
Robust Nodes Selection for Influence Maximization

Yucheng Ding
517021910416
lecklas.d@sjtu.edu.cn

Abstract

1 We consider a practical phenomenon in the influence maximization problem,
2 where a fraction of the initial nodes are not successfully activated. Such a phe-
3 nomenon would deteriorate the final influence of the selected initial nodes. To
4 overcome the mentioned problem, we study the objective of robust nodes selection
5 for the influence maximization, i.e, how to choose a set of initial nodes to maximize
6 the expected final influence when some of them may be unactivated. We then
7 propose a robust initial node selection (RNS) algorithm with a fast and accurate
8 influence estimation method (IEM). In particular, RNS use IEM, which computes
9 the activation probability of each node and estimates the expected influence, to
10 output a robust seed by selecting some influential nodes and a node seed with
11 great total influence. We extensively evaluate our algorithms over the facebook
12 social network. Evaluation results demonstrate that IEM can accurately estimate
13 the influence and RNS significantly outperforms the conventional greedy algorithm
14 in terms of final influence when a fraction of initial nodes are not activated.

15 1 Introduction

16 Nowadays, a social network, which denotes the relationships and interactions within a group of
17 individuals, plays a fundamental role as a medium for the spread of information, ideas, and influence
18 among its members [Kempe et al.(2015)Kempe, Kleinberg, and Tardos]. Applications such as
19 Facebook, Twitter, and Wechat allow people to connect and communication with each other at
20 anytime and anywhere. As a consequence, viral marketing [Bass(2004),Steffes and Burgee(2009),
21 Domingos and Richardson(2001),Mahajan et al.(1993)Mahajan, Muller, and Bass, Richardson and
22 Domingos(2002)] via social network becomes increasing popular and significant for the industrial
23 circle.

24 To handle the problem, [Kempe et al.(2015)Kempe, Kleinberg, and Tardos] formulates the influence
25 maximization (IM) problem and adopts two basic models, independent cascade (IC) and linear
26 threshold (LT), to model the influence diffusion in the social network. In particular, the problem is,
27 given a social network, select k most influential nodes to maximize the final influence size under the
28 diffusion model, where k denotes the size of the initial node seed.

29 However, in the real viral marketing, we have to consider the following practical characteristics:
30 (1) *Unsuccessful initialization*: when some of the initial clients in the viral marketing have poor
31 experience with the new products, they may not be willing to spread the information. That is, a
32 fraction of the initial nodes may not be successfully activated. (2) *Decreasing diffusion efficiency*: as
33 the information diffusion progresses, the diffusion efficiency decreases. For example, client A is an
34 initial node, B is an out-neighbour of A, and C is an out-neighbour of B. For client B, the final source
35 of information about the new product is A, who is his friend. But for C, the final source is a friend of
36 his friend, whose trust-level is less than C's "direct" friends.

37 Such two characteristics are not including in the existing works. In what follows, we propose a new
38 max-min objective function to model the unsuccessful initialization and make a new assumption



Figure 1: Influence Maximization Problem (Copied from the course PPT).

39 about information diffusion to model the decreasing diffusion efficiency in IC problem. The max-min
 40 objective aims to find such a *robust* initial seed: the successfully activated nodes have great influence
 41 even in the “worst” case, i.e, the most influential nodes are not activated. And the new information
 42 diffusion assumption has covered the decreasing diffusion efficiency by introducing an decreasing
 43 coefficient.

44 Under our objective and diffusion model, we propose Robust Node Selection (RNS) algorithm,
 45 which adopts greedy strategy and a new influence estimation method (IEM). We first prove that
 46 the influence that our approach estimation is quite close to the real influence expectation. We then
 47 analyse the approximation ratio of RNS. Finally, we evaluate IEM and RNS over the facebook dataset.
 48 Evaluation results demonstrate the accuracy of IEM and the effectiveness of RNS with a remarkable
 49 performance improvement compared with the conventional greedy algorithm.

50 2 Problem Formulation

51 We consider an influence maximization problem based on the independent cascade (IC) model. Given
 52 a social network $G(V, E)$, we need to choose a robust initial seed \mathcal{S} ($|\mathcal{S}| \leq k$) to maximize the
 53 expected number of influenced nodes even in the case that at most m ($m \leq k - 1$) nodes refuse to
 54 diffuse positive information.

55 We can formalize it into a max-min optimization problem:

$$\max_{\mathcal{S} \subseteq V, |\mathcal{S}| \leq k} \min_{\mathcal{H} \subseteq \mathcal{S}, |\mathcal{H}| \leq m} F(\mathcal{S} \setminus \mathcal{H}), \quad (1)$$

56 where $F(\mathcal{X})$ denotes the expected number of influenced nodes when the initially active seed is \mathcal{X} .

57 Namely, we call the optimization problem robust nodes selection. To better simulate the information
 58 diffusion, solve the problem, and make approximation analysis, we make the following assumptions.

59 **Assumption 1** (Independent Cascade). *The information diffuses based on the independent cascade*
 60 *model. Starting with a set of initial active nodes, the information diffuses in discrete steps. At each*
 61 *step t , the newly activated node (say, node u) independently activates its out-neighbor (say, node v)*
 62 *with some probability $p(u, v)_t$.*

63 **Assumption 2** (Decreasing Diffusion Efficiency). *As the information diffusion progresses, the*
 64 *probability of successful transmission decreases exponentially. That is, $p(u, v)_t$ is not equal to*
 65 *$p(u, v)_{t-1}$, but satisfies $p(u, v)_t = \delta * p(u, v)_{t-1}$, where $\delta < 1$.*

66 Assumption 1 is based on the independent cascade model, and Assumption 2 is based on the fact that
 67 the diffusion efficiency decreases as the information diffusion progresses, as introduced in Section 1.

68 We then formulise the information diffusion process. In the r -th step, each newly activated node u
 69 independently activates its out-neighbor v with probability $\delta^{r-1} p(u, v)_0$, where $p(u, v)_0 = p(u, v)$
 70 is the original probability of u successfully activating v , which is given as the information of the social
 71 network.

72 3 Algorithm Design

73 In this section, we propose a initial nodes selection algorithm for the robust nodes selection problem.
 74 The basic idea is that we separate the optimization into two parts: the first part is to select the

Algorithm 1 Robust Node Selection (RNS)

1: **Input:** Social influence graph $G(V, E)$, parameter k, m , influence estimation function $\hat{F}(\mathcal{X})$.
2: **Output:** Initial seed \mathcal{S} .
3: $\mathcal{S}_1, \mathcal{S}_2 \leftarrow \emptyset$;
4: **while** $|\mathcal{S}_1| < m$ **do**
5: $s \leftarrow \arg \max_{s \in V - \mathcal{S}_1} F(\{s\})$;
6: $\mathcal{S}_1 \leftarrow \mathcal{S}_1 \cup \{s\}$;
7: **end while**
8: **while** $|\mathcal{S}_2| < k - m$ **do**
9: $s \leftarrow \arg \max_{s \in \mathcal{S} - \mathcal{S}_1 - \mathcal{S}_2} F(\mathcal{S}_2 \cup \{s\})$;
10: $\mathcal{S}_2 \leftarrow \mathcal{S}_2 \cup \{s\}$;
11: **end while**
12: $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$;

75 influential nodes, denoted as \mathcal{S}_1 , and the second part is to select a node seed with great total influence,
76 denoted as \mathcal{S}_2 . The key idea is that if some of the nodes in \mathcal{S}_2 are unactivated, the influential nodes in
77 \mathcal{S}_1 will replace them, which makes up for the influence loss. In particular, we adopt a greedy strategy
78 in each part.

79 The algorithm is presented in details in Algorithm 1. In this part (from lines 4 to 7), we continuously
80 select an element with the largest influence until m elements have been selected. And set \mathcal{S}_2
81 approximate the best set $\mathcal{S} - \mathcal{S}_1$ with \mathcal{S}_1 removed from \mathcal{S} . In this part (from lines 8 to 11), we
82 continuously select an element with the largest marginal influence until $k - m$ elements have been
83 selected. And the selection of \mathcal{S}_2 is the conventional greedy algorithm, which will obtain a node seed
84 with great total influence. Finally, after \mathcal{S}_1 and \mathcal{S}_2 being selected, the algorithm will output a robust
85 initial node seed $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$.

86 4 A Fast Influence Spread Estimation

87 To implement Algorithm 1, the main challenge is how to design the influence estimation function
88 $\hat{F}(\mathcal{X})$. Many existing works estimate the influence using Monte Carlo simulation, which may repeat
89 the independent cascade process for thousands of times. To make a fast and accurate estimation, in
90 this section, we propose a new influence estimation method (IEM). Instead of making Monte Carlo
91 simulation, we estimate the probability that each node is activated. The final expected influence
92 is the sum of the activation probability over all nodes. We first show the accurate expected influence
93 calculation method, and then show that we could only focus on the first few steps to give a fast
94 estimation of the expected influence. To speed up the estimation, we only calculate the first r_0 steps.
95 And in Theorem 1, we show that for each node, the difference between the accurate activation
96 probability and the estimated activation probability can be bounded.

97 4.1 Expected Influence

98 In this section, we introduce the methods of calculating the accurate probability of each node being
99 activated. Let A_0 denote the initially active node seed, and I_v^r denote the probability of node v being
100 activated at *exactly* the r -th step. First we initialize all I_v^0 :

$$I_v^0 = \begin{cases} 1, & v \in A_0 \\ 0, & v \notin A_0 \end{cases} \quad (2)$$

101 When $r \geq 1$, for each node $v \in V - \mathcal{S}$, we have

$$I_v^r = \left(1 - \sum_{i=0}^{r-1} I_v^i \right) \left(1 - \prod_{u \in \Gamma(v)} (1 - \delta^{r-1} I_u^{r-1} p(u, v)) \right), \quad (3)$$

102 where $\Gamma(v)$ denotes the in-neighbours of v . The first part calculates the probability that v has not been
103 activated at the first $(r - 1)$ steps, and the second part calculates the probability that v is activated at
104 the r -th step by its in-neighbours, which was activated at exactly the $(r - 1)$ -th step.

Algorithm 2 Influence Estimation Method (IEM)

1: **Input:** Social influence graph $G(V, E)$, initial node seed S , decreasing coefficient δ .
2: **Output:** Estimated influence IF .
3: Initialize I_v^r with for each node v step $r(r \leq t)$;
4: **for** Each $v \in S$ **do**
5: $I_v^0 \leftarrow 1$;
6: **end for**
7: **for** $r = 1, 2, \dots, t$ **do**
8: **for** Each node $v \in V$ **do**
9: $I_v^r = \left(1 - \sum_{i=0}^{r-1} I_v^i\right) \left(1 - \prod_{u \in \Gamma(v)} (1 - \delta^{r-1} I_u^{r-1} p(u, v))\right)$;
10: **end for**
11: **end for**
12: $IF = \sum_{r=1}^t \sum_{v \in V} I_v^r$;

105 So the total probability of node v being activated is

$$\Pr(v) = \sum_{r=0}^{+\infty} I_v^r. \quad (4)$$

106 which is the sum of the probability of v being activated at all steps.

107 4.2 Influence Estimation Method

108 In this section, we propose a new influence estimation method with a error bound based on the
109 influence calculation method mentioned in the above section. Instead of calculating the probabilities
110 of all steps, we focus on the results coming from the first t steps, and show that the estimated
111 probability is close to the accurate probability. The algorithm is presented in details in Algorithm
112 2. We calculate the probability that each node v is activated at the r -th step, where $r = 1, 2, \dots, t$,
113 according to equation (3) (from lines 3 to 11), and estimate the final influence (line 12).

114 We then prove that the difference between the estimated probability and the accurate probability of
115 each node (say, node v) being activated is bounded.

116 **Theorem 1.** *When we only calculate the first t steps and $D\delta^{t/2-1} \leq 1$, where D is the maximum
117 input degree in $G(V, E)$, the difference between the real probability and estimated probability is
118 bounded by:*

$$Pr(v) - Pr_t(v) \leq \frac{\delta^{t+1}}{1 - \delta},$$

119 **Proof of Theorem 1.** When $v \in S$, it is trivial that $I_v^1 \leq \delta^0 = 1$. Then we focus on the relationship
120 between the activated probability of the neighbouring two steps.

$$\begin{aligned} I_v^r &= \left(1 - \sum_{i=0}^{r-1} I_v^i\right) \left(1 - \prod_{u \in \Gamma(v)} (1 - \delta^{r-1} I_u^{r-1} p(u, v))\right) \\ &\leq 1 - \prod_{u \in \Gamma(v)} (1 - \delta^{r-1} I_u^{r-1} p(u, v)) \\ &\leq \sum_{u \in \Gamma(v)} \delta^{r-1} I_u^{r-1} \leq D\delta^{r-1} I^{r-1}, \end{aligned}$$

121 where $I^{r-1} = \max_{u \in V} \{I_u^{r-1}\}$, so we have $I^r \leq D\delta^{r-1} I^{r-1}$. From the recurrence equation, we can
122 bound the maximum probability that a node is activated at the r -th step: $I^r \leq D^r \delta^{r(r-1)/2}$. So the
123 error, which is the probability that a node is activated after the t -th step is bounded:

$$\text{error}(u) \leq \sum_{r=t+1}^{\infty} D^r \delta^{r(r-1)/2} = \sum_{r=t+1}^{\infty} D^r \delta^{r^2/2} \delta^{-r/2}. \quad (5)$$

124 When $D\delta^{t/2-1} \leq 1$, we have $D^r \delta^{r^2/2} \delta^{-r/2} \leq \delta^r$, and

$$\text{error}(\mathbf{u}) \leq \sum_{r=t+1}^{\infty} \delta^r = \frac{\delta^{t+1}}{1-\delta} \xrightarrow{t \rightarrow \infty} 0. \quad (6)$$

125

□

126 5 Approximation Analysis

127 In this section, we analyse the approximation ratio of Algorithm 1.

128 **Theorem 2.** *Algorithm 1 will output a set \mathcal{A} satisfying*

$$\frac{\hat{F}(\mathcal{A} \setminus H^*(\mathcal{A}))}{\hat{F}(\mathcal{A}^* \setminus H^*(\mathcal{A}^*))} \geq \left(1 - \frac{1}{e}\right) \frac{1}{k-m},$$

129 where $\hat{F}(\mathcal{X})$ is the influence estimation function, $H^*(\mathcal{A})$ is the worst removal of set \mathcal{A} , \mathcal{A}^* is the
130 optimal seed, and $H^*(\mathcal{A}^*)$ is the worst removal of set \mathcal{A}^* .

131 **Proof Sketch of Theorem 2.** We now give the outline of the proof, and details are given in the
132 Appendix.

133 **Lemma 1.** *For an arbitrary instance of the information diffusion model under Assumptions 1 and 2,
134 the resulting estimated influence function \hat{F} is submodular.*

135 **Lemma 2.** *Having removed node set \mathcal{A}_1 from V , Algorithm 1 will choose a set \mathcal{A}_2 , whose estimated
136 influence is no less than $1 - \frac{1}{e}$ of the maximum estimated influence:*

$$\hat{F}(\mathcal{A}_2) \geq \left(1 - \frac{1}{e}\right) \max_{\mathcal{S} \subseteq V \setminus \mathcal{A}_1, |\mathcal{S}| \leq k-m} \hat{F}(\mathcal{S}).$$

137 **Lemma 3.** *We next prove that the maximum estimated influence after removing node set \mathcal{A}_1 is no
138 less than the estimated influence of the optimal seed \mathcal{A}^* with $H^*(\mathcal{A}^*)$ removed from it:*

$$\max_{\mathcal{S} \subseteq V \setminus \mathcal{A}_1, |\mathcal{S}| \leq k-m} \hat{F}(\mathcal{S}) \geq \hat{F}(\mathcal{A}^* \setminus H^*(\mathcal{A}^*)).$$

139 **Lemma 4.** *The estimated influence of \mathcal{A} returned by Algorithm 1 with $H^*(\mathcal{A})$ removed is no less
140 than $1/(k-m)$ of the estimated influence of \mathcal{A}_2 :*

$$\hat{F}(\mathcal{A} \setminus H^*(\mathcal{A})) \geq \frac{1}{k-m} \hat{F}(\mathcal{A}_2).$$

141 Based on the above lemmas, we can derive that

$$\begin{aligned} \hat{F}(\mathcal{A} \setminus H^*(\mathcal{A})) &\geq \frac{1}{k-m} \hat{F}(\mathcal{A}_2) \\ &\geq \frac{1}{k-m} \left(1 - \frac{1}{e}\right) \max_{\mathcal{S} \subseteq V \setminus \mathcal{A}_1, |\mathcal{S}| \leq k-m} \hat{F}(\mathcal{S}) \\ &\geq \left(1 - \frac{1}{e}\right) \frac{1}{k-m} \hat{F}(\mathcal{A}^* \setminus H^*(\mathcal{A}^*)). \end{aligned}$$

142

□

143 Theorem 2 shows that RNS obtains a set of initial nodes, whose influence is $\left(1 - \frac{1}{e}\right) \frac{1}{k-m}$ times of
144 influence of the optimal seed even in the worst removal case. So the seed is actually very robust, no
145 matter which m nodes are not activated, it always has a great total influence.

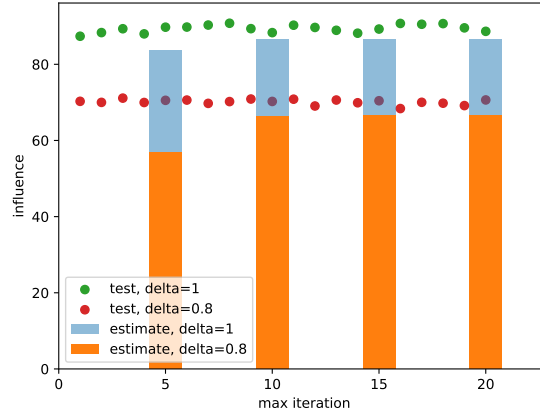


Figure 2: Influence Estimation Method Evaluation

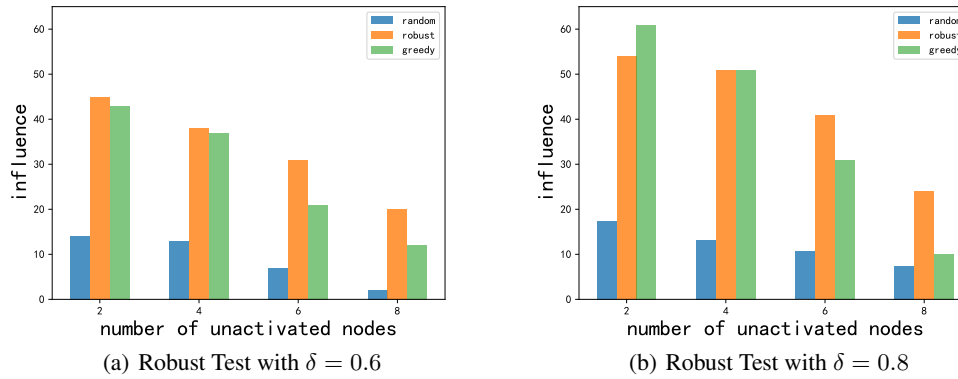


Figure 3: Robust Influence Test with Different Decreasing Coefficient δ .

146 6 Experiments

147 In this section, we evaluate our algorithms. Due to the computation constraint of the computer, we
 148 use a small dataset, which is part of the facebook network and consists of 249 nodes and 407 edges.
 149 And we randomly allocate a weight $p \in (0, 1)$ to each edge.

150 **Influence Estimation Method Evaluation.** We first evaluate our influence estimation function.
 151 Under the setting that the number of steps $t = 5, 10, 15, 20$ and the decreasing coefficient $\delta = 0.8, 1.0$,
 152 we compare the estimation result with the experimental result. We use the result of RNS as the
 153 initial nodes and evaluate the influence estimation method. We first do the influence test for 20 times,
 154 for each experiment we repeatedly execute the information diffusion process under our information
 155 diffusion model for 100 times and calculate the average amount of the final influenced nodes. Each
 156 dot in Figure 6 denotes a experiment result. We then use IEM to estimate the result, which is displayed
 157 in Figure 6 in the form of bar. Figure 6 shows that IEM can well approximate the real influence when
 158 the max step t for estimation is greater than 10.

159 **Robust Node Selection Evaluation.** We then evaluate our robust seed selection algorithm by by
 160 comparing with the randomly chosen seed and conventional greedy algorithm. We set the size of the
 161 initial node seed k to 10, the number of unactivated initial nodes m to 2, 4, 6, 8, and the decreasing
 162 coefficient δ to 0.8. Under each setting, we compare the test average influence for each algorithm.
 163 In particular, for each initial seed, we randomly remove m nodes from it and make the influence
 164 test, which is the average influence of 100 information diffusion experiment. We repeat the above
 165 operations for 5 times, and record the least test influence, which is to approximately find out the
 166 worst removal case. Figure 3 shows that RNS always outperforms the random seed, and have better

167 performance compared with the conventional greedy algorithm when the number of unactivated initial
168 nodes $m \geq 4$.

169 **7 Conclusion**

170 In this paper, we study the practical problem of robust nodes selection for the influence maximization.
171 To handle this problem, we propose RNS, which outputs a robust node seed, and a new influence
172 estimation method, IEM. We extensively evaluate our algorithms over a small facebook social network.
173 Expirical studies over the facebook dataset demonstrate the accuracy of IEM and the effectiveness
174 and robustness of RNS with a remarkable performance improvement compared with the conventional
175 greedy algorithm.

176 **References**

- 177 [Kempe et al.(2015)Kempe, Kleinberg, and Tardos] David Kempe, Jon M. Kleinberg, and Éva Tar-
178 dos. Maximizing the spread of influence through a social network. *Theory of Computing*, 11:
179 105–147, 2015.
- 180 [Bass(2004)] Frank M. Bass. A new product growth for model consumer durables. *Management*
181 *Science*, 50(12-Supplement):1825–1832, 2004.
- 182 [Steffes and Burgee(2009)] Erin M. Steffes and Lawrence E. Burgee. Social ties and online word of
183 mouth. *Internet Res.*, 19(1):42–59, 2009.
- 184 [Domingos and Richardson(2001)] Pedro M. Domingos and Matthew Richardson. Mining the
185 network value of customers. In *KDD*, pages 57–66. ACM, 2001.
- 186 [Mahajan et al.(1993)Mahajan, Muller, and Bass] Vijay Mahajan, Eitan Muller, and Frank M. Bass.
187 Chapter 8 new-product diffusion models. In *Marketing*, volume 5 of *Handbooks in operations*
188 *research and management science*, pages 349–408. North-Holland, 1993.
- 189 [Richardson and Domingos(2002)] Matthew Richardson and Pedro M. Domingos. Mining
190 knowledge-sharing sites for viral marketing. In *KDD*, pages 61–70. ACM, 2002.

191 **A Proof of the Lemmas**

192 **Proof of Lemma 1.** For an arbitrary node set T , a subset $S \subseteq T$, and any node v , we will prove that

$$\hat{F}(S \cup \{v\}) - \hat{F}(S) \geq \hat{F}(T \cup \{v\}) - \hat{F}(T), \quad (7)$$

193 where \hat{F} is the expectation of the influence in the first t round. Let A_t denote the nodes activated at
 194 the first t steps. The diffusion is a random process, and we can view the process as two stage. At
 195 each step, the first stage is that each edge decides to exist or not based on the edge probability, the
 196 second step is that information diffuses according to the connection decided in the first stage. We
 197 will prove that $A_t(T \cup \{v\}) \setminus A_t(T) \subseteq A_t(S \cup \{v\}) \setminus A_t(S)$ when the connection of edges in the
 198 first t steps in the two graphs are the same.

199 For the t - step, we have

$$A_t(T \cup \{v\}) \setminus A_t(T) = A_t(\{v\}) \setminus A_t(T),$$

200 and

$$A_t(S \cup \{v\}) \setminus A_t(S) = A_t(\{v\}) \setminus A_t(S)$$

where $A_t(S) \subseteq A_t(T)$. So we have

$$A_t(\{v\}) \setminus A_t(T) \subseteq A_t(\{v\}) \setminus A_t(S),$$

which indicates that

$$A_t(T \cup \{v\}) \setminus A_t(T) \subseteq A_t(S \cup \{v\}) \setminus A_t(S).$$

In the calculation of the expectation, the probability that two graphs reach the same connectio is equal
 and one-to-one, and for each connection,

$$|A_t(T \cup \{v\}) \setminus A_t(T)| \leq |A_t(S \cup \{v\}) \setminus A_t(S)|.$$

201 We take expectation on both sides and have equation (7), which indicates that the influence estimated
 202 function is submodular. \square

203 **Proof of Lemma 2.** According to the property of submodular functions, we have the $(1 - \frac{1}{e})$ approx-
 204 imation ratio. That is:

$$\hat{F}(\mathcal{A}_2) \geq (1 - \frac{1}{e}) \max_{S \subseteq V \setminus \mathcal{A}_1, |S| \leq k-m} \hat{F}(S).$$

205 \square

206 **Proof of Lemma 3.** We first define $C_1 = A^* \cap A_1$, and then find another set $C_2 \subseteq A^* \setminus C_1$ where
 207 $|C_1| + |C_2| = m$. Since $H^*(A_*)$ minimize the influence $\hat{F}(A^* \setminus H^*(A^*))$, we have

$$\hat{F}(A^* \setminus H^*(A^*)) \leq \hat{F}(A^* \setminus (C_1 \cup C_2))$$

208 In addition, $A^* \setminus (C_1 \cup C_2) \subseteq V \setminus A_1$ and $|A^* \setminus (C_1 \cup C_2)| = k - m$, we have

$$\max_{S \subseteq V - A_1, |S| \leq k-m} \hat{F}(S) \geq \hat{F}(A^* \setminus (C_1 \cup C_2)) \geq \hat{F}(A^* \setminus H^*(A^*)).$$

209 \square

210 **Proof of Lemma 4.** We will prove the lemma case by case:

- 211 1. If $H^*(A) = A_1$, we have $\hat{F}(A \setminus H^*(A)) = \hat{F}(A_2)$, so the lemma holds.
- 212 2. If $H^*(A) \neq A_1$, we have at least one node v in A_1 left. So we have

$$\hat{F}(A \setminus H^*(A)) \geq \hat{F}(\{v\}) \geq \frac{1}{k-m} \sum_{u \in A_2} \hat{F}(\{u\}) \geq \frac{1}{k-m} \hat{F}(A_2).$$

213 \square