

# Analítica Visual (Visual Analytics) para processos de difusão em redes complexas

**Proponente:** Dr. Renato Fabbri

**Supervisora:** Profa. Dra. Maria Cristina Ferreira de Oliveira

Instituto de Ciências Matemáticas e de Computação,

Universidade de São Paulo (ICMC-USP)

24 de maio de 2020

**Resumo:** Os processos de difusão são modelados de diversas formas, e o estudo destes processos em redes é uma área ampla. Englobam e.g. modelos inspirados em e aplicações para epidemiologia, redes de relações entre genes e proteínas, de informação em estruturas sociais humanas, de estado físico ou de sistemas computacionais. Este projeto propõe o desenvolvimento de métodos e software de análise visual (*visual analytics*) para processos de difusão em redes complexas. Em especial, há propostas de modelos potencialmente inéditos que podem ser estudados com o auxílio de simulações visuais, mesmo nos casos em que os métodos numéricos e a análise algébrica tornam-se custosos, difíceis ou de execução laboriosa. O proponente deste plano de trabalho colabora com o VICG/ICMC/USP, desde que iniciou o pós-doutorado adiantado adaptado para um Treinamento Técnico V (FAPESP XXXX), e pesquisou a caracterização de redes complexas sociais no seu doutorado. Assim, há artigos e ferramentas desenvolvidas através dos quais o pesquisador retomará a pesquisa e estabelecerá conexão entre as pesquisas realizadas e propostas. Por exemplo, foram concebidos métodos e ferramentas de visualização de redes bipartidas e de redes longitudinais (i.e. evolutivas, dinâmicas; que ganham e perdem vértices e arestas em uma sucessão de eventos). Destacam-se os métodos e ferramentas para visualização de redes bipartidas assistidos por estratégias multinível. Ao incrementar as contribuições já efetivadas, algumas interfaces e métodos em software, e artigos, estarão em melhor estado para continuidade e colaboração em outras parcerias e orientações acadêmicas.

**Palavras-chave:** Análise visual, Redes sociais, Redes complexas, Mineração de texto.

# Visual Analytics of text and topology in social networks

**Proponent:** Dr. Renato Fabbri

**Supervisor:** Prof. Dr. Maria Cristina Ferreira de Oliveira  
Institute of Mathematics and Computational Sciences,  
University of São Paulo (ICMC-USP)

24 de maio de 2020

**Abstract:**

**Keywords:** Visual analytics

# 1 Formulacao do problema

”Essentially, all models are wrong, but some are useful” (George Box, 1987)

## 1.1 Aspectos eticos

Nos trabalhos [] estao consideradas questoes eticas da pesquisa cientifica envolvendo seres humanos. Em especial, foram desenvolvidos emprestimos da antropologia para a pesquisa em fisica, sob o acompanhamento de academicos mais experientes, principalmente e nominalmente: Profa. Dra. Marilia Pisani (Filosofia, UFABC), Profa. Dra. Deborah Antunes (Psicologia, UFC), e Prof. Dr. Massimo Canevacci (Sapienza Universita di Roma). Em resumo, os aspectos eticos foram amenizados atraves do estudo das estruturas sociais do proprio pesquisador (emprestimo da “escrita de diarios”, tecnica da antropologia etnografica), e da manutencao da transparencia da pesquisa com manutencao para acesso publico de textos, dados e software (emprestimo da cultura livre: contribuicao para o legado publico da humanidade de conhecimentos e tecnologias). Com isso, foi possivel realizar alguns experimentos nas redes do proprio pesquisador, e a consideracao destes experimentos em documentos scientificos, sem o investimento de tempo para vencer burocracias e procedimentos de comites de etica.

No ambito deste plano de trabalho, o pesquisador concentrara esforcos no desenvolvimento de modelos, ferramntas em software, simulacoes, e relatoria cientifica. Este aspecto da pesquisa nao implica na necessidade da aprovacao de qualquer comite de etica. Mesmo assim, caso haja espaco e pertinencia para experimentos em estruturas sociais, as propostas dos experimentos serao consideradas pelos pesquisadores proponente e pela responsavel para submissao aos comites de etica cabiveis.

## 1.2 Sincronizacao

Simplificadamente , a sincronizacao eh tratada em conjunto com ou como um efeito permanente na rede [?]. Um pouco por reuso do vocabulario e pelos modelos de interesse, chamamos aqui de sincronizacao uma classe de processos e modelos de difusao.

Alguns modelos desta classe foram criados para desenvolvimento tecnologico pessoal do pesquisador, e visa a propagacao de informacao e concordancias sociais. Considere um grafo (i.e. uma rede)  $G = (V, E)$ , em que  $V = \{v_i\}_0^{N-1}$  sao vertices, e  $E = \{e_i = (v_j, v_k)\}_{i=0}^{z-1}$  sao arestas. Os pressupostos gerais sao:

- alguns vertices  $v_i$ ,  $|v_i| = a \ll N$ , são ativados para iniciar a propagação, chamados sementes ou vertices iniciais.
- Cada vertice, ao ser ativado, ativa  $b$  (constante) vizinhos e não pode mais ser ativado. Caso o vertice não possua  $b$  vizinhos não-ativados, ativará todos os vizinhos possível.

Os critérios de escolha das sementes,  $a$ ,  $b$ , são arbitrários. Além disso, há variantes, e.g. os vertices ativados podem ser novamente ativados, ou  $b$  pode ser variante no tempo, randomico ou dependente de características dos vertices (e.g. grau). Pode-se também impor restrições adicionais, dentre os quais a mais comum é que todos os vertices sejam ativados (sincronização completa).

As experiências descritas em [?, ?, ?, ?, ?] resultaram na imposição de pouca atividade de cada vertice, e propagação desta atividade, pois demoraram meses alguns procedimentos (e.g. os realizados em Dez/2012 até Mar/2013). Os experimentos mais efêmeros eram baseados em comunicação (i.e. ativação) de alguns poucos vertices e não trouxeram mudança nítida na estrutura e funcionamento da rede (assuntos, trocas, etc). Também foi obtida a heurística de começar pelos vertices menos conectados, os periféricos [1, ?]. O detalhamento das justificativas e evidências que corroboram estas características foge ao escopo deste documento e encontra-se na literatura citada.

Por fim, pode-se utilizar técnicas avançadas de ciência de redes. Considere o modelo básico resultante dos pressupostos gerais e considerações acima. Ou seja,  $a$  sementes,  $b$  ativações por ativação,

Para exemplificar um processo de sincronização, considere uma estratégia multinível, em uma rede original  $G_0$  é representada como uma sucessão de redes  $M = \{G_i\}_0^{m-1}$  tal que se  $i < j$ ,  $N_i \leq N_j$  e  $z_i \leq z_j$ . Para a obtenção de  $G_j$  a partir de  $G_i$ , são determinados os conjuntos de vertices que serão colapsados em supervertices através de algum dentre os vários algoritmos de *matching*.  $V_j$  resulta dos supervertices e dos vertices que não fizeram parte de nenhum conjunto colapsado.  $E_j$  resulta das arestas implicadas pela rede  $G_i$  e pela correspondência entre  $V_i$  e  $V_j$ , um algoritmo que depende das características da rede (simples, com peso, direcionada, bipartida, multicamada, heterogênea, etc).

Com o objetivo de obter sementes iniciais e a sequência de vizinhos a serem ativados,

Um conjunto de vertices para colapso consistirá na escolha dos  $b$  vizinhos mais conectados de um vertice. Este vertice será o vertice mais conectado que ainda não pertence a um destes conjuntos. O processo continuará até na rede menor restem somente  $b$  vertices que não participaram de nenhum colapso.

A introdução de ruído na escolha dos vertices e vizinhos a serem colapsados

permite a comparacao e escolha entre as diferentes colecoes de sementes  $S = \{v_i\}_0^{s-1}$  e arestas  $A$  relacionadas a cada vertice  $A = \{(v_i, \{v_j\})_0^{a_i-1}\}_0^{a-1}$  a serem exercitadas na propagacao da ativacao.

Na propagacao de um patogeno, o contagio eh muitas vezes modelado sendo igualmente possivel por cada aresta. Este eh realmente o caso quando o patogeno eh assintomatico. Ja em quadros mais complicados e obitos, o individuo pode contagiar diretamente apenas um numero limitado de pessoas, pois retira-se para tratamento ou internacao. Alem disso, dado que a rede observada eh integrada a outras redes, o patogeno eh introduzido na rede  $G$  por redes nao observadas, sem a necessidade de arestas em  $E$ .

Estas ultimas consideracoes sao uteis para estudos epidemiologicos (e.g. SARs-Covid-19). O modelo em que  $A$  eh determinado adequa-se melhor para aplicacoes comerciais (e.g. Marketing Multinivel) e para *crowdsourcing* de informacao/dados (e.g. democracia liquida e formacao de redes).

### 1.3 Trabalhos relacionados de outros autores

Os processos de difusao em redes complexas tem recebido crescente atencao na literatura cientifica. De fato, os modelos sao utilizados para uma classe vasta de fenomenos que ganharam mais relevancia na ciencia recente, como redes de interacoes entre proteínas e genes em celulas, de contagio na epidemiologia, de informacao, opiniao e fofoca em plataformas e estruturas sociais.

Eh possivel, neste contexto, discernir duas abordagens paradigmaticas. Ha a abordagem abstrata, em que os vertices e arestas nao sao definidos para alem de suas definicoes matematicas ou em que sao minimizadas as particularidades dos sistemas de interesse. Ha abordagens bastante especializadas para os sistemas de interesse (e.g. biologico, social, fisico, tecnologico), e nesta linha provavelmente se destacam as “analises de difusao baseadas em redes” e as redes medicas.

Interessa-nos em especial adaptar os modelos para utilizacao de estrategias multinivel, e redes em que ha uma modificacao sequencial nos vertices e arestas considerados. A caracterizacao de novos modelos, como o descrito na secao anterior, seguirah a pertinencia como reforcada pelos trabalhos dos pesquisadores envolvidos ou por lacunas encontradas na literatura.

Como esperado, o propontente acompanha a pesquisa em andamento na ciencia da complexidade, redes complexas, e redes de difusao. Este acompanhamento inclui os artigos e livros da bibliografia, e fontes que se mostraram uteis no decorrer da pos-graduacao e projeto FAPESP (). Tais fontes incluem MOOCs<sup>1</sup>,

---

<sup>1</sup>a

software<sup>2</sup>, assinatura e contribuicoes em paginas da Wikipedia<sup>3</sup>, e notas publicadas por instituicoes competentes<sup>4</sup>

## 2 Objetivo

O objetivo deste plano de trabalho eh o desenvolvimento de modelos de processos de difusao em redes, e a implementacao de interfaces de visualizacao

### 2.1 Objetivos especificos

## 3 Justificativa

### 3.1 Historico de pesquisa especializada do proponente

## 4 Metodologia

O trabalho sera acompanhado de constante aprofundamento teorico na area. Com base na experiencia do pesquisador na area, ha pertinencia na visita a MOOCs<sup>5</sup>, e.g. *Network Dynamics of Social Behavior* (<https://www.coursera.org/learn/networkdynamics/home/welcome>), este segmento do *Social and Economic Networks: Models and Analysis* (<https://www.coursera.org/lecture/social-economic-networks/5-1-diffusion-0d7yv>). Ha varios materiais (inclusive MOOCs) excelentes e relevantes para este trabalho, por exemplo intro ao NetLogo Complexity explorer<sup>6</sup>.

Serao garantidos um minimo de 12 horas de dedicacao semanais com foco absoluto neste plano de trabalho, a serem distribuidos entre aprofundamentos teoricos, desenvolvimento de metodos e modelos, desenvolvimento e manutencao de software, articulacao com outros pesquisadores e instituicoes, e escrita de artigos e relatorios scientificos. Este eh o numero minimo de horas semanais previstas na Resolucao CoPq N<sup>o</sup> 7413. O pesquisador estarah mantendo outras atividades de pesquisa e engenharia de software relacionadas aa ciencia de redes, habito mantido ha mais de 10 anos. Caso novas tecnologias e metodos desenvolvidos sejam diretamente relacionados a este plano de trabalho, o pesquisador entrarah em contato com a Pro-Reitoria de Inovacao, como previsto na resolucao supracitada.

---

<sup>2</sup>graphology, networkx

<sup>3</sup>[https://en.wikipedia.org/wiki/Network-based\\_diffusion\\_analysis](https://en.wikipedia.org/wiki/Network-based_diffusion_analysis)

<sup>4</sup><https://www.santafe.edu/research/projects/transmission-sfi-insights-covid-19>.

<sup>5</sup>Sigla de *Massive Online Open Courses*

<sup>6</sup><https://www.complexityexplorer.org/>

## 5 Cronograma de execucao

1. Concepcao do plano de trabalho e estabelecimento da parceria.
2. Aprofundamento do conhecimento de teoria de redes complexas e processos de difusao em redes, em conformidade com a Secao .
3. Acréscimos aos modelos atuais de analítica visual e visualização de dados aplicados à análise de redes sociais, com o foco no participante da rede, nos pesquisadores em potencial e na classificação/tipologia de redes e participantes.
4. Implementação computacional. Estamos já implementando layouts para grafos no ccNetViz.
5. Escrita e publicação dos resultados em artigos. Esta etapa está já em andamento pois possuímos diversos escritos com resultados relacionados à mineração e visualização de dados de redes sociais que estão sendo submetidos para publicação.
6. Trocas com pesquisadores externos, estabelecimento de colaborações.
7. Elaboração do relatório científico final.

Atividade	2017	2018		2019
	2°	1°	2°	1°
1	•	•		
2	•	•		
3	•	•	•	•
4	•	•	•	•
5	•	•	•	•
6				•

Tabela 1: Cronograma de atividades ao longo dos semestres.

## A Listagens

Itens selecionados dentre os ja lidos e por ler, para registro do pesquisador e suporte na avaliacao.

## A.1 Artigos

Redes adaptativas melhores com estrategia de difusao do que por consenso [2].

Estimacao de influencia em redes de difusao em tempo continuo: algoritmo menos caro, aplicacao, usa sementes [3].

Arcabouco unificado para redes de difusao: doencas infecciosas, processos economicos e sociais, com foco na semelhanca encontrada para difusao online [4]. Autores excelentes.

Contagio simples e complexo, tratamento matematico. Usa sementes mas nao chama de sementes [5].

Arqueologia da difusao para reconstrucao do historico da progressao [6].

Inferencia de links e taxa de trasmissao em redes de difusao atraves de cascatas de estados dos vertices [7].

Estimacao e otimizacao de influencia em redes de difusao em tempo contínuo [8].

Esparsidade promovendo controle ótimo de consenso e sincronização [9].

Sincronizaco de redes cerebrais quando em descanso [10].

## A.2 Wikipedia

Difusao de informacao, cita redes em aspectos bem definidos, mas nao tem uma  
soh figura de rede, pode ter informacao valiosa para novos modelos e visualiza-  
coes: [https://en.wikipedia.org/wiki/Diffusion\\_of\\_innovations](https://en.wikipedia.org/wiki/Diffusion_of_innovations).

Redes de sincronizacao [https://en.wikipedia.org/wiki/Synchronization\\_networks](https://en.wikipedia.org/wiki/Synchronization_networks).

### A.3 Extra

Possibilita paralelos com os experimentos já realizados a criação de modelos parametrizados com os conceitos apresentados: [https://www.slideshare.net/kevinstrowbridge/diffusion-of-innovations-overview?next\\_slideshow=1](https://www.slideshare.net/kevinstrowbridge/diffusion-of-innovations-overview?next_slideshow=1).

Verbete em handbook da Oxford sobre redes de difusao: <https://www.oxfordhandbooks.com/view/10.1093/oxfordhb/9780199948277.001.0001/oxfordhb-9780199948277-001-0001>

*STRUCTURE AND DYNAMICS OF DIFFUSION NETWORKS*, tese de doutorado em filosofia <https://people.mpi-sws.org/~manuelgr/pubs/thesis-manuelgr-final.pdf>.



## Referências

- [1] R. Fabbri, R. Fabbri, D. C. Antunes, M. M. Pisani, and O. N. Oliveira Jr, “Temporal stability in human interaction networks,” *arXiv preprint arXiv:1310.7769*, 2013.
- [2] S.-Y. Tu and A. H. Sayed, “Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks,” *IEEE Transactions on Signal Processing*, vol. 60, no. 12, pp. 6217–6234, 2012.
- [3] N. Du, L. Song, M. Gomez Rodriguez, and H. Zha, “Scalable influence estimation in continuous-time diffusion networks,” in *Advances in Neural Information Processing Systems 26* (C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, eds.), pp. 3147–3155, Curran Associates, Inc., 2013.
- [4] S. Goel, D. J. Watts, and D. G. Goldstein, “The structure of online diffusion networks,” in *Proceedings of the 13th ACM Conference on Electronic Commerce, EC 12*, (New York, NY, USA), pp. 623–638, Association for Computing Machinery, 2012.
- [5] G. Ghasemiesfeh, R. Ebrahimi, and J. Gao, “Complex contagion and the weakness of long ties in social networks: Revisited,” in *Proceedings of the Fourteenth ACM Conference on Electronic Commerce, EC 13*, (New York, NY, USA), pp. 507–524, Association for Computing Machinery, 2013.
- [6] E. Sefer and C. Kingsford, “Diffusion archeology for diffusion progression history reconstruction,” *Knowl. Inf. Syst.*, vol. 49, pp. 403–427, Nov. 2016.
- [7] M. Gomez-Rodriguez, D. Balduzzi, and B. Schölkopf, “Uncovering the temporal dynamics of diffusion networks,” in *Proceedings of the 28th International Conference on International Conference on Machine Learning, ICML 11*, (Madison, WI, USA), pp. 561–568, Omnipress, 2011.
- [8] M. Gomez-Rodriguez, L. Song, N. Du, H. Zha, and B. Schölkopf, “Influence estimation and maximization in continuous-time diffusion networks,” *ACM Transactions on Information Systems (TOIS)*, vol. 34, no. 2, pp. 1–33, 2016.
- [9] X. Wu and M. R. Jovanović, “Sparsity-promoting optimal control of consensus and synchronization networks,” in *2014 American Control Conference*, pp. 2936–2941, IEEE, 2014.
- [10] A. Ponce-Alvarez, G. Deco, P. Hagmann, G. L. Romani, D. Mantini, and M. Corbetta, “Resting-state temporal synchronization networks emerge

from connectivity topology and heterogeneity,” *PLoS computational biology*, vol. 11, no. 2, 2015.