# 𝔓akistan 𝔍ournal of 𝔖tatistics and 𝔒peration 𝔕esearch

# Bibliomining and Comparison of Q4 and ESCI WoS Indexed journals under "Statistics and Probability" Category

Nadeem Shafique Butt[1*]

* Corresponding Author

1. Department of Family and Community Medicine, Faculty of Medicine in Rabigh, King Abdulaziz University, Kingdom of Saudi Arabia. nshafique@kau.edu.sa

**Abstract**

The field of 'Statistics and Probability' has expanded its scope over the last few decades and have become an integral part of many fields with continuously increasing demand. This manuscript aimed for at a bibliometric analysis and comparison of all published documents during 2015 – 2019, from journals in the study topic category of 'Statistics and Probability' for Q4 Impact Factor (IF) journals and Emerging Source Citation Index (ESCI) of Web of Science (WoS). Sources with incomplete data for study timeframe were excluded and 31 sources from Q4 IF and 32 from ESCI journals were selected yielding 12808 and 4294 documents respectively. After data extraction from WoS, the bibliometric analysis at; source, author and document levels, were performed using "Bibliometrix" R-package. Q4-IF sources produced around 3 times more documents than ESCI sources. Articles were the main document type for both categories. China and USA were leading countries for Q4-IF while India, USA and Korea were dominant among ESCI documents. Two authors, namely, 'Cordeiro GM' and 'Alizadeh M' were among the 10 most productive authors in both categories. Sources "Communications in Statistics-Theory and Methods" and "Korean Journal of Applied Statistics" were leading contributors for Q4-IF and ESCI category respectively. For both categories, mainly similar trends were observed for keywords and topic coverage. In both Q4-IF and ESCI journals 'Maximum likelihood' and 'Ordered statistics' were observed to be most predominant keywords. A consistent publication trend with few similarities was observed in terms of documents production over the years for these two categories.

**Key Words:** Statistics and Probability, Bibliomining, Bibliometrics, Web of Science, ESCI and Q4 Journals

## 1. Introduction

Statistics and Probability has become an integral subject that offer approaches to deal with structure and give insights through data. Additionally, big data is establishing new challenges to researchers and mainly statisticians (Secchi, 2018). Furthermore, recent computational methods advancement has grasped wide attention from researchers and readers of several disciplines toward this subject. This field has expanded its scope over the last few decades and have become an integral part of many fields (Donoho, 2017; Drummond & Tom, 2011). This attribute has increasing demand and contribution for major areas of research such as Arts & Humanities, Life Sciences & Biomedicine, Physical Sciences, Social Sciences and Technology. Notably, despite known significance, the number of journals in this specific subject are assumed to be far less as compared to number of journals in other subjects. Nevertheless, there is limited literature that has research productivity and evolution of the 'Statistics and Probability' as a subject (Butt, Forthcoming 2021). However, some positive wave has been observed on the issue over the last few years.

Several databases are available and used by researchers and academicians. Yet, Web of Science (WoS), a Clarivate Analytics (Formerly Thomson Reuters) maintained platform, is considered as one of the most comprehensive and precise source for scientific exploration and appraisal with highest quality indexing. It is also assumed to be more appropriate to evaluate the research output of different regions, authors or organizations (Jelercic, Lingard, Spiegel,

Pichlhöfer, & Maier, 2010; Ronda-Pupo, Díaz-Contreras, Ronda-Velázquez, & Ronda-Pupo, 2015). It encompasses search across salient search databases, disciplines and document types along with more than one billion searchable cited references (WoS).

Although, for any journal it is not common to be included in one subject category as usually most journals show overlap in terms of their coverage context yet WoS has made certain specific subject categories and subsequently each published document inherits all subject categories given to the parent journal. 'Statistics and Probability' is one of such category in WoS with around 124 journals categorized in four quartiles (Q1 to Q4) according to their impact (MJL-WoS). Additionally, in late 2015, WoS launched 'Emerging Sources Citation Index' (ESCI), with more than 7000 journals covering scientific, social science, and humanities literature (ESCI-WoS, 2015). Journals indexed in the ESCI do not obtain Impact Factors. Though, Journal Citation Reports (JCR) citation counts includes ESCI citations and consequently contributing to other journals Impact Factors. Moreover, in the continuously growing, dynamic and diverse literature, ESCI provides WoS users with extended possibilities to explore emerging research areas.

Terms, 'Bibliomining' and 'Bibliometrics' are interchangeably used and provide a gateway to evaluate such proceedings and fill the knowledge gap (Abramo & D'Angelo, 2011; Shieh, 2010). In the field of statistics, various bibliometric studies have explored different aspects such as; citation patterns in the journals of statistics and probability (Stigler, 1994), communications between statistical methodology and applied statistics (Eto, 2000), most-cited statistical papers (Ryan & Woodall, 2005), decade of research in statistics (De Battisti, Ferrara, & Salini, 2015), statistical modelling of citation exchange between statistics journals (Varin, Cattelan, & Firth, 2016), and the importance of being clustered: uncluttering the trends of statistics (Anderlucci, Montanari, & Viroli, 2019).

Prominently, over the last few decades, WoS has been one of the widely used source for bibliometric analysis in various other scientific fields (Hossain, 2020; Merigó & Yang, 2017; Shukla, Muhuri, & Abraham, 2020; Yu & He, 2020; Yu, Xu, Pedrycz, & Wang, 2017). Foe ESCI, still no metrices and performance evaluation support is provided. Moreover, to the best of our knowledge, very limited literature has explored and compared the performance of statistics journals in any quartile and the ESCI category of WoS. Thus, this study aimed at a bibliometric analysis and comparison of all published documents during 2015 – 2019, from journals in the study topic category of 'Statistics and Probability' under categories of Q4-IF and Emerging Source Citation Index (ESCI) of Web of Science (WoS).

## 2. Methodology

All indexed sources from the WoS 'Statistics and Probability' classification under Q4 Impact Factor (IF) and Emerging Source Citation Index (ESCI) during 2015 – 2019 were selected (Figure 1). Later all identified journals (sources) were verified individually from the actual list provided by WoS in study category and added in "advanced search" through field tag: SO= Publication Name [Index]. For further analysis, we selected 32 out of 38 journals in the same category as shown in Fig 1. For further analysis, all 31 sources from Q4-IF and 32 out of 38 ESCI journals were selected, 6 titles were excluded from analysis because of incomplete coverage in WoS as shown in Figure 1.

'Cogent mathematics statistics' (newly added to ESCI with data available for only 2019), 'Journal of Japanese Society of Computational Statistics' and 'Journal of Modern Applied Statistical Methods' (no data for the year 2019), 'Modern Stochastics Theory and Applications' and 'Probability Uncertainty and Quantitative Risk' (no data for the year 2015) and 'Statistical Inference for Stochastic Processes' (no data for the year 2015 and 2016).

Data was extracted from WoS in plain text format and converted to bibliometric format and later bibliometric analysis at source level, author level and at document level were performed using R "Bibliometrix" package (Aria & Cuccurullo, 2017). At source and author level, impact was assessed by h-index and g-index. The h-index measure both impact of citations and publications productivity (Hirsch, 2005; Vílchez-Román, 2014). While the g-index measure  productivity based on the distribution of citations received by a researcher's publications (Egghe, 2006). Average Citation per document is an indicator showing citations per publication to quantify the author's impact, countries, and journals (Yi, Ao, & Ho, 2008). It is calculated as total number of citations received by total number of documents published by a journal to assess the yearly impact and provides fairer evaluation for author and journal activity (Harzing, 2010). Dendrogram was planned to evaluate keywords.
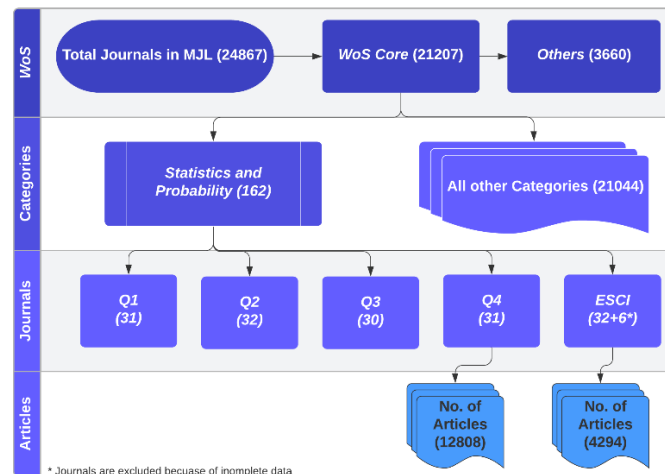
**Figure 1:** Flowchart of data extraction process from WoS (Dated: 19th March 2019)

## 3. Results

The search strategy yielded a total of 12808 documents from 31 Q4-IF journals and 4294 documents form 32 ESCI journals by 16121 and 6374 authors respectively. Similar collaborative index was observed for Q4-IF (1.41) and ESCI (1.61) journals. Higher average citations per documents 1.719 is observed in Q4-IF journals as compare to ESCI. Details of database bibliometrics characteristics given in Table 1.

**Table 1:** Bibliometric summary of extracted data (2015-2019)

| Description | Q4-IF | ESCI |
|---|---|---|
| **MAIN INFORMATION ABOUT DATA** | | |
| Sources (Journals) | 31 | 32 |
| Documents | 12808 | 4294 |
| Average citations per documents | 1.719 | 1.016 |
| Average citations per year per doc | 0.4047 | 0.2244 |
| References | 157207 | 69745 |
| Countries | 114 | 106 |
| **DOCUMENT TYPES** | | |
| Articles | 12401 (97%) | 4042 (94%) |
| Review | 143 | 14 |
| Correction | 68 | 10 |
| Editorial | 150 | 122 |
| Letter | 22 | 95 |
| Others | 24 | 11 |
| **DOCUMENT CONTENTS** | | |
| Keywords Plus (ID) | 9224 | 3464 |
| Author's Keywords (DE) | 31116 | 11606 |
| **AUTHORS** | | |
| Authors | 16121 | 6374 |
| Author Appearances | 30277 | 10193 |
| Authors of single-authored documents | 2106 | 725 |
| Authors of multi-authored documents | 14015 | 5649 |
| **AUTHORS COLLABORATION** | | |
| Single-authored documents | 2890 (23%) | 947 (22%) |
| Documents per Author | 0.794 | 0.674 |
| Authors per Document | 1.26 | 1.48 |
| Co-Authors per Documents | 2.36 | 2.37 |
| Collaboration Index | 1.41 | 1.69 |

**Table 2:** Top 10 most productive authors with impact

| Category | Author | No. of Publications | As Corresponding Author (%) | As First Author (%) | *h-index* | *g-index* | Total Citations |
|---|---|---|---|---|---|---|---|
| Q4 | Balakrishnan N | 90 | 17 (19) | 32 (36) | 7 | 9 | 239 |
| | Nadarajah S | 89 | 31 (35) | 15 (17) | 7 | 10 | 181 |
| | Cordeiro GM | 56 | 4 (7) | 14 (25) | 10 | 15 | 295 |
| | Aslam M | 39 | 31(79) | 14 (36) | 7 | 10 | 166 |
| | Qin H | 31 | 16 (52) | 3 (10) | 7 | 10 | 135 |
| | Alizadeh M | 29 | 10 (34) | 7 (24) | 8 | 10 | 153 |
| | Khoo MBC | 26 | 7 (27) | 1 (4) | 7 | 10 | 160 |
| | Kundu D | 25 | 15 (60) | 5 (20) | 7 | 11 | 143 |
| | Elsawah AM | 16 | 12 (75) | 14 (88) | 7 | 10 | 117 |
| | Vieu P | 11 | 1 (9) | 1 (9) | 8 | 10 | 120 |
| ESCI | Hamedani GG | 48 | 16 (33) | 15 (31) | 6 | 9 | 131 |
| | Cordeiro GM | 30 | 2 (7) | 5 (17) | 5 | 7 | 75 |
| | Alizadeh M | 27 | 6 (22) | 5 (19) | 6 | 9 | 106 |
| | Yousof HM | 25 | 8 (32) | 9 (36) | 6 | 10 | 122 |
| | Afify AZ | 19 | 5 (26) | 5 (26) | 6 | 11 | 131 |
| | Butt NS | 19 | 1 (5) | 1 (5) | 5 | 9 | 82 |
| | Gupta S | 17 | 5 (29) | 4 (24) | 4 | 5 | 31 |
| | Sabelfeld KK | 15 | 12 (80) | 10 (67) | 5 | 6 | 47 |
| | Al-omari AI | 14 | 13 (93) | 11 (79) | 4 | 6 | 48 |
| | Ghosh I | 13 | 6 (46) | 5 (38) | 5 | 6 | 50 |

For Q4-IF category, authors 'Balakrishnan N' and 'Nadarajah S' were most productive. Authors, 'Balakrishnan N' and 'Cordeiro GM' showed maximum citations. While for ESCI, authors 'Hamedani GG' and 'Cordeiro GM' were most productive. Authors, 'Afify AZ' and 'Cordeiro GM' showed maximum citations. Two authors, 'Cordeiro GM' and 'Alizadeh M' were among the 10 most productive authors in both categories. (Table 2).

**Table 3:** Top 5 Country appearances, Corresponding Author Country and local cited papers

| | Countries (Appearances) | Corresponding Author Countries | Highly Local Cited |
|---|---|---|---|
| Q4-IF | China (6590) | China (2358) | Torrado N, 2015, J Appl Probab (17) |
| | USA (6057) | USA (2148) | Wang Xj, 2015, Statistics-Abingdon (16) |
| | France (1675) | France (614) | Li C, 2015, Stat Probabil Lett (12) |
| | Iran (1666) | Iran (683) | Elsawah Am, 2015, J Stat Plan Infer (11) |
| | India (1625) | India (742) | Aly Aa, 2015, Commun Stat-Simul C (11) |
| ESCI | India (1086) | India (657) | Yousof HM, 2015, Pak J Stat Oper Res (21) |
| | USA (948) | USA (468) | Hamedani GG, 2018, Pak J Stat Oper Res-a (14) |
| | Korea (813) | Korea (573) | Hamedani GG, 2018, Pak J Stat Oper Res-a,b (14) |
| | Iran (415) | Iran (218) | Afify AZ, 2015, Pak J Stat Oper Res-a (13) |
| | France (388) | France (194) | Afify AZ, 2015, Pak J Stat Oper Res (12) |

Table 3 shows the top 5 country appearances, corresponding author country and local cited papers. China and India were found to have most authors appearances and corresponding authors.

In Q4-IF journals, China, Iran and USA showed the leading contributions. Ferdowsi univ mashhad showed strong collaboration between China and Iran. King Abdulaziz Univ showed strong linkage with china and few with other countries. Macmaster Univ showed strong collaboration of China and Canada while Nankai Univ showed greater collaboration with china and turkey.

In ESCI journals Korea, Egypt and India were leading contributors. Benha Univ showed strong linkages of researchers from Egypt, India and USA while King Abdulaziz university collaboration was observed with Egypt, India and Pakistan. Maximum likelihood estimation was observed most occurring author keyword in both Q4-IF and ESCI journals. Details given in Figure 2 Sankey Diagram.
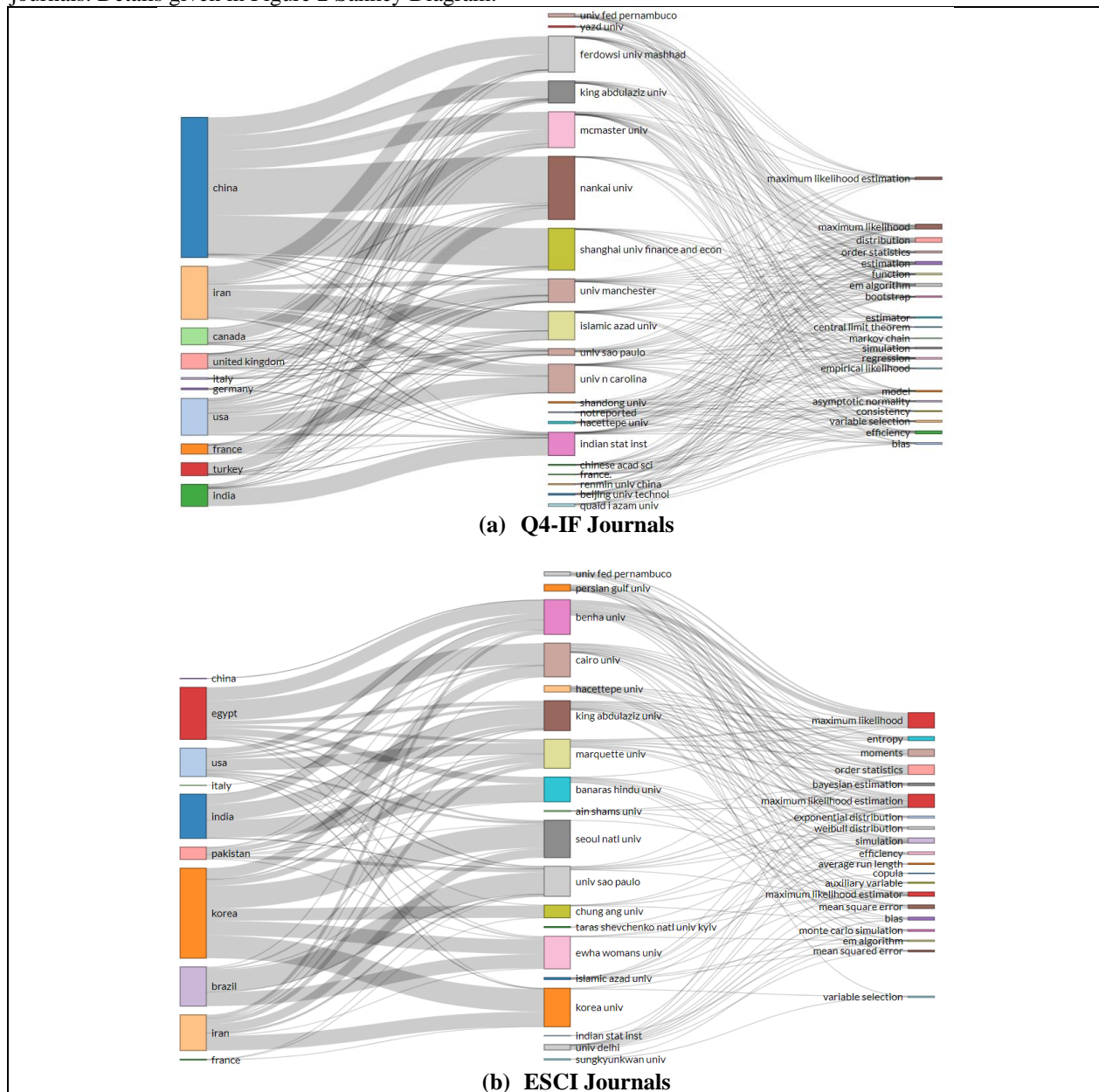


**(a)  Q4-IF Journals**



**(b)  ESCI Journals**

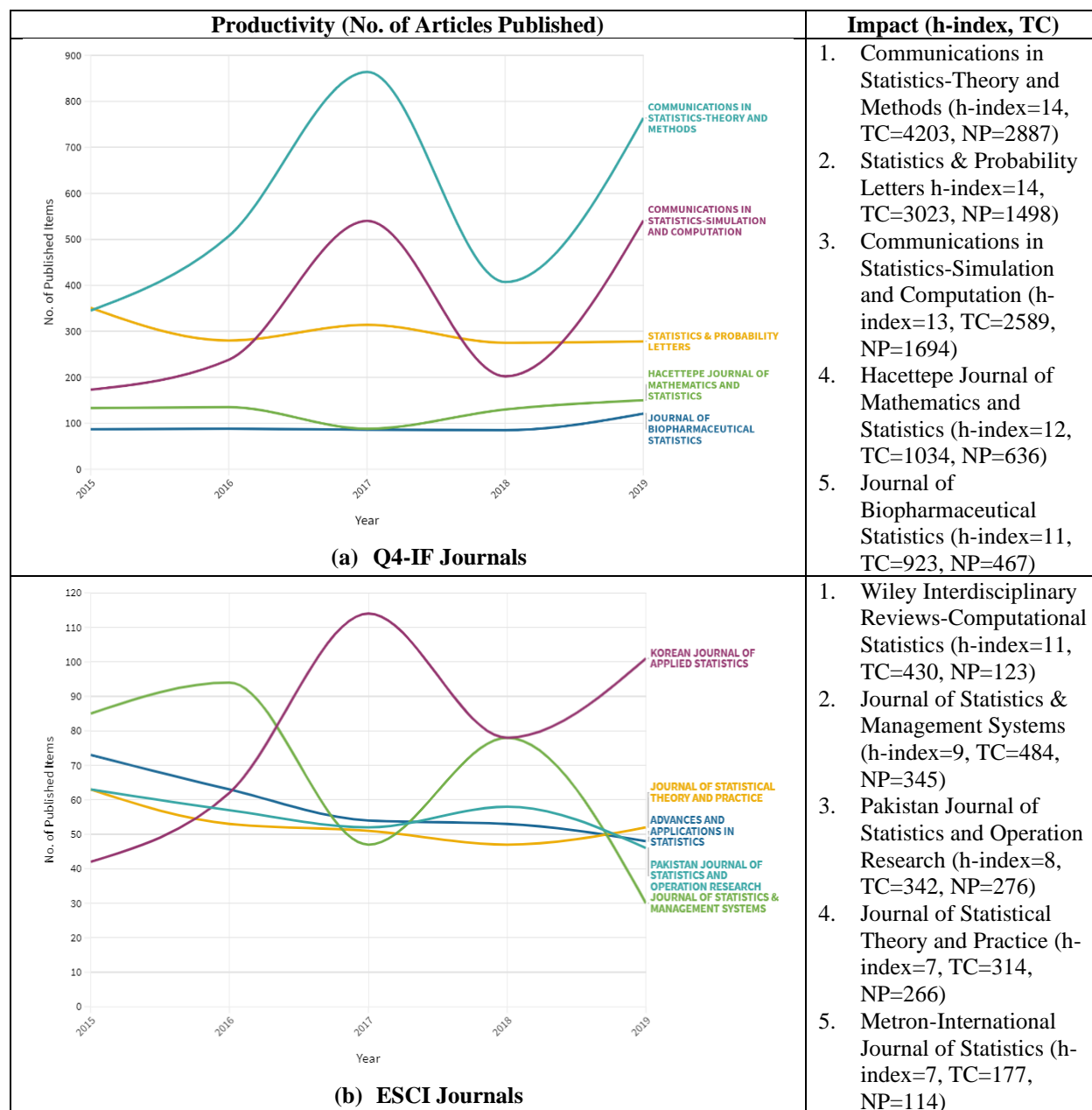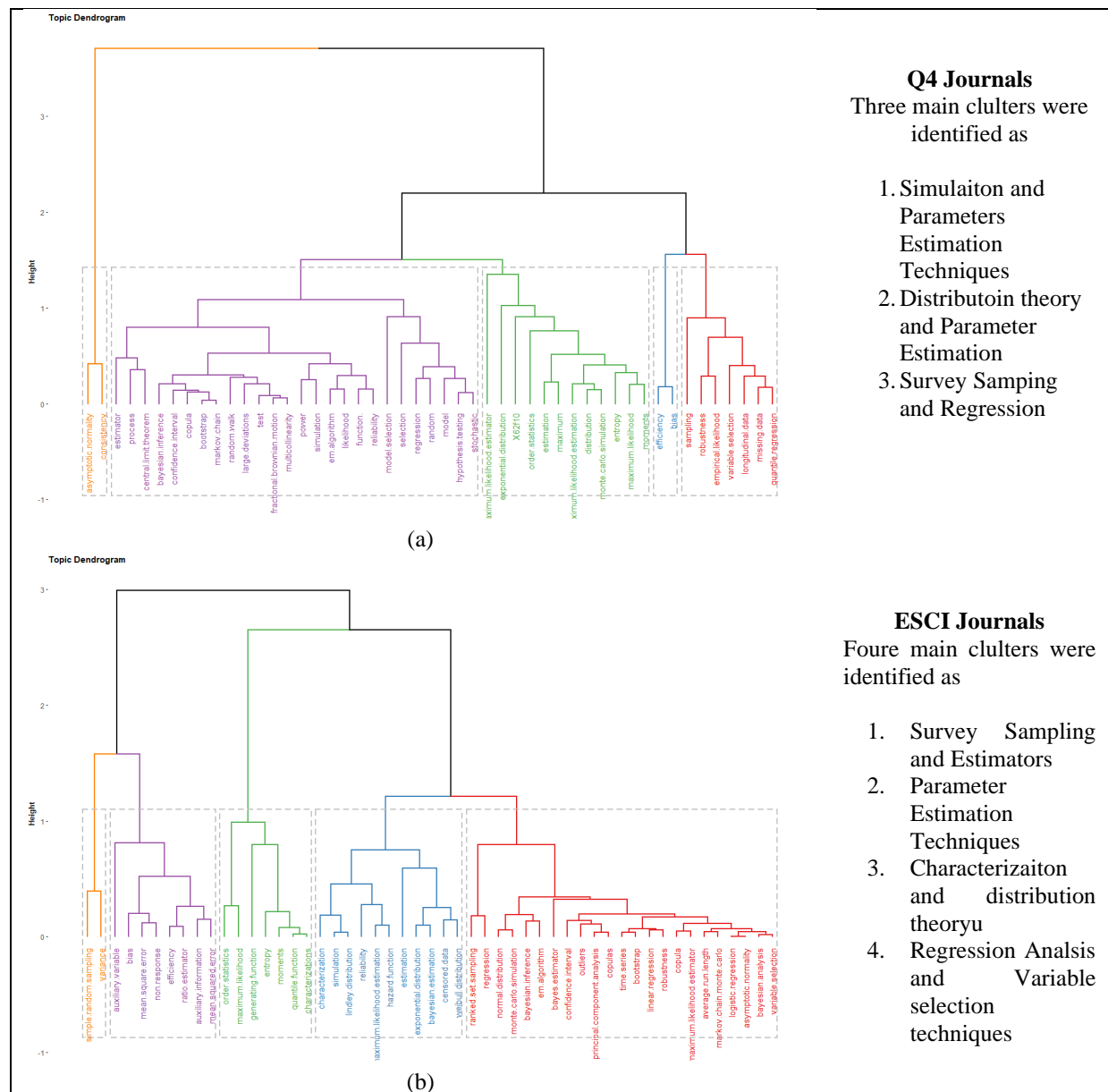**Figure 2:** Sankey Diagram based on top countries, institutes and keywords

**Figure 3:** Top 5 Sources trend by year and Impact

Among Q4 journals "Communications in Statistics-Theory and Methods" was the leading contributor as shown in the figure 3 (a). Same five journal also appeared in top 5 journals with respect to impact. Figure 3 (b) shows the top 5 ESCI journals where "Korean Journal of Applied Statistics" was found to be the most productive.

In Q4-IF Journals 'Distribution', 'Estimation', 'Asymptotic normality', 'Maximum likelihood' and 'Order statistics' were top trending author's keywords. While 'Maximum Likelihood Estimation', 'Order statistics', 'Simulation', 'Bias' and 'Maximum likelihood' were observed to be the most trending author keywords. Collectively in both Q4-IF and ESCI journals 'Maximum likelihood' and 'Ordered statistics' were observed to be most predominant keywords.

**Figure 4:** Author's Keywords dendrogram

For Q4-IF journals, three main clusters identified were 'Simulation and Patameter Estimation Techniques', 'Distribution theory and Parameter Estimation' and 'Survey Sampling and Regression'. While for ESCI journals, four main clusters were identified including; 'Survey Sampling and Estimators', 'Parameter Estimation Techniques', 'Characterizaiton and distribution theory', 'Regression Analysis and Variable selection techniques.

## 4. Discussion

This article shares the bibliometric analysis for all documents published in Q4-IF and ESCI journals in WoS category 'Statistics and probability' between 2015-2019. In general, a uniform trend in terms of numbers of publications each year with  majority as articles in both categories was observed. Almost equal number of journals were found in each category, but the number of articles/documents produced by 31 Q4-IF journals were nearly 3 times than that of 32 ESCI journals in same time frame. Similar trends were observed for most of the other variables including number of authors. Articles were the most common document type followed by editorials and single author documents were < 1/4th for both categories. Though, majority of the published documents were multi-author documents yet relatively higher number of single-authored documents, 23% and 22% was observed both in Q4-IF and ESCI journals.

For Q4-IF category, authors 'Balakrishnan N' and 'Nadarajah S' were most productive and authors, 'Balakrishnan N' and 'Cordeiro GM' showed maximum citations. While for ESCI, authors 'Hamedani GG' and 'Cordeiro GM' were most productive and authors, 'Afify AZ' and 'Cordeiro GM' showed maximum citations. Two authors, 'Cordeiro GM' and 'Alizadeh M' were among the 10 most productive authors in both categories. Interestingly, both of these prolific authors showed relatively less contributions as corresponding and/or first authors as compared to other leading authors.

In terms of authors' country appearances, China, USA, France, Iran and India were leading countries for Q4-IF documents while for ESCI documents, India, USA, Korea, Iran and France were leading. Mostly similar trends were observed for corresponding author countries for both Q4-IF and ESCI documents. The USA and India were common among top 5 contributor for both Q4-IF and ESCI categories. Among affiliations, Ferdowsi univ mashhad showed relatively more collaboration between China and Iran. King Abdulaziz Univ showed strong linkage with china and few with other countries. Macmaster Univ showed strong collaboration of China and Canada while Nankai Univ showed greater collaboration with china and turkey. While for ESCI journals Korea, Egypt and India were leading contributors. Benha Univ showed strong linkages of researchers from Egypt, India and USA while King Abdulaziz university collaboration was observed with Egypt, India and Pakistan. Maximum likelihood estimation was observed most occurring author keyword in both Q4-IF and ESCI journals. Study findings also suggest that though developed countries with affiliations were more among top contributors in Q4-IF than in ESCI category, yet many Asian countries and affiliations were dominant contributors in either category.

Among Q4-IF journals "Communications in Statistics-Theory and Methods" was the leading contributor throughout all 5 years with relatively higher h-index and total citations, followed by "Communications in Statistics-Simulation and Computation" and peak of publication frequency was observed in 2017 for both journals. However, "Journal of Biopharmaceutical Statistics", "Statistics & Probability Letters" and "Hacettepe Journal of Mathematics and Statistics" showed steady trend over time. Same five journal also appeared in top 5 journals with respect to impact. For ESCI journals, "Korean Journal of Applied Statistics" was found to be the most productive. While a sharp decline in publication frequency was observed for "Journal of Statistics & Management Systems". While, journals "Advances and Applications in Statistics", "Journal of Statistical Theory and Practice" and "Pakistan Journal of Statistics and Operation Research" showed steady trend over time. With respect to impact "Wiley Interdisciplinary Reviews-Computational Statistics" and "Metron-International Journal of Statistics" were shown to be predominant.

Mainly three and four clusters of keywords/topics were found in Q4-IF and ESCI categories with similarities and overlap. Diverse but mainly similar trends of keywords and topic coverage were observed for both Q4-IF and ESCI journals while 'Maximum likelihood' and 'Ordered statistics' were observed to be most predominant topic collectively.
Although limited available literature and data for comparison was a limitation yet it also suggests for further and continuous exploration of trends and relevant analysis.

In terms of other limitations, analysis was conducted only on WoS- Q4-IF and ESCI journals in the "Statistics & Probability" category with limited timeframe of 2015-2019 that may limit the generalizability of finding to the category in general. Secondly, limitations in WoS database may have some unidentified issues, however the findings shared here for the leading contributors were manually verified. Additionally, continuous changes and updates may show different publications data to be analyzed depending upon date of search and timeframe. Metadata from other sources might be beneficial to complement this study and provide comprehensive context on the subject.

**Conclusion:**
Considering scarcity of literature on 'Statistics & Probability' publication trends, despite its significance in research and academics, this paper assists to fill the gap by providing overview and salient trends in WoS-Q4-IF and ESCI "Statistics & Probability" category (2015-2019). A consistent publication trend was observed in terms of documents production but Q4-IF productivity was relatively much higher. Articles were the major type of document for both. Among prolific authors, only two were common between Q4-IF and ESCI categories. Overall, 114 countries contributed to the selected Q4-IF journals led by China and USA while India, USA and Korea were leading for ESCI journals among 106 countries. Mainly, similar trends were observed for impact and contributions as corresponding and/or first authors. Diverse but mainly similar trends of keywords and topic coverage were observed for both

categories. Although limited available literature and data for comparison was a limitation yet it also suggests for further and continuous exploration of trends and relevant analysis. In conclusion, the bibliometric findings of this study can benefit relevant stakeholders and particularly researchers to better understand the performance and trends of study subject and plan with better informed decisions with the help of these findings.

**Conflict of Interest:**
Authors declare no conflict of interest.

## References

1. Abramo, G., & D'Angelo, C. A. (2011). Evaluating research: from informed peer review to bibliometrics. *Scientometrics, 87*(3), 499-514.
2. Anderlucci, L., Montanari, A., & Viroli, C. (2019). The Importance of Being Clustered: Uncluttering the Trends of Statistics from 1970 to 2015. *Statistical Science, 34*(2), 280-300.
3. Aria, M., & Cuccurullo, C. J. J. o. i. (2017). bibliometrix: An R-tool for comprehensive science mapping analysis. *11*(4), 959-975.
4. Butt, N. S., Malik, A. A. & Shahbaz, M. Q.(2021). Bibliometric Analysis of Statistics Journals Indexed in Web of Science under Emerging Source Citation Index *SAGE Open, 11(1)* 1-8.
5. De Battisti, F., Ferrara, A., & Salini, S. (2015). A decade of research in statistics: A topic model approach. *Scientometrics, 103*(2), 413-433.
6. Donoho, D. (2017). 50 years of data science. *Journal of Computational and Graphical Statistics, 26*(4), 745-766.
7. Drummond, G. B., & Tom, B. D. (2011). Statistics, probability, significance, likelihood: words mean what we define them to mean. *Advances in physiology education, 35*(4), 361-364.
8. Egghe, L. (2006). Theory and practise of the g-index. *Scientometrics, 69*(1), 131-152.
9. ESCI-WoS. (2015). Emerging Sources Citation Index (ESCI), ISI Web of Knowledge by Clarivate Analytics (formerly known as Thomson Reuters). https://clarivate.com/webofsciencegroup/solutions/webofscience-esci/
10. Eto, H. (2000). Bibliometric distance between methodology and application in statistics. *Scientometrics, 48*(1), 85-97.
11. Harzing, A.-W. (2010). *The publish or perish book*: Tarma Software Research Pty Limited.
12. Hirsch, J. E. (2005). An index to quantify an individual's scientific research output. *Proceedings of the National academy of Sciences, 102*(46), 16569-16572.
13. Hossain, M. M. (2020). Current status of global research on novel coronavirus disease (Covid-19): A bibliometric analysis and knowledge mapping. *Available at SSRN 3547824*.
14. Jelercic, S., Lingard, H., Spiegel, W., Pichlhöfer, O., & Maier, M. (2010). Assessment of publication output in the field of general practice and family medicine and by general practitioners and general practice institutions. *Family practice, 27*(5), 582-589.
15. Merigó, J. M., & Yang, J.-B. (2017). A bibliometric analysis of operations research and management science. *Omega, 73*, 37-48.
16. MJL-WoS. Master Journal List (MJL), Web of Science Group.  Retrieved 11.02.2020 https://mjl.clarivate.com/search-results
17. Ronda-Pupo, G. A., Díaz-Contreras, C., Ronda-Velázquez, G., & Ronda-Pupo, J. C. (2015). The role of academic collaboration in the impact of Latin-American research on management. *Scientometrics, 102*(2), 1435-1454.
18. Ryan, T. P., & Woodall, W. H. (2005). The most-cited statistical papers. *Journal of Applied Statistics, 32*(5), 461-474.
19. Secchi, P. (2018). On the role of statistics in the era of big data: A call for a debate. *Statistics & Probability Letters, 136*, 10-14.
20. Shieh, J. C. (2010). The integration system for librarians' bibliomining. *The Electronic Library*.

21. Shukla, A. K., Muhuri, P. K., & Abraham, A. (2020). A bibliometric analysis and cutting-edge overview on fuzzy techniques in Big Data. *Engineering Applications of Artificial Intelligence, 92*, 103625.
22. Stigler, S. M. (1994). Citation patterns in the journals of statistics and probability. *Statistical Science*, 94-108.
23. Varin, C., Cattelan, M., & Firth, D. (2016). Statistical modelling of citation exchange between statistics journals. *Journal of the Royal Statistical Society. Series A,(Statistics in Society), 179*(1), 1.
24. Vílchez-Román, C. (2014). Bibliometric factors associated with h-index of Peruvian researchers with publications indexed on Web of Science and Scopus databases. *TransInformação, 26*(2), 143-154.
25. WoS. Clarivate Analytics (Formerly Thomson Reuters), Web of Science. Retrieved from https://clarivate.com/webofsciencegroup/solutions/web-of-science/
26. Yi, H., Ao, X., & Ho, Y.-S. (2008). Use of citation per publication as an indicator to evaluate pentachlorophenol research. *Scientometrics, 75*(1), 67-80.
27. Yu, D., & He, X. (2020). A bibliometric study for DEA applied to energy efficiency: Trends and future challenges. *Applied Energy, 268*, 115048.
28. Yu, D., Xu, Z., Pedrycz, W., & Wang, W. (2017). Information Sciences 1968–2016: a retrospective analysis with text mining and bibliometric. *Information Sciences, 418*, 619-634.