

Proposta do Trabalho Final PAA: Encontrando os Top-K Influenciadores em uma Rede Social

Vagner Clementino¹

¹Departamento de Ciência da Computação
Universidade Federal de Minas Gerais (UFMG)

vagnercs@dcc.ufmg.br

1. Contextualização

Tradicionalmente as campanhas de marketing se baseiam em determinar um conjunto de consumidores, denominado público-alvo, e posteriormente focam suas ações naquela grupo [Hughes 1996]. Neste contexto, a mineração de dados desempenha um papel fundamental por permitir a construção de modelos que tentam prever o comportamento de um cliente baseado em seu histórico de compras [Kumar et al. 1999]. Nos casos em que esta abordagem têm sucesso, foi possível perceber um aumento na lucratividade [Piatetsky-Shapiro and Masand 1999]. Contudo, este tipo de abordagem possui uma limitação básica: ela considera que a decisão de compra de uma pessoa é independente dos demais consumidores, desconsiderando o impacto que demais clientes, especialmente aqueles mais “próximos”, por ventura possam exercer.

O efeito que os demais consumidores possuem sobre a decisão de compra de um cliente é conhecido em Economia como *externalidade da rede*. Com a expansão da Internet e do uso das redes sociais, este “efeito da rede” têm se mostrado de suma importância em diversos setores, especialmente naqueles ligados diretamente à informação (software, imprensa, telecomunicações e etc.) [Shapiro and Varian 2013].

Neste contexto, imagine que você trabalhe em uma empresa de marketing digital que pretende divulgar um novo produto *A* da marca *X* para o maior número possível de usuários em determinada rede social. Uma primeira estratégia seria divulgar o novo produto para cada usuário da rede, que é conhecida como marketing de massa. Tal opção é cara e baixa escalabilidade. Uma segunda alternativa seria apresentar o produto apenas aos usuários que *seguem* a marca *X*, utilizando, deste forma, um marketing direcionado. Esta estratégia peca pela sua abrangência, tendo em vista que não se tem a garantia que as informações do novo produto chegará aos demais usuários da rede social. Uma terceira via seria identificar um grupo de usuários tais que *a partir deles é possível alcançar qualquer outro usuário da rede*. Ao escolher este grupo de usuários, também denominados “sementes”, haverá uma maior probabilidade de que a informação chegue aos demais usuários. Este trabalho tem o foco nesta última abordagem.

A maneira que uma informação é propagada em uma rede é conhecida como *Modelo de Propagação*. Neste sentido, a “probabilidade” de propagação da informação dependerá naturalmente do modelo de difusão existente na rede. Contudo, neste trabalho, parte-se da premissa que quanto maior o número de usuários que a informação pode alcançar maior será a propagação independente do modelo de propagação da rede. Na subseção 2.2 iremos descrever o Modelo de Propagação adotado neste trabalho.

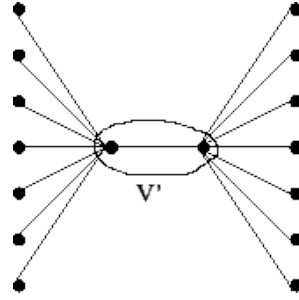


Figura 1. Cobertura de Vértice

Uma rede social pode ser modelada como um grafo não direcionado $G(V, E)$, onde o conjunto de vértices V representa os usuários e o conjunto de arestas E representa os relacionamentos entre os mesmos. Na área de *Teoria dos Grafos*, o problema de definir o conjunto $V' \subseteq V$, de *tamanho mínimo*, tal que para todo vértice em V é possível alcançar qualquer vértice em V é conhecido como **Cobertura de Vértice**. A subseção 2.1 define o problema formalmente.

Em um primeiro momento o conjunto V' (Cobertura de Vértice) poderia ser definida como as sementes da ação de marketing. Contudo, na prática, o tamanho de V' pode ser grande o suficiente de modo a inviabilizar o orçamento da campanha. Desta forma, se faz necessário encontrar um subconjunto de $W \subset V'$ que alcance um maior número de pessoas sem, todavia, estourar o orçamento. Normalmente o tamanho deste subconjunto de W' deverá ser igual a k , onde k é o número máximo de usuários que a campanha conseguirá patrocinar. O problema de encontrar W tal que $|W| = k$ é conhecido na literatura como **Máxima Influência** e será definido formalmente na subseção 2.2.

2. Definição Formal dos Problemas

2.1. O Problema da Cobertura de Vértice

O problema da *Cobertura de Vértices* [Garey and Johnson 1979] pode ser definido formalmente como segue: dado um grafo $G = (V, E)$ e um número inteiro $K \leq |V|$ verificar se existe uma cobertura de vértice de tamanho K ou menos para G , isto é, um subconjunto $V' \subseteq V$ tal que $|V'| \leq K$ e, para cada vértice $\{u, v\} \in E$, pelo menos um, u ou v , pertence a V' . A figura 2.1 ilustra um grafo com destaque para sua cobertura de vértice V' de tamanho 2.

Encontrar a Cobertura de Vértice mínima para um grafo $G = (V, E)$ qualquer é *NP-completo* [Garey and Johnson 1979, Cormen et al. 2009]. Desta forma, até o momento, não existe algoritmo de tempo polinomial capaz de resolver o problema. Diante da inexistência de um algoritmo que retorna a solução exata em tempo polinomial, vem sendo propostos na literatura diversos *algoritmos aproximativos*. Dentre eles, o mais utilizado, consegue executar em $O(V + E)$ em um nível de aproximação igual a 2 [Cormen et al. 2009], ou seja, seja C^* uma cobertura de vértice ótima, o algoritmo retornará uma cobertura de vértice C , tal que $|C^*| \leq |C| \leq 2 \times |C^*|$.

2.2. O Problema da Máxima Influência

Uma definição formal para a *Máxima Influência* pode ser descrita como: dado uma cobertura de vértice V' do grafo $G = (V, E)$, encontrar $W \subset V'$, de modo que $|W| \leq K$ e

$\bigcup_{j=1}^k w_j$ é maximizado, onde $w_j \in V'$ e $1 \leq j \leq k$. No contexto deste trabalho, maximizar $\bigcup_{j=1}^k w_j$ significa influenciar um conjunto maior de usuários.

Conforme exposto anteriormente, a fim de mensurar o valor $\bigcup_{j=1}^k w_j$ devemos definir um *Modelo de Propagação*. No trabalho proposto, será utilizado o modelo conhecido como *Linear Threshold Model* [Granovetter 1978, Schelling 2006]. Neste modelo, um vértice v é influenciado por cada um dos seus vizinhos w de acordo com um peso $b_{v,w}$ que respeita a Equação 1.

$$\sum_{w \text{ vizinho de } v} b_{v,w} \leq 1 \quad (1)$$

Para cada vértice v é definido aleatoriamente um *threshold* θ_v no intervalo $[0, 1]$; este limiar representa a fração dos vizinhos que deve se tornar ativo para que o vértice v torne-se ativo. Dada uma escolha aleatória dos *threshold* e um conjunto inicial de nós ativos A_0 (sendo os demais nós inativos), o processo de difusão ocorre da seguinte forma: na etapa t , todos os nós que estavam ativos na etapa $t - 1$ permanecem ativos, e um vértice v qualquer se torna *ativo* se o somatório dos pesos de seus vizinhos *ativos* sejam $\geq \theta_v$, conforme Equação 2. Cabe ressaltar que o conceito de *ativação* de um vértice depende do contexto que o modelo está sendo aplicado. Neste trabalho considera que um vértice v foi ativado se ele comprou um determinado produto anunciado.

$$\sum_{w \text{ é um vizinho ativo de } v} b_{v,w} \geq \theta_v \quad (2)$$

Conforme pode ser observado, ao se aplicar *Linear Threshold Model* o problema se resume em encontrar o conjunto A_0 tal que $|A_0| = k$ e que maximize o total de vértice ativos, denominado conjunto S . É fácil verificar que $|S| = k$ antes da execução do algoritmo.

Para modelo de propagação que foi considerado, o problema de determinar um valor de k que maximize a influência é NP-Difícil [Garey and Johnson 1979]. Existe na literatura um algoritmo guloso que consegue um fator de aproximação da ordem de $1 - \frac{1}{e}$ [Hochbaum 1997], onde e é a base do logaritmo natural. Desta forma, no melhor caso existe uma garantia de desempenho um pouco acima de 63%. O algoritmo guloso que consegue este é baseado na abordagem apresentada em [Domingos and Richardson 2001]. Assim qualquer trabalho que proponha em resolver o problema da Máxima Influência pode considerar aquele valor como linha de base (“base-line”).

3. Modelagem de Trabalho

O problema da *Máxima Influência* pode ser modelado como um grafo não direcionado $G(V, E)$, onde o conjunto de vértices V representa os possíveis consumidores e as arestas E representam os relacionamentos entre os mesmos. O processo de propagação das informações utilizará *Linear Threshold Model*. Desta forma, para cada consumidor $v \in V$ será atribuído aleatoriamente um *threshold* θ_v no intervalo $[0, 1]$. As arestas em E serão ponderadas com os valores $b_{v,w}$ respeitando à Equação 1. Cada vértice $v \in V$ terá um

atributo determinando se ele está ativo, sendo que ante da execução do algoritmo todos estarão com o valor *não ativo*.

A heurística proposta (Seção 5) deverá determinar o conjunto A_0 , também conhecido como sementes, tal que $A_0 \in V$ e $|A_0| \leq k$. A partir de A_0 será executado o modelo de propagação *Linear Threshold Model* que ao final da execução resultará em um conjunto $S \in V$ que contém os vértices que foram ativados durante o processo. O tamanho de S para diversos valores de k serão comparados com o baseline descrito na Seção 4.

4. Algoritmo Exato

Em [Kempe et al. 2003] foi proposto um algoritmo guloso que adiciona vértices ao conjunto A_0 que maximiza o conjunto S . Os autores não entraram em detalhes sobre como escolheram vértices que efetivamente aumentam o tamanho de S . Naquele trabalho, a heurística proposta foi comparada com três outras:

- Na primeira abordagem os vértices em A_0 foram escolhidas randomicamente;
- Na segunda heurística, foram escolhidos k vértices em ordem decrescente de seus graus d_v ;
- A terceira abordagem utilizou o conceito de *Distance centrality* [Scott 2012], que é uma medida influência largamente utilizada na sociologia, e parte do pressuposto de que vértices mais próximos terão mais chances de influenciar uns aos outros.

A heurística proposta por [Kempe et al. 2003] teve um desempenho satisfatório comparada com as demais. Desta forma, este trabalho utilizará estas quatro heurísticas como baseline. Para tanto será utilizado o mesmo *dataset* (vide Seção 6) e os mesmos valores para o atributo k .

5. Heurística Proposta

Este trabalho parte de seguinte premissa para definir o conjunto A_0 : se a partir de um vértice $v \in V$ é possível um grande número de vértices, logo v é um bom candidato para estar em A_0 . Naturalmente os vértices pertencente à *Cobertura de Vértices* do grafo que modela a rede social seriam a primeira alternativa. Conforme exposto, definir uma cobertura de vértice é NP-Completo e a heurística existente consegue retornar uma Cobertura com um nível de aproximação igual a 2. Na prática, o tamanho da Cobertura de Vértice será muito maior do que k , contudo, uma possível alternativa é encontrar dentro da Cobertura encontrada os k vértices que alcancem o maior número de vértices no grafo.

E é justamente seguindo esta linha que a heurística ora proposta determina o conjunto A_0 . Seja $C \in V$ a cobertura de vértice de um $G(V, E)$ obtida utilizando a heurística proposta em [Cormen et al. 2009]. O tamanho de C é no máximo $2 \times C^*$, onde C^* é a Cobertura de Vértice ótima para o grafo $G(V, E)$. Caso $|C| = k$ podemos naturalmente definir $A_0 = C$. Do contrário, devemos encontrar no máximo k vértice em C para fazer parte de A_0 . O Algoritmo 1 apresenta como o conjunto A_0 é definido.

Conforme pode ser observado o Algoritmo 1 utilizada uma estratégia gulosa, com base na Cobertura de Vértices, para definir o conjunto A_0 . A Cobertura de Vértice é definida pelo método FIND-VERTEX-COVER que é baseado em [Cormen et al. 2009]. Contudo, a escolha gulosa é baseado na medida *degree acesss* que o total de vértices que é coberto por um vértice v mas que ainda não foi coberto por nenhum outro vértice em C .

Algorithm 1: FIND-SEEDS retorna o conjunto semente A_0 com base na Cobertura de Vértice de um grafo.

Input: Um grafo não direcionado e não ponderado $G(V, E)$
um inteiro k correspondente a primeiro índice de A ;
um inteiro r correspondente ao último índice de A
Output: Um conjunto $A_0 \in V$ tal que $1 \leq |A_0| \leq k$

```
1  $C \leftarrow \emptyset$ 
2  $A_0 \leftarrow \emptyset$ 
3  $Q \leftarrow \emptyset$   $Q$  é uma fila
4  $C \leftarrow \text{FIND-VERTEX-COVER}(G(V, E))$ 
5 if  $|C| = k$  then
6    $A_0 \leftarrow C$ 
7   return  $A_0$ 
8 else
9    $\text{CALCULE-DEGREE-ACESSS}(G(V, E), C)$ 
10   $\text{SORT}(C)$  Ordenando o conjunto  $C$  em ordem decrescente ao grau
    acessibilidade.
11   $Q \leftarrow C$  Atribuindo o conjunto  $C$  para uma fila
12  while  $|A| < k$  or  $Q$  is not  $\emptyset$  do
13     $v \leftarrow \text{DEQUEUE}(Q)$ 
14    if  $v.\text{degreeAcess} > 0$  then
15       $A_0 \leftarrow A_0 \cup v$ 
16  return  $A_0$ 
17 return  $A_0$ 
```

O *degree acesss* pode ser visto como o tamanho da contribuição exclusiva de um vértice para um possível cobertura de um grafo. O calculo do *degree acesss* é realizado pelo método CALCULE-DEGREE-ACESSS.

Naturalmente, por se tratar de um heurística, a escolha gulosa proposta no Algoritmo 1 não resultará em uma solução ótima global. Todavia, espera-se que o conjunto A_0 definido desta forma resultará em um maior número de vértices ativos (conjunto S).

6. Plano de Experimentos

A fim de avaliar a heurística proposta se faz necessário utilização de um conjunto de dados que apresente as características estruturais de uma rede social de grande escala. Desta forma, de maneira análoga ao trabalho de [Kempe et al. 2003], será utilizado um grafo de colaboração obtido a partir de coautorias em publicações de física [J. Leskovec and Faloutsos 2015]. Têm sido mostrado que redes de coautoria são capazes de capturar as principais características das redes sociais de modo mais geral [Newman 2001].

O gráfico colaboração contém um vértice para cada pesquisador que tem pelo menos um artigo como coautor banco de dados da *arXiv*¹. Caso um artigo tenha dois ou mais

¹<http://arxiv.org/>

autores, o dataset possui uma aresta para cada par de autores. Observe que isso resulta em arestas paralelas quando dois pesquisadores estão como coautores de vários artigos. O gráfico resultante possui 9877 vértices e 25998 arestas e compreende de artigos publicados no período de janeiro de 1993 a abril de 2003 (124 meses). O data está disponível para download em <https://snap.stanford.edu/data/ca-HepTh.html>.

As baterias de teste consistirão da execução da heurística proposta por 10 mil vezes. Este valor se justifica por ser o número de execução realizada pelo trabalho utilizado como baseline do trabalho. A cada nova execução os valores dos *threshold* θ_v de cada vértice serão recalculados aleatoriamente com objetivo de remover qualquer viés que possa prejudicar os resultados do trabalho. Os dados que serão utilizados para comparação será o valor médio do conjunto S obtido de todas as execuções.

Referências

- Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). *Introduction to Algorithms, Third Edition*. The MIT Press, 3rd edition.
- Domingos, P. and Richardson, M. (2001). Mining the network value of customers. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 57–66. ACM.
- Garey, M. R. and Johnson, D. S. (1979). *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., New York, NY, USA.
- Granovetter, M. (1978). Threshold models of collective behavior. *American journal of sociology*, pages 1420–1443.
- Hochbaum, D. S. (1997). Approximation algorithms for np-hard problems. chapter Approximating Covering and Packing Problems: Set Cover, Vertex Cover, Independent Set, and Related Problems, pages 94–143. PWS Publishing Co., Boston, MA, USA.
- Hughes, A. M. (1996). *The complete database marketer: second-generation strategies and techniques for tapping the power of your customer database*. McGraw-Hill.
- J. Leskovec, J. K. and Faloutsos, C. (2015). SNAP Datasets: Stanford large network dataset collection. <https://snap.stanford.edu/data/ca-HepTh.html>.
- Kempe, D., Kleinberg, J., and Tardos, É. (2003). Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 137–146. ACM.
- Kumar, R., Raghavan, P., Rajagopalan, S., and Tomkins, A. (1999). Extracting large-scale knowledge bases from the web. In *VLDB*, volume 99, pages 639–650. Citeseer.
- Newman, M. E. (2001). The structure of scientific collaboration networks. *Proceedings of the National Academy of Sciences*, 98(2):404–409.
- Piatetsky-Shapiro, G. and Masand, B. (1999). Estimating campaign benefits and modeling lift. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 185–193. ACM.
- Schelling, T. C. (2006). *Micromotives and macrobehavior*. WW Norton & Company.
- Scott, J. (2012). *Social network analysis*. Sage.

Shapiro, C. and Varian, H. R. (2013). *Information rules: a strategic guide to the network economy*. Harvard Business Press.