

March 2017

Maintaining intellectual diversity in data science

Richard P MANN ^{a,1} and Olivia WOOLLEY-MEZA ^b

^a*Department of Statistics, School of Mathematics, University of Leeds, UK*

^b*Computational Social Science, ETH Zurich, Zurich, Switzerland*

Abstract. Data science is a young and rapidly expanding field, but one which has already experienced several waves of temporarily-ubiquitous methodological fashions. In this paper we argue that a diversity of ideas and methodologies is crucial for the long term success of the data science community. Towards the goal of a healthy, diverse ecosystem of different statistical models and approaches, we review how ideas spread in the scientific community and the role of incentives in influencing which research ideas scientists pursue. We conclude with suggestions for how universities, research funders and other actors in the data science community can help to maintain a rich, eclectic statistical environment.

Keywords. collective intelligence, diversity, contagion networks

1. Introduction

In 2012 the Harvard Business Review declared Data Scientist to be the ‘sexiest job of the 21st century’ [8]. The last five years have borne out that pronouncement. As illustrated in Figure 1, global interest in data science has increased exponentially, at least as measured by the number of related searches on Google. There has been a huge increase in the number of universities offering courses in ‘data science’ or ‘data analytics’, led by student demand in response to a rapid growth in the number of well-paid ‘data scientist’ job positions. Doubtless, some of this is a relabeling of previously extant activities; much of what we teach our data analytics students has previously been covered in statistics and machine-learning programs, while some companies advertising data scientist positions would have once advertised for statisticians. Nonetheless, within these fields there has also been a substantial increase in activity. To illustrate, the Annual Conference on Neural Information Processing Systems (NIPS), regarded as the leading venue for publishing and discussing research in machine-learning, has seen huge year-on-year increases in registrations over the past five years, with nearly 4000 delegates attending in 2015. Overall, it is clear that there has been a step change in the number of individuals and organizations involved in performing sophisticated analyses of large and complex data. And there is no sign yet that this growth in data analytics in industry, government and academia is slowing down. Within scientific research, this huge increase in data analytic capabilities and activity presents us with an important question: how can we ensure that

¹Corresponding Author: R. P. Mann E-mail: r.p.mann@leeds.ac.uk

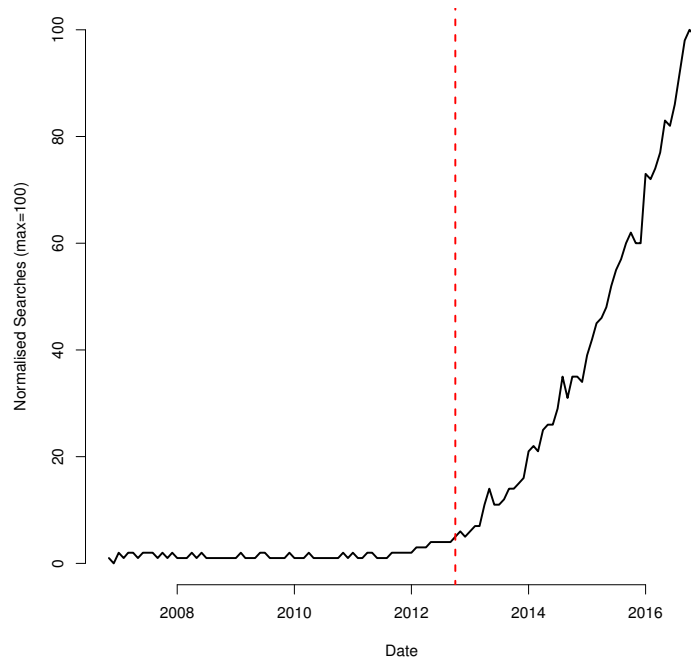


Figure 1. The monthly number of searches for ‘data science’, as extracted from Google Trends, between October 2006 and October 2016. Counts are normalized to a maximum of 100 over the time period. The vertical dashed line indicates the date on which The Harvard Business Review declared Data Scientist to be ‘the sexiest job of the 21st century’.

all this work is creating genuine new knowledge? It may seem natural to assume that if more people are trained to analyze data using exciting new tools like Random Forests, Deep Neural Networks and Gaussian Processes, that we should expect more insightful, robust analyses of data to result, and therefore obtain more knowledge from our scientific endeavors. However, in many cases these tools are being applied to data that is statistically troublesome: observational data, often unstructured, subject to strong selection biases, without controls and with many interacting factors potentially affecting the outcome of interest. Examples include text and behavior extracted from social media [31], hospital admission data [20], and store loyalty card records. With rare exceptions [17] these do not permit anything resembling the classic randomized, controlled trials that are the gold standard of causal inference. Moreover, many methods typically employed in machine-learning and industrial data analytics are primarily focused on *predictive* accuracy, rather than inference and interpretation of underlying causal structures. Finally, and importantly for this perspective, data analytic techniques are subject to ‘bubbles’ of interest with the scientific community. In the 1980’s artificial neural networks were firmly at the forefront of machine-learning and artificial intelligence research. The popularity of these waned in 1990’s and 2000’s in favor of Gaussian processes [27], Random Forests [6] and other non-parametric methods, before a resurgence led by new techniques for training many-layered neural networks, termed ‘Deep Learning’, in the 2010’s. These

March 2017

waves of interest in one technique or another have a pronounced effect on the collective scientific enterprise, as they reduce the diversity of statistical models and approaches being used to investigate similar data sets. As such, any problems inherent to a single type of statistical model can be amplified if that technique becomes popular within the community. Meanwhile, the particular advantages of other methods can become lost as their popularity wanes. It is far from clear that the popularity of one method over another is strictly related only to its analytical power; instead there are strong undercurrents of fashion and conformism in the methods researchers are expected to use. Furthermore, conformisms and fashion are further magnified by network effects that give a small number of researchers and methods exponentially more visibility.

Lack of diversity in analytical approach is a missed opportunity. There is substantial cross-disciplinary evidence for the important role diversity plays in collective intelligence [33,24]: experimentally in human [37] and animal behavior [2], in theoretical models of collective behavior [38,14,21], and specifically in the successes of statistical ensemble models. Ensemble models use combinations of multiple, often very many, different statistical models to perform data analysis. For example, the popular Random Forests model [6] is an ensemble model: many distinct decision tree models are generated from a single data set, before being aggregated to make predictions from new data. The advantage of an ensemble approach is that each model within the ensemble may identify and utilize different features and patterns within the data. Aggregated together, the errors made by each model cancel out to some degree, resulting in a collectively accurate prediction. The power of aggregating distinct models for a given data set became evident to many at the conclusion of the Netflix Prize competition [4]. This competition set participants the challenge of improving on the accuracy of Netflix's own algorithm for predicting future film watching choices by at least 10%. For some time no single team working with one statistical model was able to achieve this mark. The competition was eventually concluded when several teams came together, combining their different models so as to improve their overall predictions [5,34,26]. This outcome demonstrated that, in statistical modeling, the search for one 'true' or 'best' model is often misguided. Instead, as a community we should seek the best *combination* of approaches, especially when faced with complex, multi-dimensional phenomena. Since it is rarely possible or desirable to centrally coordinate a search for a good collection of statistical models, we must instead consider how the incentives individuals face and the networks they inhabit influence the type and variety of statistical research that they perform.

1.1. Collective wisdom, collective madness

The ability of a group to exhibit intelligence superior to any of its constituent members is well established. One of the first to study how collective intelligence emerges was Condorcet, who considered the case of an idealized jury [9]. Consider n jurors, tasked with deciding whether a defendant is guilty or innocent. Each juror is individually only accurate in making this determination with a relatively low probability, say 60%. That is little better than guessing at random. However, assuming that the jurors make up their minds independently, Condorcet showed that collective group decision (determined by a simple vote) is far more likely to be accurate than a single individual, growing quickly with the number of jurors, n . Francis Galton made similar observations regarding the accuracy of collective estimation [11], which is ultimately predicated on the Law of Large

March 2017

Numbers: independent errors made by many individuals tend to cancel out in aggregate, making the group much “smarter” than a single individual. As noted above, these early observations have been replicated since within data science, in ensembles of models and research teams.

However, collective wisdom is contingent on diversity and independence between members of a group [33]. The counterpoint to the collective wisdom of Condorcet’s idealized jury is the collective madness we see when individuals are too strongly influenced by the collective mood. This is exhibited in famous examples of damaging group think, such as Tulip Fever, the South Sea Bubble and other stock market booms and busts. We may also observe on a daily level in our own lives instances where social conformity and peer pressure lead individuals and groups to behave sub-optimally. The community of researchers in data science, statistics, machine learning and artificial intelligence is also subject to social forces that discourage a diversity of approaches. Certain statistical models become fashionable, are accepted as the new big thing and are soon ubiquitous. Senior researchers train their students in the methods that they know. At any one time, certain types of model are, generally, more accurate and/or efficient than others, and researchers tend to gravitate towards these. The culture of benchmarking one’s new method against the state of the art in terms of accuracy necessitates that researchers utilize the best currently available methodologies if they wish to get their work published. Previous research has shown that when individuals are rewarded solely on the basis of the accuracy of their individual predictions, this tends to lead to a severe lack of diversity and a subsequent catastrophic reduction in collective wisdom, putting the whole collective on a par with a single individual [14,21].

In the light of these incentives and mechanisms that discourage diversity, and in the knowledge that diversity is critical to the collective wisdom of any community, it behooves us to consider carefully how ideas spread in the research community, and how the incentives and structures inherent to the scientific community can be used to encourage a wider diversity of research approaches.

2. Understanding the spread of ideas

Diverse ideas and methods are combined through a decentralized emergent process, as scientists become aware of information, communicate it, “adopt” or use these ideas and create new ones. This process is driven by the dynamic and multi-layered structure of interactions between scientists and the mechanisms via which ideas are taken up. Here we will discuss insights from the study of complex networks [3,23,10] and spreading processes on these networks [16,25] that can shed light on how to structure scientific interactions in a way that sustains a diverse pool of methods in data science and encourages researchers to combine and integrate them in productive ways.

2.1. Networks of scientific interaction

Analyzing the networks of scientific can reveal important information about the current divisions between different methods and the groups of researchers that use them. Furthermore, these network structures also yield a better understanding of the potential for these methods to be better integrated through new interactions. There are many rich

March 2017

sources of data that have recently become available that can be used to characterize the interaction of scientists in different contexts. For example, the Web of Science or other bibliometric sources can be used to construct citation networks. Although this is only an imperfect measure of scientific ideas and how they spread, there is quantitative evidence that scientific ideas “spread” through citations [18] and that these networks can be used to predict their spread [30]. Using bibliometric data we can also extract the structure of scientific collaborations and co-authorship networks. These networks can be used to understand the potential for spread of ideas and methods between researchers. Recent analysis shows that individuals are more likely to be cited by those closer to them in a network of scientific co-authorships [29]. There are of course other, less formal forms of scientific interaction which are becoming easier to measure. For example, measurements of face-to-face interactions at conferences [32] and exchange and “friendship” on online social platforms such as Twitter and ResearchGate can also be used to map the structure of interactions between scientists.

All of these quantifications of scientific interaction contain different information at different temporal resolutions. However, two pervasive characteristics are community structure and a high inequality, or broad distribution of the centrality (most basically measured through connectedness) of different network components, whether they be researchers or methods.

2.2. *Community structure, individual centrality and integrating diverse methods*

Modules or communities are, intuitively speaking, the sub units of a network made up of individuals (e.g. researchers or methods) that interact more strongly with each other than with the rest of the network. One interesting approach to integrating diverse perspectives would be to combine methods from distant communities through ensemble methods.

Taking a more decentralized long-term approach, communication channels between researchers operating in these distant communities also need to be established. However, there are many reasons to be weary of a naive approach which simply encourages more unstructured interaction. To retain useful diversity, interaction across communities must not compromise an basic degree of isolation and independence between communities.

The importance of community structure is better understood if we consider the mechanisms that seem to govern the spreading of ideas and innovations. The best studied model is fractional threshold contagion [35], where the probability that an individual becomes “infected” through an infective contact is dependent on the *fraction* of other infected individuals that it interacts with (in network terms, its neighbors). This rule captures the idea that adoption is a social process: there is pressure to conform, or there are added synergistic benefits to adopting if others whom we interact with adopt. Given such a process, communities serve as incubators for new ideas, through local reinforcement. But communities also have the opposite effect slowing the adoption ideas that originate outside of them. Increasing connectivity randomly throughout a network could decrease the diversity of ideas, allowing only the most contagious, or those that start in the most well connected places, to persist.

Theoretical work indicates that networks that have modular structure but sustain intermediate levels of connectivity between modules provide optimal conditions for the global uptake of ideas [22]. Furthermore, recent work shows that the most effective way to transfer ideas between communities is through connections made between individuals

March 2017

that are more peripheral rather than through those better connected [15] (see Fig. 2 for a schematic of the different network topologies we discuss).

The dynamics of contagion through a fractional threshold mechanisms are driven by the inverse relationship between the centrality of an individual (according to network degree) and its susceptibility to ideas. This highlights the key role that the more peripheral individuals play in the sustaining a rich and diverse set of scientific ideas. The sensitivity of peripheral researchers becomes even more valuable if we look beyond models that only consider one idea spreading in isolation. Ideas are part of an ecosystem, and they interact with each other through their scientific “hosts”, akin to the dynamics of co-infection and super-infection studied in evolutionary epidemiology [1]. Competition for limited human time and attention is the most studied interaction [36,12,13]. However, there are also synergistic effects between ideas. Most obviously, ideas that are similar and consistent each other are more likely to be adopted by the same scientist. These interaction effects accentuate rich-get-richer feedbacks in the system and thus reduce diversity. Furthermore, since the most established scientists tend to be the most connected, they experience more information overload and therefore the competition for their attention is greater, and at the same time, they are less likely to adopt novel ideas that contradict the establishment that has secured their privileged position.

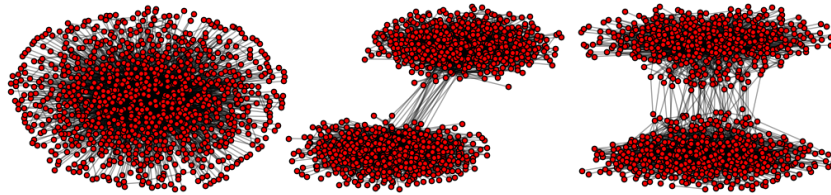


Figure 2. Three different network topologies that can sustain different diversity of ideas: All networks exhibit stratification, that is, most individuals have a low degree but there is a small number of highly connected individuals. The network on the left has no community structure. This network cannot sustain much diversity. The network in the middle has a strong community structure with inter-community connections through central individuals. In this network the diverse ideas that spark in the different communities cannot spread globally. The network on the right has strong community structure, but with inter-community connections through peripheral individuals. This is the topology that best sustains the global penetration of diverse ideas that are fostered locally.

2.3. *The role of incentives*

The type of research that individuals pursue and the methods they use are influenced by their peers and networks, as we have seen above. However, the decisions scientists make in this regard should not be viewed simply as a passive result of community pressure. Instead, these decisions are active choices that are made both to satisfy the researcher’s curiosity, but also to achieve their professional objectives such as promotion, funding and recognition. How scientists are rewarded for their research will affect, in positive or negative ways, the diversity of ideas and therefore the collective wisdom of the scientific community. Scientific ideas and publications exist in a quasi-market, where some are accorded a high value and attract high rewards. A typical researcher is unlikely to pursue highly novel but unpopular ideas if there is no potential for them to be recognized through

March 2017

the awarding of prizes, promotion or further research funding if those ideas later turn out to be fruitful. Likewise, if rewards are systematically allocated to incremental and low-risk research, this will tend to attract more researchers to these areas.

When individual rewards are oriented solely towards simple criteria of success this tends to suppress diversity and risk-taking [14,21]. Conversely, there is recent evidence that useful diversity can be encouraged by rewarding accurate *minority* predictions [21]. These are occasions on which an individual or a model predicts correctly, while the majority of others predict incorrectly. This creates an incentive to focus on less exploited sources of information, or more niche features of a data set, since the individual cannot win any reward by simply imitating what the majority of others are already doing. For the individual, this may make their model less accurate overall. But it makes their model far *more* useful to the community, as it contributes additional information not already presented by others.

Similar ideas are already applied in established methods of ensemble creation. For example, the technique of Boosting [28] is a meta-algorithm for assembling ensembles of weak classifiers that act as a strong classifier in aggregate. A common way to achieve this is to iteratively add weak classifiers to an ensemble, re-weighting the data set under consideration after each new classifier is added. Examples within the data that are currently poorly classified are given higher weights, while those which are already well classified are given lower weights. In this way, additional weak classifiers are ‘incentivized’ to focus on accurately classifying the examples that are currently poorly modeled.

3. Conclusion

Given the great uncertainty about the mechanisms driving information dynamics in science, we have to be cautious in suggesting what interaction structures can best sustain diversity of methods in data science and their productive and eclectic integration. However, some guiding principles are clear.

Diversity thrives in a network structure that allows communities to work in partial isolation. Work carried out in the in the borders between communities and at the periphery of the establishment can lead to important innovations, thus enough funding and other forms of incentives need to be allocated to these regions. Detecting existing communities of methods and promoting their integration is an opportunity to improve the power of methods. Most straightforwardly, this can be done through the ensemble methods we have discussed.

Beyond planned combination of methods, connections that sustain interchange between separate communities, and from the periphery to the core of the research community are necessary. However, these connections must be well timed, since premature competition with more established ideas can be counterproductive. Furthermore, the individuals in the core, who accumulate connections and prestige, are not necessarily the most effective integrators. These individuals suffer from the most acute information overload and they have the most to lose when established approaches are overturned. Thus, increasing connections between scientists working at the periphery, in communities that are typically distant, could be a promising new way of fostering a diverse set of ideas and integrating them for innovative science. There is also an important role for funders, conference organizers and university hiring committees in protecting small, potentially unfashionable research areas for the benefit of the wider scientific ecosystem.

March 2017

We also must not lose sight of the fact that the allocation of funding and other rewards is not simply a mechanism for enabling researchers with various interests, but also acts as a driver of those interests. Most researchers wish for some degree of recognition and reward for their work, and will gravitate to areas that offer this. By reforming rewards to encourage diversity, for instance by explicitly favoring minority research ideas, we can avoid wasteful and potentially damaging group think and maintain the rich variety of data analytic approaches that has enabled the building of ensemble efforts.

Much research remains ahead of us. We need better mapping of the structure of scientific interactions by aggregating information from the different sources available. Of special interest is measuring and characterizing communities that capture the time-varying structure of scientific interaction. This will enable a faster and more precise identification of the scientists and methods emerging in new communities and in community boundaries. We also need to better understand the distinct roles that content and social forces play in willingness to communicate and adopt ideas. This requires both a theoretical effort and more in depth and rigorous empirical validation and model selection. The greatest challenge is perhaps in arriving at an understanding of the mechanisms that drive productive combination and true innovation. The data and theories necessary to do this are becoming available, and the information technologies necessary to implement the insight garnered are also feasible. We could miss a big opportunity if we don't invest in this direction. Even worse, these same tools could lead to a decrease in diversity and scientific productivity if we allow them to work unchecked. However, we must remember the limits of predicting scientific innovation [7]. In contrast to the typical setting studied in collective intelligence, scientific approaches cannot be evaluated entirely on their immediate or short-term performance. As first postulated by Kuhn [19], paradigm changing ideas are accepted because of their yet unverified potential and because the social dynamics support change. Inflexible and over-engineered incentives and communication channels will obstruct dynamic and creative data science.

References

- [1] S. Alizon. Co-infection and super-infection models in evolutionary epidemiology. *Interface focus*, 3(6):20130031, 2013.
- [2] L. M. Aplin, D. R. Farine, R. P. Mann, and B. C. Sheldon. Individual-level personality influences social foraging and collective behaviour in wild birds. *Proc. Roy. Soc. B*, 281(1789):20141016, 2014.
- [3] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- [4] R. M. Bell and Y. Koren. Lessons from the Netflix prize challenge. *ACM SIGKDD Explorations Newsletter*, 9(2):75–79, 2007.
- [5] R. M. Bell, Y. Koren, and C. Volinsky. The bellkor solution to the netflix prize, 2007.
- [6] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [7] A. Clauset, D. B. Larremore, and R. Sinatra. Data-driven predictions in the science of science. *Science*, 355(6324):477–480, 2017.
- [8] T. H. Davenport and D. Patil. Data scientist: The sexiest job of the 21st century. *Harvard Business Review*, 90(10), 2012.
- [9] N. De Condorcet. *Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix*. Cambridge University Press, 2014.
- [10] D. Easley and J. Kleinberg. *Networks, Crowds, and Markets: Reasoning About a Highly Connected World*. Cambridge University Press, 2010.
- [11] F. Galton. Vox populi (the wisdom of crowds). *Nature*, 75(7):450–451, 1907.

March 2017

- [12] J. P. Gleeson, K. P. O’Sullivan, R. A. Baños, and Y. Moreno. Effects of network structure, competition and memory time on social spreading phenomena. *Physical Review X*, 6(2):021019, 2016.
- [13] N. O. Hodas and K. Lerman. How visibility and divided attention constrain social contagion. In *Privacy, Security, Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Conference on Social Computing (SocialCom)*, pages 249–257. IEEE, 2012.
- [14] L. Hong, S. E. Page, and M. Riolo. Incentives, information, and emergent collective accuracy. *Managerial and Decision Economics*, 33(5-6):323–334, 2012.
- [15] J. Huisman and O. Woolley-Meza. Ultra-peripheral links drive structural instability in complex contagion. *in preparation*, 2017.
- [16] M. J. Keeling and K. T. Eames. Networks and epidemic models. *Journal of the Royal Society Interface*, 2(4):295–307, 2005.
- [17] A. D. Kramer, J. E. Guillory, and J. T. Hancock. Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences*, 111(24):8788–8790, 2014.
- [18] T. Kuhn, M. Perc, and D. Helbing. Inheritance patterns in citation networks reveal scientific memes. *Physical Review X*, 4(4):041036, 2014.
- [19] T. S. Kuhn and D. Hawkins. The structure of scientific revolutions. *American Journal of Physics*, 31(7):554–555, 1963.
- [20] R. Mann, F. Mushtaq, A. White, G. Cervantes, T. Pike, D. Coker, S. Murdoch, T. Hiles, C. Smith, D. Berridge, et al. The problem with big data: Operating on smaller datasets to bridge the implementation gap. *Frontiers in Public Health*, 2016.
- [21] R. P. Mann and D. Helbing. Minorities report: optimal incentives for collective intelligence. *ArXiv e-prints*, Nov. 2016.
- [22] A. Nematzadeh, E. Ferrara, A. Flammini, and Y.-Y. Ahn. Optimal network modularity for information diffusion. *Physical review letters*, 113(8):088701, 2014.
- [23] M. Newman. *Networks: an introduction*. OUP Oxford, 2009.
- [24] S. E. Page. *The Difference: How the Power of Diversity Creates Better Groups, Firms, Schools, and Societies*. Princeton University Press, 2008.
- [25] R. Pastor-Satorras, C. Castellano, P. V. Mieghem, and A. Vespignani. Epidemic processes in complex networks. *arXiv:1408.2701*, 2014.
- [26] M. Potté and M. Chabbert. The pragmatic theory solution to the netflix grand prize. *Netflix prize documentation*, 2009.
- [27] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. The M.I.T Press, 2006.
- [28] R. E. Schapire. The boosting approach to machine learning: An overview. In *Nonlinear estimation and classification*, pages 149–171. Springer, 2003.
- [29] C. Schulz., O. Woolley-Meza, B. Uzzi, and D. Helbing. A citation impact indicator based on author network distances. *in preparation*, 2017.
- [30] F. Shi, J. G. Foster, and J. A. Evans. Weaving the fabric of science: Dynamic network models of science’s unfolding structure. *Social Networks*, 43:73–85, 2015.
- [31] V. Spaiser, T. Chadeaux, K. Donnay, F. Russmann, and D. Helbing. Communication power struggles on social media: A case study of the 2011-12 russian protests. *Journal of Information Technology & Politics*, 2017.
- [32] J. Stehlé, N. Voirin, A. Barrat, C. Cattuto, V. Colizza, L. Isella, C. Régis, J.-F. Pinton, N. Khanafer, W. Van den Broeck, et al. Simulation of an seir infectious disease model on the dynamic contact network of conference attendees. *BMC medicine*, 9(1):87, 2011.
- [33] J. Surowiecki. *The Wisdom of Crowds*. Random House LLC, 2005.
- [34] A. Töschler, M. Jahrer, and R. M. Bell. The bigchaos solution to the netflix grand prize. *Netflix prize documentation*, 2009.
- [35] D. J. Watts. A simple model of global cascades on random networks. *Proceedings of the National Academy of Sciences*, 99(9):5766–5771, 2002.
- [36] L. Weng, A. Flammini, A. Vespignani, and F. Menczer. Competition among memes in a world with limited attention. *Scientific Reports*, 2, 2012.
- [37] A. W. Woolley, C. F. Chabris, A. Pentland, N. Hashmi, and T. W. Malone. Evidence for a collective intelligence factor in the performance of human groups. *Science*, 330(6004):686–688, 2010.
- [38] A. Zafeiris and T. Vicsek. Group performance is maximized by hierarchical competence distribution.

March 2017

Nature Communications, 4, 2013.