

Caste in the News: A Computational Analysis of Indian Newspapers

António Filipe Fonseca¹ , Sohhom Bandyopadhyay²,
Jorge Louçã¹, and Jaison Manjaly²

Social Media + Society
October–December 2019: 1–7
© The Author(s) 2019
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/2056305119896057
journals.sagepub.com/home/sms


Abstract

Conflicts involving caste issues, mainly concerning the lowest caste rights, pervade modern Indian society. Caste affiliation, being rigorously enforced by the society, is an official contemporary reality. Although caste identity is a major social discrimination, it also serves as a necessary condition for affirmative action like reservation policy. In this article, we perform an original and rigorous analysis of the discourse involving the theme “caste” in India newspapers. To this purpose, we have implemented a computational analysis over a big dataset of the 2016 and 2017 editions of three major Indian newspapers to determine the most salient themes associated with “caste” in the news. We have used an original mix of state-of-the-art algorithms, including those based on statistical distributions and two-layer neural networks, to detect the relevant topics in the news and characterize their linguistic context. We concluded that there is an excessive association between lower castes, victimization, and social unrest in the news that does not adequately cover the reports on other aspects of their life and personal identity, thus reinforcing conflict, while attenuating the vocality and agency of a large section of the population. From our conclusion, we propose a positive discrimination policy in the newsroom.

Keywords

content analysis, computational social science, natural language processing, caste, news media

Introduction

Caste is a form of social stratification characterized by endogamy, hereditary transmission of a lifestyle, which often includes an occupation, status in a hierarchy, customary social interaction, and exclusion (Scott & Marshall, 2009). The caste social system in modern India is a paradigmatic ethnographic example of such a social stratification. Norms and conventions that name, group and sometimes rank people can be traced far back into distant Indian past, “varnas,”¹ and later “jatis,”² originated in the later Vedic society (circa 1500–500 BCE). Over the centuries, social discrimination has contributed toward the hardening of caste identities. Social unrest in India during colonialism prompted the British to start positive discrimination by reserving a certain percentage of government jobs for the lower castes. After India achieved independence, the policy of caste-based reservation of jobs was formalized with lists of Scheduled Castes and Scheduled Tribes. The government officially recognized historically discriminated communities of India such as the Untouchables and certain economically backward castes as Other Backward Class. These policies of affirmative action aimed at reducing the inequality but paradoxically have also created an incentive to keep this

stratification alive. Over the time lower castes remained locked in unskilled, low-paying occupations.

Many important questions remain to be answered by the social scientist: Why do caste discrimination remain active in society, despite the Constitution and all the legal and political efforts? Why, despite recent economic growth, is caste conflict also rising (Deshpande, 2011; Thorat & Neuman, 2012)? How does caste divisive power articulate with nationalism? How will strict caste endogamy evolve in a modern liberal and open society? How are media contributing to propagate such ideas?

Including quantitative social scientists in our multi-disciplinary team, with the ability to analyze big quantities of data, our effort was directed to provide a better understanding of the expression of caste discrimination and conflict in the public sphere, more specifically in the news media,

¹ISCTE—Lisbon University Institute, Portugal

²IIT Gandhinagar, India

Corresponding Author:

António Filipe Fonseca, ISTAR, ISCTE-Lisbon University Institute, Edifício ISCTE-IUL (Edifício II), Sala D0.10, Av. das Forças Armadas, 1649-026 Lisboa, Portugal.
Email: ajffa@iscte-iul.pt



using a quantitative approach. To this purpose, our main task was to extract patterns of word usage and expression from newspaper texts and then to interpret the coexistence of those patterns. Our goal in this analysis will be to clarify objectively the nature of the written discourse that is presented to the general public, taking into account the lack of vocality of the lowest and underprivileged castes in India. We will try to assert with a relevant degree of certitude the relations between caste and signs of autonomy and agency of SC/ST communities as they are described in the news.

In the next section, we will present an overview of how the literature considers the “caste” question in India, focusing on how castes are represented in the media. Our quantitative study contributes for a better understanding of this specific matter. We start by discussing some of the relevant problems when studying newspaper texts and performing statistics over word usage. We will consider the advantages and the perils of supporting an analysis on this method. Then, the quantitative methods used in the context of our research are described, including an original mix of state-of-the-art algorithms that comprise procedures based on Latent Dirichlet allocation (LDA) distributed variables, and two-layer neural network Word2Vec. The results point to the relevant information that was obtained applying these methods. Finally, we will discuss the results obtained from this quantitative analysis and present the conclusions drawn against the already known readings of the “caste” problem.

Caste in India

Some surveys of nationally representative samples show that there has been convergence between the upper castes and the lower castes on education and occupations over the past decades in India (Hnatkowska & Lahiri, 2011). Caste plays an important role in economic mobility (Deshpande, 2011). On the one hand, by business and employment facilitation, caste networks can be very economically advantageous for lower disadvantaged castes (Damodaran, 2008). On the other hand, these same networks of acquaintances can restrict the economic mobility: they can, for example, enforce rural-urban wage gaps by keeping rural individuals receiving city migrant’s remittances (Munshi & Rosenzweig, 2016) or restricting labor market access to high caste occupations. Deshpande (2011) notices that overall, the evidence indicates that contours of caste disadvantage are changing in terms of some dimensions of relative distance between OBCs and upper castes. However, there is no indication of reversal of the broad historical caste hierarchies. Although they have gained substantially after the 90s economic reform, SC/ST individuals are still at the bottom of the income and social ladders (Thorat & Neuman, 2012).

Being legally penalized and prohibited, the maltreatment of Dalits in India (Article 17 of India Constitution abolished “untouchability” and forbade its practice in any form) has led some authors to classify caste discrimination as “India

hidden apartheid” (Human Rights Watch, 2007). Some criticize this view point citing substantial improvements in the position of Dalits in post-independence India, namely the implementation of the rights and privileges present in the Constitution of India, as also the Protection of Civil Rights Act of 1955 and the Prevention of Atrocities Act of 1989. These critics argue that the practice has disappeared in urban public life. Many other authors, however, have pointed out that although there is a sense in which caste has been banished from public discourse, this unspeakability only exists in civil society, it has not disappeared from society at large: away from the watchful gaze of the modern elite, caste is an important category that frames the common ways of seeing and being and living in India (Deshpande, 2011; Teltumbde, 2018; Thorat & Neuman, 2012). To the oppressed castes, this public repression appears as a conspiracy of upper castes to deprive them of their voice. One of the central elements to the recalcitrance of caste in contemporary Indian politics is the search for a past, a cultural legacy, a history and a sense of Self (Nigam, 2004).

According to Human Rights Watch, in 2016, although having taken some important strides with legal reforms with respect to the treatment of vulnerable populations, the Indian government failed to implement laws to protect Dalits and tribal groups from discrimination and violence. Despite a right to education law that mandated free and compulsory elementary education for all children, discrimination against children from Dalit, tribal and Muslim communities led to high dropout rates among these children (Sifton, 2016). Violence and conflict are frequent between castes. Violence against lower castes and outcastes is rendered banal by being woven into the fabric of everyday life, it is also conducted through spectacular acts. Dalits are raped and murdered for daring to aspire to land, electricity, drinking water and to non-Dalit partners (Loomba, 2016). The 2016 report of the National Crimes Record Bureau depicts crimes against SC/ST that are those registered under SC/ST Prevention of Atrocities Act (PoA Act). The report concerning 19 large cities with population over 2 million, not only rural areas where caste conflict are known as more common, shows that urban caste violence is not negligible (Deshpande, 2017). In many cases government officials usually are not held accountable and impunity persists for police and other security personnel who are shielded by laws (Sifton, 2016).

The Indian Constitution imposes equal treatment for every citizen irrespective of caste or creed, however the most important part of the lives of as much as one-quarter of the population belonging to the lower caste do not exist on the news. Very frequently, as audience demands it, news media report passionate episodes involving violence and conflict associated with caste and religious fights, eventually with some political tone and usually, they vocally signal the existence of social discrimination. However stories reporting the individuality, the normal life or any autonomous view of the lower 25% of the Indian population are unlikely to be known.

Much less is broadcast or written about. Dalits and lower caste have a problem of agency and self-expression within mainstream media. They may appear as violated or humiliated but their voice, will and hope, does not exist because newsrooms do not have SC/ST representatives.

Caste in the News

When the journalist B. N. Uniyal (1996) surveyed the scene in 1996, he found not a single Dalit accredited journalist in Delhi. In 2001, Professor Robin Jeffrey at La Trobe University wrote a strong article in which he reported that “Dalits are too vulnerable either to proclaim their ‘Dalitness’ to their newspapers when they do have jobs or to start newspapers of their own,” that “Readers will continue to receive newspapers which were the outcome of Dalits’ ‘not being there,’ and which reinforced ideologically the material weaknesses from which Dalits suffer. Newspapers tell readers that Dalits were poor, naive, ignorant; victims, and sometimes perpetrators, of violence; in short, people who almost never travel, eat or get married” (Jeffrey, 2001).

In 2006, a Centre for the Study of Developing Societies in New Delhi survey (Chamaria et al., 2006) reported that not even one of the 315 key decision-makers in mainstream news media belonged to the Scheduled Castes or Scheduled Tribes. The share of upper caste Hindus or “dwijas” in the upper echelons of the media was 85%. OBCs, who were estimated to constitute around 40% of the population, accounted for a mere 4% of top media jobs. In the English print media, OBCs accounted for just 1% of top jobs and in the Hindi print media 8%.

It is not evident that the absence of Dalit journalists is the result of conscious discrimination. Informal journalist networks, like economic networks as we mentioned before, are conditioned by caste mechanisms and influence the recruitment process. In Indian society, human networks mostly function within the formula of caste (Balasubramaniam, 2011). Policies of reservation, affirmative action, targeted expenditure, and investment could be positive in this context.

Against the positive discrimination policy tentatively fostered by the law, the absence in the news of a discourse about common life of the underrepresented, a discourse of autonomy and individuality, will diminish identity and self-affirmation of Dalits and lower castes. However, the permanent elaboration on conflict and violent episodes is bound to increase gaps, faults, and cleavages between caste communities, leading to further revolt and conflict. This imbalance of discourse leads to a narrow definition of Indian culture. It will unavoidably reinforce negative stereotypes and perceptions of superiority and inferiority among the different members of society: Brahmin will always be mostly associated with progress, law and culture and Dalit with violence, conflict and difficulty.

When writing news about castes, journalists contribute to the interpretation of the caste system in media, and ultimately in public opinion. If media are unevenly controlled by certain caste, as we saw in the introduction, an apology of those caste values will inevitably emerge in the media. This possibility makes it difficult to have an unbiased reading of the many diverse caste issues that pervade Indian society. Also, any initially conflicting or divisive nature of some caste problems may grow disproportionately due to the media reports and increase the cleavage between casts. The same phenomenon happens nowadays in the context of social media. For the purpose of bigger consumption, social media companies using automated editorial algorithms can sort their audiences ever more efficiently into groups of the like-minded individuals. This way they create echo chambers that amplify their particular views. It’s no accident that on some occasions, people of different political views, being so disparately radicalized, cannot even understand each other. It’s also no surprise that terrorist groups have been able to exploit with big efficiency social media to deadly effect (Sunstein, 2018).

For this reason, computational quantitative analysis of the theme of castes in newspapers can bring a more objective view of the unbalance between casts in Indian society. Our analysis merely identifies the most salient topics in the universe of the “caste” discourse in a completely subject agnostic way. The analysis is exclusively based in statistics over words thus avoiding any preconceptions, or prior knowledge about the theme. A pos-interpretation of the statistical results is nevertheless necessary. However, this pos-analysis work is much less constrained by prior experience of the researcher, and much more informed by objective and measurable objective properties latent in the huge quantity of written text and impossible to manually analyze.

When implementing a computational analysis based on statistics over a “bag of words” approach, it is very difficult to separate the use of words like “caste” or “dalit” as element of information, debate and argumentation on news texts, its semantic power and its meaning, from its simple value as a signifier, as a Saussurean sign. The interpretation of any quantitative analysis should be backed up by interpretation of the use of the string of alphabet characters that make up the words. When applying topic detection, with the necessary interpretation of a set of words that should correspond to a certain detected topic, the only subjective intervention of the researcher is the mapping between words and topics. The analysis becomes much more objective in comparison with an extensive content analysis of the text. Although perhaps lacking nuance and detail, it nevertheless can provide enough evidence to the assurance of solid claims and conclusions.

Methods

To attain our research goals, we applied Natural Language Processing algorithms to a large dataset of news collected

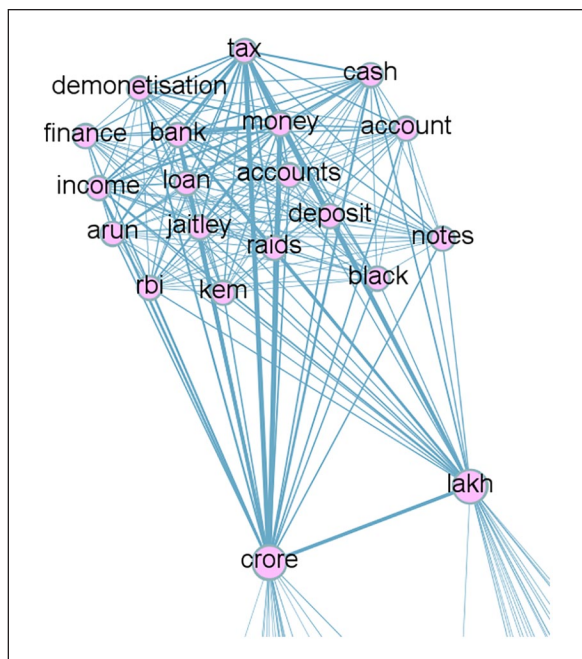


Figure 2. Word representation of the topic “money and finance” as an example.

words represent co-occurrence in the same topic. The thickness of the link between two words represents the sum of their respective probabilities in the topic.

To quantify the robustness of each topic, a measure of “topic coherence” (Schnabel et al., 2015) was calculated. Measures of topic coherence quantify how much the words characterizing the topic support each other. These measures are mostly based on frequency of word co-occurrence in sliding windows over the text dataset. We calculated the “Cv” coherence measure being the most adequate measure as reported in the literature (Röder et al., 2015). Computing the coherence for each topic, we have iteratively applied the LDA algorithm with different initial number of topics to maximize coherence. The best results were obtained for the detection of 25 topics, either for the “caste” articles as also for the rest of the dataset. In Supplemental Annex B, we report several examples of topic word sets.

Results

The resulting two set of topics, for the set of news articles regarding the caste subject and for the rest of the dataset, are represented in Table 1, including corresponding coherence measures. These topics represent, on the left, the set of “caste” articles, and on the right the remaining dataset of articles for one run of the LDA algorithm. Other runs did not present very different results. On the left of each topic there is its coherence value. For the “caste” articles, three of the least coherence topics did not present a clear definition.

Regarding the subset of articles related to the “caste” subject, the most relevant topics the status of institutionalized castes, politics, controversy, violence, justice and religion. By examining in detail the articles concerning family, agriculture, and students, it is clear that these articles report such themes in the context of political, judicial, or protest events.

Concerning the rest of the dataset, we see that they present a broader scope of topics and they are much better defined than in the former set. In Supplemental Annex B, you can examine in more detail some of these topics. We can see that the main subject are several kinds of politics—home and international, government and parties—economics, entertainment (arts, celebrity, sports, and film), and judicial affairs. With the exception of women violence, there are no topics directly concerning personal cases involving repression, family or student injustice and controversy, as in the case of the “caste” subset of articles.

In an overall assessment, we can say that the “caste” related articles in the newspapers on our dataset are mostly associated with underprivileged, repressed or victimized, human situations in comparison with the general news articles present in the same collection. The subject “caste” is almost exclusively focused on prejudice, violence, and conflict situations.

Discussion

In 1978, the American Society of News Editors (ASNE), to represent the voice of minorities in news media, made a pledge to achieve parity with the percentage of people of color in the general population by the year 2000. Since then, their survey, the Newsroom Diversity Survey,³ has the purpose of quantifying to what extent minorities are represented in American newsrooms. Data show many American newsrooms are still nowhere close to achieving this goal of fair parity. The most recent survey of 2017 (see Note 3) reported that minority individuals (Black, Asian, Hispanic, Native American, or other), especially in masthead functions, being much better than in 1978, still misrepresent the common society proportions.

Similar analysis was performed regarding casts in India (Chamaria et al., 2006; Jeffrey, 2001; Uniyal, 1996). In 2006, “Twice born” Hindus (Brahmin, Kshatriya, Vaishya, others) had a share of over 85% in newsrooms, representing not more than 16% in the population. The complementary share of “Intermediary” Hindu castes (Jat, Reddy, Maratha, Patel, etc.), Hindu OBC, Muslim, Christian, Sikhs, but more importantly SC and ST, which accounted for the rest 84% of the population, was thus terribly misrepresented.

Conclusion

With our computational analysis we saw that the prevalent discourse in the “caste” related articles is almost entirely focused on prejudice, violence, and conflict situations.

Table 1. List of the Topics.

Caste articles		Remaining articles	
0.75	Reservation demand	0.61	Film
0.72	Reservation demand	0.57	Trending
0.67	Reservation	0.55	Judicial
0.66	Caste violence	0.52	Political parties and elections
0.66	Controversy	0.52	Celebrity and television
0.64	Popular culture	0.47	International affairs
0.63	Uttar Pradesh politics	0.47	Insurgency, Pakistani politics
0.62	Election	0.45	Cricket
0.55	Political controversy	0.45	Fundamentalism
0.55	Religious fundamentalism	0.44	Violence
0.54	International affairs	0.41	Bollywood
0.52	Controversy	0.40	Women specific crime (dowry, infanticide)
0.52	Status of SC/ST/OBC	0.40	General politics
0.48	Family	0.40	Economics and Finance (transportation and energy)
0.45	Religion	0.39	Economics and Finance (land and agriculture)
0.45	Social and political violence	0.39	Regional
0.42	Politics and cricket	0.38	National politics
0.41	Christianity	0.37	Demonetization
0.40	Agriculture	0.33	Inter religious tension
0.39	Justice for SC/ST	0.33	Medical services
0.39	Student politics	0.32	Political competition
0.38	Regional politics	0.32	Festival culture
	unclear	0.31	West Bengal
	unclear	0.29	Governmental politics
	unclear	0.28	International politics

Although this perspective is mostly commendable, once the underprivileged castes do in fact suffer prejudice and violence, this monotonic discourse does not in fact reflect and report the multitude of aspects necessarily present in the life of the lower caste population. Every community has its tradition, habits, private life, and the most dear or problematic life events that contribute to its happiness and reinforce its identity. Life and events that deserve to be reported and known about within the public sphere. A lack of agency, autonomy, and vocality of the lower caste population, which also is evident in the newsrooms, is severely affected and augmented with this type of discourse, in which underprivileged people is inevitably seen merely as a source of trouble and problems.

With this work, which largely complement the seminal work of R. Jeffrey (2001), we hope to reinforce his claim and to contribute to a better understanding of the Indian public sphere in which concerns caste discrimination, underlining the importance of news media in the improving of democracy but also the important role Natural Language Processing and computational methods can play in social science research and communications studies.

Declaration of Conflicting Interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

The author(s) received no financial support for the research, authorship, and/or publication of this article.

ORCID iD

António Filipe Fonseca  <https://orcid.org/0000-0001-5483-421X>

Supplemental Material

Supplemental material for this article is available online.

Notes

1. Varna literally means type, order, color, or class and was a framework for grouping people into classes, first used in Vedic Indian society and frequently referred in the ancient Indian texts. The four classes were the Brahmins (priestly people), the Kshatriyas (also called Rajanyas, who were rulers, administrators, and warriors), the Vaishyas (artisans, merchants, tradesmen, and farmers), and Shudras (laboring classes). The varna categorization implicitly had a fifth element, being those people deemed to be entirely outside its scope, such as tribal people and the untouchables (Bayly, 2001; Fowler, 1997; Stanton et al., 2012).
2. Jati, meaning birth, is mentioned much less often in ancient texts, where it is clearly distinguished from varna. The jatis are complex social groups that lack universally applicable definition or characteristic, and have been more flexible and

diverse than was previously often assumed (Bayly, 2001; Fowler, 1997).

3. ASNE Newsroom Diversity Survey 2017, <https://www.asne.org/diversity-survey-2017> accessed in October 2018.

References

- Balasubramaniam, J. (2011). Dalits and a lack of diversity in the newsroom. *Economic and Political Weekly*, 46(11), 21–23.
- Bayly, S. (2001). *Caste, society and politics in India from the eighteenth century to the modern age* (Vol 3). Cambridge University Press.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Chamaria, A., Kumar, J., & Yogendra, Y. (2006). *Survey of the social profile of key decision makers in the national media* [Unpublished report]. Center for the Study of Developing Societies.
- Damodaran, H. (2008). *India's new capitalists: Caste, business, and industry in a modern nation*. Springer.
- Deshpande, A. (2011). *The grammar of caste: Economic discrimination in contemporary India*. Oxford University Press.
- Deshpande, A. (2017, December 11). The ugly reality of caste violence and discrimination in urban India. *The Wire*. <https://thewire.in/caste/ugly-reality-caste-violence-discrimination-urban-india>
- Fowler, J. D. (1997). *Hinduism: Beliefs and practices*. Sussex Academic Press.
- Hnatkovska, V., & Lahiri, A. (2011, August). *Convergence across castes* [IGC working papers]. International Growth Center. <https://econ.washington.edu/sites/econ/files/old-site-uploads/2014/09/Convergence-Across-Castes.pdf>
- Human Rights Watch. (2007, February). *Hidden apartheid caste discrimination against India's "Untouchables"* (Vol. 19). <https://www.hrw.org/reports/2007/india0207/>
- Jeffrey, R. (2001). [NOT] being there: Dalits and India's newspapers. *South Asia: Journal of South Asian Studies*, 24(2), 225–238.
- Loomba, A. (2016). The everyday violence of caste. *College Literature*, 43(1), 220–225.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient estimation of word representations in vector space*. <https://arxiv.org/pdf/1301.3781.pdf>
- Munshi, K., & Rosenzweig, M. (2016). Networks and misallocation: Insurance, migration, and the rural-urban wage gap. *American Economic Review*, 106(1), 46–98.
- Nigam, A. (2004). Caste politics in India. *South Asian Journal*, 4.
- Röder, M., Both, A., & Hinneburg, A. (2015). Exploring the space of topic coherence measures. In *Proceedings of the eighth ACM international conference on web search and data mining* (pp. 399–408). ACM.
- Schnabel, T., Labutov, I., Mimno, D., & Joachims, T. (2015). Evaluation methods for unsupervised word embeddings. In *Proceedings of the 2015 conference on empirical methods in natural language processing* (pp. 298–307). Association for Computational Linguistics.
- Scott, J., & Marshall, G. (Eds.). (2009). *A dictionary of sociology*. Oxford University Press.
- Sifton, J. (2016, June 7). Challenges & opportunities: The advancement of human rights in India. *Human Rights Watch*. <https://www.hrw.org/news/2016/06/07/challenges-opportunities-advancement-human-rights-india-0>
- Stanton, A. L., Ramsamy, E., Seybolt, P. J., & Elliott, C. M. (Eds.). (2012). *Cultural sociology of the Middle East, Asia, and Africa: An encyclopedia*. SAGE.
- Sunstein, C. R. (2018). *#Republic: Divided democracy in the age of social media*. Princeton University Press.
- Teltumbde, A. (2018). *Republic of caste: Thinking equality in the time of neoliberal Hindutva*. Navayana.
- Thorat, S., & Neuman, K. S. (2012). *Blocked by caste: Economic discrimination in modern India*. Oxford University Press.
- Uniyal, B. N. (1996). In search of a Dalit journalist. *The Pioneer*, 16.

Author Biographies

António Filipe Fonseca (PhD) obtained in 2015 his PhD in complexity sciences from ISCTE-IUL Lisbon University Institute and the Faculty of Sciences of the University of Lisbon, Portugal. He develops his research at ISTAR Information Sciences, Technologies and Architecture Research Center. He teaches at the Doctoral Program in Complexity Sciences in Lisbon and at the Indian Institute of Technology Gandhinagar India. He is interested in Computational Social Sciences, Scale and Sustainability, Networks and Communication Studies.

Sohhom Bandyopadhyay (M.Sc) is currently pursuing a PhD in Cognitive Science at the Indian Institute of Technology Gandhinagar, India. He finished his Masters in Cognitive Science in 2017, and has a Bachelors in Computer Science.

Jorge Louçã (PhD) obtained in 2000 his PhD in informatics from the Université Paris Dauphine, Paris, France and the Faculty of Sciences of the University of Lisbon, Portugal. He is Professor of Informatics at Lisbon University Institute. He founded the Master and Doctoral Programmes in Complexity Sciences, joint academic programs from the ISCTE-IUL and the University of Lisbon. He presently develops his research at ISTAR—Information Sciences and Technologies and Architecture Research Centre, at ISCTE-IUL. He is an active member of the Complex Systems Society (CSS), and collaborates with the UNESCO UniTwin Complex Systems Digital Campus. He is the coordinator of The Observatorium research group, aiming to real-time monitoring of multi-level network structures for the study of knowledge generation and opinion dynamics on the Internet.

Jaison A Manjaly (PhD), is Jasubhai Memorial Chair Associate Professor of Philosophy & Cognitive Science at Indian Institute of Technology Gandhinagar, India. He currently heads the Humanities and Social Science program at IIT Gandhinagar.