this into
338 a convention hypothesis testing problem with linear coefficients, the statement can be
339 rephrased in the following form

$$\begin{aligned} H_0 : \beta_1 = 0, & \quad \text{there is no coefficient,} \\ & \quad \text{hence the square footage should not be in the model,} \\ & \quad \text{so it's not important} \\ H_1 : \beta_1 \neq 0, & \quad \text{the coefficient should not be zero,} \\ & \quad \text{hence the square footage should be in the model,} \\ & \quad \text{so it's important} \end{aligned} \quad (10)$$

340 To execute the hypothesis test, it is important to find the test statistics which is computed
341 as the follows

$$t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\beta_1)} \quad (11)$$

342 and the statistics follow standard normal distribution under the null hypothesis. This
343 means that if the null hypothesis is correct the scientist should expect the test statistics
344 t to be small. Alternatively, if the test statistics is sufficient large, it gives the scientist
345 more confident to say that there is enough evidence to reject the null hypothesis. Hence,
346 it is possible to make inference which conclusion is more confident to be used. This
347 helps the end-user to establish basic intuition of whether the square footage is important.
348 It will be important if the test statistics is high. In addition, it is also available to end-user
349 the marginal impact the square footage has on housing price, because that information is
350 registered in the linear coefficient. The value of the estimator $\hat{\beta}_1$ indicates the numerical
351 changes on Y given that X changes by one unit. This is the crucial information to have
352 and it is fortunate by product coming with the model. Hence, it is theoretically sound to
353 say that linear regression models are interpretable. There is no additional work required
354 to help end-user to understand the internal structure after the model is fitted.

355 In addition, an example that is tailored to the ArXiv Explorer can be shown and discussed
356 in the following. Suppose the goal is to predict how much a user likes an article given the
357 features describing the article. Then the "interest" can be rated in 1-10 stars. In this case,
358 it can be treated as a continuous variable. Suppose the feature that the data scientist is
359 interested in is the length of the abstract. Then the length of the abstract can be X while
360 the rating of the article can be Y . A linear model can be built using $Y = \beta_0 + \beta_1 X$
361 and the linear coefficients can be computed. To answer the question: does length of the
362 abstract affect the user's rating of the article? It is convenient to compute the standard
363 error of the β_1 estimator and a hypothesis test can be established. The test statistics
364 $\hat{\beta}_1/\text{SE}(\beta_1)$ can be computed and then it can be compared with the z-score (the statistics
365 under standard normal distribution).

366 3.3 VPN for Download

367 The download attribute is not supported via VPN. This could be a problem where in
368 some countries the internet server may or may not have access to certain websites.

369 A diagram of the VPN working procedure is provided in Figure 5. The user device is
370 connected to the internet service provider. There are two scenarios provided in Figure

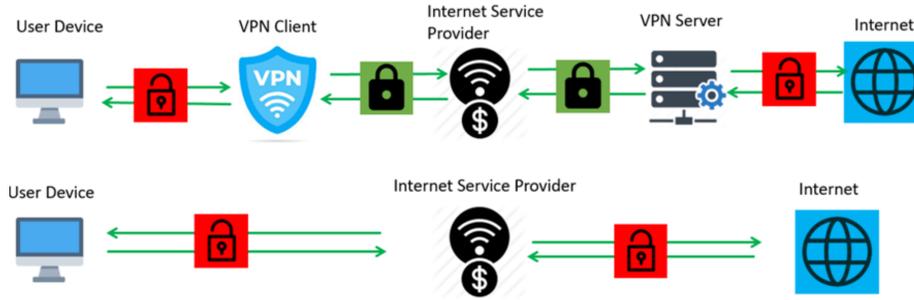


Figure 5: **VPN diagram.** How does it work?

371 5. The one in the bottom is the conventional scenario without the aid of VPN to secure
 372 the information. Then the information goes out of the internet service provider and is
 373 released to the internet. In this scenario, the user device directly communicates with the
 374 internet service provider and the information is not encrypted.

375 When VPN is provided, it serves as an additional client server to filter the information,
 376 of which is the information will be encrypted before it communicates in and out. Specif-
 377 ically, the VPN is set between the user device and the internet provider. The information
 378 coming out of the internet will be encrypted before it goes into the user device.

379 3.4 Search Algorithm

380 The current download function and the category function do not have any sorting or
 381 ranking system provided inside the application. This is an area of improvement, because
 382 the query produced from the Arxiv website may or may not provide desired order for
 383 the user. As the first-hand UI to interact with the human user, it is essential to learn
 384 and understand what the user needs and intends to read. This is where a potential
 385 PageRank-like algorithm can be born.

386 PageRank Algorithm, the founding algorithm for Google, is a linkage and sorting analysis
 387 algorithm that assigns a numerical weight score to each entry of a list of documents
 388 provided from the database. In this case, the database is the World Wide Web or WWW.

389 For simplicity, let us say there are only 3 documents in this small toy universe. Let us
 390 denote them to be A , B , and C . The initialization of the PageRank (denoted as \mathcal{R}) is
 391 equal weight. Hence, we have $\mathcal{R}(A) = \mathcal{R}(B) = \mathcal{R}(C) = 1/3$. Now suppose B and
 392 C has links to A , then each of them would transfer their value to A according to the
 393 following formula

$$\mathcal{R}(A) = \mathcal{R}(B) + \mathcal{R}(C) \quad (12)$$

394 and hence A would be valued at $1/3 + 1/3 = 2/3$.

395 In a separate situation, suppose the document B has a link to both C and A while the
 396 document C has a link to document A only. Then B has to allocate points to A and the
 397 same thing goes for C . However, to value the PageRank for A , the scores need to be
 398 divided by the number of outbound links \mathcal{L} . Thus, in this case, the PageRank for A is
 399 stated as the following

$$\mathcal{R}(A) = \frac{\mathcal{R}(B)}{\mathcal{L}(B)} + \frac{\mathcal{R}(C)}{\mathcal{L}(C)} = \frac{1}{2} + \frac{1}{1} = \frac{3}{2} \quad (13)$$

400 In general situation, we can formally write the PageRank score for any document d to be
 401

$$\mathcal{R}(d) = \sum_{d' \in \Pi} \frac{\mathcal{R}(d')}{\mathcal{L}(d')} \quad (14)$$

402 where Π is the document space and d' is a dummy variable to trace every document in
 403 the document space.

404 A sample PageRank algorithm in python is provided below.

```

405 """PageRank algorithm with explicit number of iterations.
406
407 Returns
408 -----
409 ranking of nodes (pages) in the adjacency matrix
410
411 """
412
413 import numpy as np
414
415 def pagerank(M, num_iterations: int = 100, d: float = 0.85):
416     """PageRank: The trillion dollar algorithm.
417
418     Parameters
419     -----
420     M : numpy array
421         adjacency matrix where M_i,j represents the link from 'j'
422         ↳ to 'i', such that for all 'j'
423         sum(i, M_i,j) = 1
424     num_iterations : int, optional
425         number of iterations, by default 100
426     d : float, optional
427         damping factor, by default 0.85
428
429     Returns
430     -----
431     numpy array
432         a vector of ranks such that v_i is the i-th rank from [0,
433         ↳ 1],
434         v sums to 1
435
436     """
437     N = M.shape[1]
438     v = np.random.rand(N, 1)
439     v = v / np.linalg.norm(v, 1)
440     M_hat = (d * M + (1 - d) / N)
441     for i in range(num_iterations):
442         v = M_hat @ v
443     return v
444
445 M = np.array([[0, 0, 0, 0, 1],
446              [0.5, 0, 0, 0, 0],
447              [0.5, 0, 0, 0, 0],
448              [0, 1, 0.5, 0, 0],
449              [0, 0, 0.5, 1, 0]])
450 v = pagerank(M, 100, 0.85)
451

```

453 **References**

- 454 Beel, J., Gipp, B., Langer, S., and Breitingner, C. (2016). Paper recommender systems: a
455 literature survey. *International Journal on Digital Libraries*, 17(4):305–338.
- 456 Chen, J. and Zhuge, H. (2019). Automatic generation of related work through summariz-
457 ing citations. *Concurrency and Computation: Practice and Experience*, 31(3):e4261.
- 458 Ebesu, T. and Fang, Y. (2017). Neural citation network for context-aware citation
459 recommendation. In *Proceedings of the 40th international ACM SIGIR conference on*
460 *research and development in information retrieval*, pages 1093–1096.
- 461 Ghosh, S., Das, D., and Chakraborty, T. (2016). Determining sentiment in citation text
462 and analyzing its impact on the proposed ranking index. In *International Confer-*
463 *ence on Intelligent Text Processing and Computational Linguistics*, pages 292–306.
464 Springer.
- 465 He, Q., Pei, J., Kifer, D., Mitra, P., and Giles, L. (2010). Context-aware citation
466 recommendation. In *Proceedings of the 19th international conference on World wide*
467 *web*, pages 421–430.
- 468 Mohammad, S., Dorr, B., Egan, M., Hassan, A., Muthukrishnan, P., Qazvinian, V., Radev,
469 D., and Zajic, D. (2009). Using citations to generate surveys of scientific paradigms.
470 In *Proceedings of human language technologies: The 2009 annual conference of*
471 *the North American chapter of the association for computational linguistics*, pages
472 584–592.
- 473 Teufel, S., Siddharthan, A., and Tidhar, D. (2006). Automatic classification of citation
474 function. In *Proceedings of the 2006 conference on empirical methods in natural*
475 *language processing*, pages 103–110.