

Asymptotic Performance Analysis of Majority Sentiment Detection in Online Social Networks

Tian Tong¹ and Rohit Negi²

Abstract—We analyze the problem of majority sentiment detection in Online Social Networks (OSN), and relate the detection error probability to the underlying graph of the OSN. Modeling the underlying social network as an Ising Markov random field prior based on a given graph, we show that in the case of the empty graph (independent sentiments) and the chain graph, the detection is always inaccurate, even when the number of users grow to infinity. In the case of the complete graph, the detection is inaccurate if the connection strength is below a certain critical value, while it is asymptotically accurate if the strength is above that critical value, which is analogous to the phase transition phenomenon in statistical physics.

I. INTRODUCTION

Online social networks (OSN), such as Facebook and Twitter [1], have a significant influence on people and society. The massive data embedded in these networks have turned OSNs into a gold mine for politicians, economists, and sociologists alike to collect, analyze, and understand the views of people. Therefore, detecting and analyzing the sentiments of OSN users is of great interest in recent machine learning and sociology research [2], [3].

We focus on the problem of *majority sentiment detection*, also known as the vote detection, where the majority sentiment is estimated based on *noisy measurements* of user sentiments. (Majority sentiment is the one among two binary sentiments, such as ‘approve’/‘disapprove’, which predominates among the users.) Related research abounds under various topics such as public opinion studies [4], voting theory [5], and opinion mining [6]. The basic assumptions are that users (or members) in the network are connected by some relationships, such as the friend relationship in Facebook, or follower/followee relationship in Twitter, and that two connected users are more probable to share the same sentiment, since the network typically signifies affinity of opinion. At the same time, automated language processing tools that measure individual sentiments from posts or tweets suffer from noise, due to short size of the text or the inability to recognize sarcasm [7].

In this paper, we attempt to answer an interesting question: how does the network (graph¹ of OSN) influence the error performance of (automated) majority sentiment detection? In particular, we wish to investigate whether such detection is asymptotically accurate, i.e., whether the error probability becomes arbitrarily small as the network size, in terms

of number of users, grows. As we will show, the error performance is strongly related to the graph of the OSN. It involves two levels of influence: the graph structure and the strength of the connections.

To analyze the performance of majority sentiment detection under various network topologies, we model the network by an Ising Markov Random Field (MRF) model [8]. This model was first introduced in statistical physics to interpret the paramagnetic-ferromagnetic phase transition phenomenon. We first provide general upper and lower bounds on the asymptotic detection error probability. Next, we consider special cases of networks to illustrate the phase-transition-like phenomenon in the error probability behavior, where the detection is either inaccurate or is asymptotically accurate. Specifically, we show that in the case of empty graph and chain graph, both of which are weakly connected networks, *the detection is always inaccurate, even when the number of users grows to infinity*. This result appears to be of interest in its own right, since one may naively expect accurate performance in the limit of infinite users. On the other hand, in the complete graph (with standard scaling down of the connection strength), *there exists a critical value for connection strength* (analogous to the critical temperature in statistical physics), below which the detection is inaccurate while above which the detection is asymptotically accurate.

In Section II, we introduce the Ising model of OSN sentiments and analyze majority sentiment detection in the independent sentiment case. In Section III, we obtain bounds on the asymptotic sentiment detection error probability for arbitrary graphs of the Ising model. In Section IV, we consider special cases of graphs to calculate these bounds. Section V shows numerical results while Section VI concludes the paper.

II. SYSTEM MODEL

The social network structure is modeled as an undirected graph, as shown in Fig. 1. Let $\mathbf{X} = (X_1, \dots, X_n)^T \in \{-1, 1\}^n$ denote the vector of binary sentiments of n members, where 1 or -1 denote positive or negative sentiments respectively. These sentiments are unknown to observers. What is observed is a noisy measurement of \mathbf{X} , called \mathbf{Y} . $\mathbf{Y} = (Y_1, \dots, Y_n)^T \in \{-1, 1\}^n$ is modeled as conditionally independent binary measurements of \mathbf{X} each with cross-over probability p , i.e., the output of a binary symmetric channel with input \mathbf{X} . Without loss of generality, we assume that $p < \frac{1}{2}$. This model was first introduced in [7] for the ‘latent sentiment detection’ problem, with more details about the probability mass function of \mathbf{X} following in Section II-B.

*This research was partially supported by NSF awards CCF1422193 and CNS1218823.

Electrical and Computer Engineering Department, Carnegie Mellon University, ¹ttong1@andrew.cmu.edu, ²negi@ece.cmu.edu

¹We use the terms ‘graph of OSN’ and ‘network’ interchangeably.

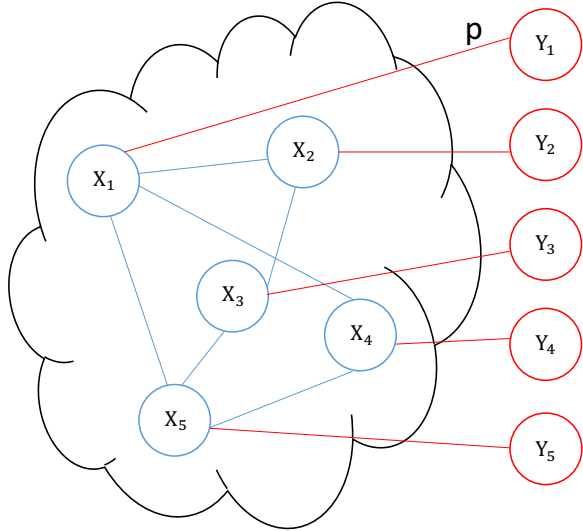


Fig. 1: Markov Random Field model of sentiment detection.

The majority sentiment is defined as

$$m = \text{sign}(\mathbf{1}^T \mathbf{X}), \quad (1)$$

where $\mathbf{1}$ denotes the vector of all ones. We assume that n is odd to avoid trivial ambiguity.

We will use the majority vote detector as our estimator for the majority sentiment:

$$\hat{m} = \text{sign}(\mathbf{1}^T \mathbf{Y}). \quad (2)$$

This detector estimates the majority sentiment as the majority of noisy measurements. In general, it is not the optimal Maximum A Posteriori Probability (MAP) detector. However, it will be sufficient to illustrate the key insights of this paper, promised in Section I, as will be shown below.

The detection error probability (equivalently, classification error probability of the two sentiment class problem) is

$$P_e^{(n)} = \mathbb{P}(m \neq \hat{m}). \quad (3)$$

We wish to investigate whether the majority sentiment detector is asymptotically accurate, i.e., whether $P_e^{(n)}$ becomes arbitrarily small when n is sufficiently large.

In the paper, we denote $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ as the sample average, X_{lim} as a variable distributed as the limiting distribution of $\sqrt{n} \bar{X}_n$, and \xrightarrow{d} as convergence in distribution.

A. I.I.D Case

Before discussing interesting networks, let us first consider the simplest case: X_i s are independent and identical distributed (i.i.d.) random variables, taking values -1 or 1 , each with probability $1/2$. In this example, the members are not connected, or in other words the network is an empty graph. This is the case in contemporary vote detection schemes where voters are assumed to independently vote, and the underlying network is not modeled [5]. For this case, the

detector we adopt in (2) is actually the optimal MAP detector, so that there is no loss of error performance.

For this case, we can derive the exact value of the asymptotic error probability, as shown in Proposition 1.

Proposition 1 *In the i.i.d. (empty graph) case,*

$$\lim_{n \rightarrow \infty} P_e = \frac{2}{\pi} \arcsin \sqrt{p} > 0.$$

Proof: Since (X_i, Y_i) s are i.i.d. random variables, with means zero, and variances $\mathbb{E}[X_i^2] = \mathbb{E}[Y_i^2] = 1$, $\mathbb{E}[X_i Y_i] = 1 - 2p$, the multidimensional central limit theorem tells that:

$$\begin{pmatrix} \sqrt{n} \bar{X}_n \\ \sqrt{n} \bar{Y}_n \end{pmatrix} \xrightarrow{d} \begin{pmatrix} X_{\text{lim}} \\ Y_{\text{lim}} \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 1 - 2p \\ 1 - 2p & 1 \end{pmatrix} \right). \quad (4)$$

Since by definition (3), $P_e^{(n)} = \mathbb{P}(\sqrt{n} \bar{X}_n \sqrt{n} \bar{Y}_n < 0)$, we can calculate its limit from (4) as

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \mathbb{P}(X_{\text{lim}} Y_{\text{lim}} < 0) = \frac{2}{\pi} \arcsin \sqrt{p},$$

where the limit exists thanks to the convergence in distribution (4). ■

In this special case, the detection error probability tends to some positive constant, and thus never reduces to 0, even with infinite number of sentiments. Inherently, this inaccuracy results from the absence of connections. *This result is perhaps counterintuitive, since one may have expected that the majority sentiment can be detected accurately when the user size tends to infinity, as typically happens in the case of single parameter estimation.*

B. Network Model

The network model characterizes the probability mass function of \mathbf{X} . In general, it could be any prior implying that members with connections are more probable to have the same sentiment. In this paper following [7], we adopt a homogeneous Ising MRF prior as

$$p(\mathbf{x}) = \frac{\exp(\theta \mathbf{x}^T A \mathbf{x})}{Z_n(\theta)}. \quad (5)$$

Here A denotes the symmetric graph adjacency matrix, with $A_{ij} = 0/1$ denoting absence/presence of an edge respectively. $\theta > 0$ is called the inverse temperature parameter. We remark that θ characterizes the connection strength in the network, namely, connected members are more probable to share the same sentiment in case of larger θ . The normalizer $Z_n(\theta) = \sum_{\mathbf{x} \in \{-1, 1\}^n} \exp(\theta \mathbf{x}^T A \mathbf{x})$ is called the partition function.

The joint distribution of (\mathbf{x}, \mathbf{y}) can be written as another Ising MRF as

$$p(\mathbf{x}, \mathbf{y}) = \frac{\exp(\theta \mathbf{x}^T A \mathbf{y} + \varepsilon \mathbf{y}^T \mathbf{x})}{Z_n(\theta) (2 \cosh \varepsilon)^n}. \quad (6)$$

Here ε is defined by $p = \frac{\exp(-\varepsilon)}{\exp(\varepsilon) + \exp(-\varepsilon)}$. In the next section, we analyze the asymptotic error probability under the above network model. For this purpose, we assume that as n grows,

there is a given sequence of graphs of n vertices that models the OSN as (5).

III. ERROR PROBABILITY UNDER NETWORK MODEL

In this section, we will first derive an upper bound for the detection error probability in Theorem 1, and two exact asymptotic results in Theorems 2 and 3, assuming a given sequence of graph adjacency matrices A . After that, we will show that the asymptotic performance of the detection error probability is related to the concentration behavior of $\sqrt{n} \overline{X}_n$ in Theorem 4.

Theorem 1

$$P_e^{(n)} \leq \mathbb{E} \left[\exp \left(\frac{-(1-2p)^2}{8(1-p)^2} (\sqrt{n} \overline{X}_n)^2 \right) \right].$$

Here the expectation is taken over \mathbf{X} .

Proof: Let $Z_i = Y_i - (1-2p)X_i$. Since Z_i s are conditionally independent given \mathbf{X} , with $\mathbb{E}[Z_i|\mathbf{X}] = 0$ and $|Z_i| \leq 2(1-p)$, Hoeffding's inequality tells that average \overline{Z}_n satisfies $\mathbb{P}(\sqrt{n} \overline{Z}_n > \epsilon \mid \mathbf{X}) \leq \exp\left(\frac{-\epsilon^2}{8(1-p)^2}\right)$ and $\mathbb{P}(\sqrt{n} \overline{Z}_n < -\epsilon \mid \mathbf{X}) \leq \exp\left(\frac{-\epsilon^2}{8(1-p)^2}\right)$ for any $\epsilon > 0$. By definition (3):

$$\begin{aligned} P_e^{(n)} &= \mathbb{P}(\sqrt{n} \overline{X}_n \sqrt{n} \overline{Y}_n < 0) \\ &= \mathbb{P}(\sqrt{n} \overline{X}_n \sqrt{n} \overline{Z}_n < -(1-2p)(\sqrt{n} \overline{X}_n)^2) \\ &= \mathbb{E} \left[\mathbb{P}(\sqrt{n} \overline{X}_n \sqrt{n} \overline{Z}_n < -(1-2p)(\sqrt{n} \overline{X}_n)^2 \mid \mathbf{X}) \right] \\ &\leq \mathbb{E} \left[\exp \left(\frac{-(1-2p)^2}{8(1-p)^2} (\sqrt{n} \overline{X}_n)^2 \right) \right]. \end{aligned} \quad (7)$$

Theorem 2

$$\liminf_{n \rightarrow \infty} P_e^{(n)} = \liminf_{n \rightarrow \infty} \mathbb{E} \left[Q \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} |\sqrt{n} \overline{X}_n| \right) \right],$$

where $Q(\cdot)$ is the tail probability of standard normal distribution: $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp(-t^2/2) dt$.

Theorem 3 (a) If $\sqrt{n} \overline{X}_n \xrightarrow{d} \Phi$, where Φ is a distribution, then

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \int_{-\infty}^{\infty} Q \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} |x| \right) \Phi(dx).^2$$

(b) Specifically, if $\sqrt{n} \overline{X}_n \xrightarrow{d} N(0, \sigma^2)$, then

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \frac{1}{\pi} \operatorname{arccot} \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} \sigma \right) > 0.$$

To prove Theorems 2 and 3, we first derive a central limit theorem type result in Lemma 1. This will follow from the conditional independence of Y_i s given \mathbf{X} , and reveals that $\sqrt{n} (\overline{Y}_n - (1-2p)\overline{X}_n)$ tends to a normal distribution

conditioned on \mathbf{X} ; besides that, when $\sqrt{n} \overline{X}_n$ converges in distribution, the joint limiting distribution of $\sqrt{n}(\overline{X}_n, \overline{Y}_n)$ can also be obtained.

Lemma 1 (a) *Conditional convergence:* For all \mathbf{X} ,

$$\sqrt{n}(\overline{Y}_n - (1-2p)\overline{X}_n) \mid \mathbf{X} \xrightarrow{d} N(0, 4p(1-p)).$$

(b) *Unconditional convergence:*

$$\sqrt{n}(\overline{Y}_n - (1-2p)\overline{X}_n) \xrightarrow{d} N(0, 4p(1-p)).$$

(c) *Joint convergence:* If $\sqrt{n} \overline{X}_n \xrightarrow{d} \Phi$, then

$$\sqrt{n}(\overline{X}_n, \overline{Y}_n - (1-2p)\overline{X}_n) \xrightarrow{d} (\Phi, N(0, 4p(1-p))),$$

where the two limiting distributions are independent.

Proof: Let $Z_i = Y_i - (1-2p)X_i$. By Lévy's continuity theorem [9], it is equivalent to prove the pointwise convergence of characteristic functions:

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) \mid \mathbf{X}] &= \exp(-4p(1-p)\beta^2/2), \\ \lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n)] &= \exp(-4p(1-p)\beta^2/2), \\ \lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\omega\sqrt{n}\overline{X}_n + j\beta\sqrt{n}\overline{Z}_n)] &= \\ &= \phi(\omega) \exp(-4p(1-p)\beta^2/2), \end{aligned}$$

where $\phi(\cdot)$ denotes the characteristic function of Φ .

For part (a), since Z_i s are conditionally independent given \mathbf{X} , with $\mathbb{E}[Z_i|\mathbf{X}] = 0$, $\mathbb{E}[Z_i^2|\mathbf{X}] = 4p(1-p)$, and $|Z_i| \leq 2(1-p)$, the Lindeberg condition: $\forall \epsilon > 0$,

$$\lim_{n \rightarrow \infty} \frac{1}{n4p(1-p)} \sum_{i=1}^n \mathbb{E} \left[Z_i^2 I(|Z_i| \geq \epsilon\sqrt{n4p(1-p)}) \mid \mathbf{X} \right] = 0$$

is satisfied, where $I(\cdot)$ is the indicator function. Lindeberg-Feller central limit theorem [9] tells that:

$$\lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) \mid \mathbf{X}] = \exp(-4p(1-p)\beta^2/2), \quad (8)$$

where RHS of (8) is the characteristic function of $N(0, 4p(1-p))$.

For part (b), extend the result in (a) to the unconditional version. Notice that $|\mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) \mid \mathbf{X}]| \leq 1$, thus

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n)] &= \lim_{n \rightarrow \infty} \mathbb{E}[\mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) \mid \mathbf{X}]] \\ &= \mathbb{E} \left[\lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) \mid \mathbf{X}] \right] \\ &= \exp(-4p(1-p)\beta^2/2), \end{aligned}$$

where the expectation and limit are exchanged by Lebesgue's dominated convergence theorem [9].

For part (c), define $\Delta_n(\mathbf{X}) = \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) \mid \mathbf{X}] - \exp(-4p(1-p)\beta^2/2)$. By equation (8), $\lim_{n \rightarrow \infty} \Delta_n(\mathbf{X}) =$

²This denotes Lebesgue integral with respect to probability measure Φ .

0, $\forall \mathbf{X}$. So,

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \mathbb{E} [\exp(j\omega\sqrt{n}\overline{X}_n + j\beta\sqrt{n}\overline{Z}_n)] \\
&= \lim_{n \rightarrow \infty} \mathbb{E} [\exp(j\omega\sqrt{n}\overline{X}_n) \mathbb{E}[\exp(j\beta\sqrt{n}\overline{Z}_n) | \mathbf{X}]] \\
&= \lim_{n \rightarrow \infty} \mathbb{E} [\exp(j\omega\sqrt{n}\overline{X}_n) \exp(-4p(1-p)\beta^2/2)] \\
&\quad + \lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\omega\sqrt{n}\overline{X}_n) \Delta_n(\mathbf{X})] \quad (9) \\
&= \phi(\omega) \exp(-4p(1-p)\beta^2/2).
\end{aligned}$$

In equation (9), since $|\Delta_n(\mathbf{X})| \leq 2$ is bounded, $\lim_{n \rightarrow \infty} \mathbb{E}[\exp(j\omega\sqrt{n}\overline{X}_n) \Delta_n(\mathbf{X})] = 0$ by Lebesgue's dominated convergence theorem [9]. ■

Now we use Lemma 1 to prove Theorems 2 and 3.

Proof: [Theorem 2] Start from equation (7). Define $\varepsilon_n(\mathbf{X}) = \mathbb{P}(\sqrt{n}\overline{X}_n\sqrt{n}\overline{Z}_n < -(1-2p)(\sqrt{n}\overline{X}_n)^2 | \mathbf{X}) - Q\left(\frac{(1-2p)}{\sqrt{4p(1-p)}} | \sqrt{n}\overline{X}_n\right)$. From the limiting distribution in Lemma 1 part (a), $\lim_{n \rightarrow \infty} \varepsilon_n(\mathbf{X}) = 0, \forall \mathbf{X}$. So,

$$\begin{aligned}
\liminf_{n \rightarrow \infty} P_e^{(n)} &= \liminf_{n \rightarrow \infty} \mathbb{E}[\mathbb{P}(\sqrt{n}\overline{X}_n\sqrt{n}\overline{Z}_n < \\
&\quad -(1-2p)(\sqrt{n}\overline{X}_n)^2 | \mathbf{X})] \\
&= \liminf_{n \rightarrow \infty} \mathbb{E} \left[Q \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} | \sqrt{n}\overline{X}_n \right) \right] \\
&\quad + \liminf_{n \rightarrow \infty} \mathbb{E}[\varepsilon_n(\mathbf{X})] \\
&= \liminf_{n \rightarrow \infty} \mathbb{E} \left[Q \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} | \sqrt{n}\overline{X}_n \right) \right].
\end{aligned}$$

Here $\liminf_{n \rightarrow \infty} \mathbb{E}[\varepsilon_n(\mathbf{X})] = \mathbb{E}[\liminf_{n \rightarrow \infty} \varepsilon_n(\mathbf{X})] = 0$ by Lebesgue's dominated convergence theorem [9], since $|\varepsilon_n(\mathbf{X})| \leq 1$. ■

Proof: [Theorem 3] Mimic the proof of Theorem 2, with all 'lim inf' replaced by 'lim'.

Finally, based on Theorems 1, 2 and 3, we prove Theorem 4: whether the detection error probability tends to 0 is exactly determined by whether $\sqrt{n}\overline{X}_n$ asymptotically stays away from 0, in probability.

Theorem 4 (a) If $\forall B > 0, \lim_{n \rightarrow \infty} \mathbb{P}(|\sqrt{n}\overline{X}_n| \leq B) = 0$, then $\lim_{n \rightarrow \infty} P_e^{(n)} = 0$.
(b) If $\exists B > 0$ s.t. $\liminf_{n \rightarrow \infty} \mathbb{P}(|\sqrt{n}\overline{X}_n| \leq B) > 0$, then $\liminf_{n \rightarrow \infty} P_e^{(n)} > 0$.

As a comparison to statistical physics, part (a) corresponds to the ferromagnetic phase where spins are mostly in one direction; part (b) corresponds to the paramagnetic phase where spins are nearly equal in both directions.

Proof: For part (a), given any $\epsilon > 0$, choose B large enough such that $\exp\left(\frac{-(1-2p)^2 B^2}{8(1-p)^2}\right) \leq \epsilon/2$, then choose n large enough such that $\mathbb{P}(|\sqrt{n}\overline{X}_n| \leq B) \leq \epsilon/2$. Theorem 1 tells that

$$\begin{aligned}
P_e^{(n)} &\leq \mathbb{E} \left[\exp\left(\frac{-(1-2p)^2}{8(1-p)^2} (\sqrt{n}\overline{X}_n)^2\right) \right] \\
&\leq \mathbb{P}(|\sqrt{n}\overline{X}_n| \leq B) + \exp\left(\frac{-(1-2p)^2}{8(1-p)^2} B^2\right) \leq \epsilon.
\end{aligned}$$

Thus $\lim_{n \rightarrow \infty} P_e^{(n)} = 0$.

For part (b), choose a $B > 0$, such that $\liminf_{n \rightarrow \infty} \mathbb{P}(|\sqrt{n}\overline{X}_n| \leq B) > 0$. Theorem 2 tells that

$$\begin{aligned}
\liminf_{n \rightarrow \infty} P_e^{(n)} &= \liminf_{n \rightarrow \infty} \mathbb{E} \left[Q \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} | \sqrt{n}\overline{X}_n \right) \right] \\
&\geq \liminf_{n \rightarrow \infty} \mathbb{P}(|\sqrt{n}\overline{X}_n| \leq B) Q \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} B \right) > 0.
\end{aligned}$$

■

IV. NETWORK EXAMPLES

Given a network topology (i.e., graph adjacency matrix A), theorems in the previous section characterize the asymptotic performance of the detection error probability. In this section, we will discuss two extreme network examples of Fig. 1: the chain graph (1-D Markov chain), and the complete graph. We will show that in the chain graph (as well as the empty graph as shown in Section II-A), the detection error probability is asymptotically positive; in contrast, *in the complete graph there exists a phenomenon of phase transition*, where the detection error tends to 0 when the connection strength is more than some critical value, while it is asymptotically positive when the connection strength is less than that critical value. *This is similar to the paramagnetic-ferromagnetic phase transition in statistical physics!*

A. Chain Graph

In a chain graph, each vertex is connected to two neighbors, forming a chain. For convenience, we adopt a periodical boundary condition (PBC), namely, the n th member is connected with the first member; nevertheless, boundary conditions asymptotically make no difference. It can be viewed as a Markov chain as: $X_{i+1} = \begin{cases} X_i, & \text{w.p. } \frac{\exp(\theta)}{\exp(\theta) + \exp(-\theta)} \\ -X_i, & \text{w.p. } \frac{\exp(-\theta)}{\exp(\theta) + \exp(-\theta)} \end{cases}, i = 1, \dots, n$, with X_{n+1} treated as X_1 . We will prove that the detection error probability is asymptotically positive, in Proposition 2 below.

Proposition 2 *In the chain graph,*

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \frac{1}{\pi} \operatorname{arccot} \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} e^\theta \right) > 0.$$

Proof: In the chain graph, the partition function at field b , defined by $Z_n(\theta, b) = \sum_{\mathbf{x} \in \{-1, 1\}^n} \exp(\theta(x_1 x_2 + \dots + x_{n-1} x_n + x_n x_1)) + b \mathbf{1}^T \mathbf{x}$, is [10]

$$\begin{aligned}
Z_n(\theta, b) &= e^{n\theta} \left(\cosh b + \sqrt{\sinh^2 b + e^{-4\theta}} \right)^n \\
&\quad + e^{n\theta} \left(\cosh b - \sqrt{\sinh^2 b + e^{-4\theta}} \right)^n. \quad (10)
\end{aligned}$$

Plug $b = 0$ into (10) to get $Z_n(\theta) \doteq Z_n(\theta, 0) = (2 \cosh \theta)^n + (2 \sinh \theta)^n$.

The characteristic function of $\sqrt{n} \overline{X}_n$ is

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \mathbb{E} \left[\exp(j\omega \sqrt{n} \overline{X}_n) \right] \\
&= \lim_{n \rightarrow \infty} \frac{1}{Z_n(\theta)} \sum_{\mathbf{x} \in \{-1,1\}^n} \exp(\theta(x_1 x_2 + \dots + x_n x_1) \\
&\quad + j\omega \sqrt{n} \overline{X}_n) \\
&= \lim_{n \rightarrow \infty} \frac{Z_n(\theta, \frac{j\omega}{\sqrt{n}})}{Z_n(\theta)} \\
&= \lim_{n \rightarrow \infty} \frac{e^{n\theta} \left(\cosh \frac{j\omega}{\sqrt{n}} + \sqrt{\sinh^2 \frac{j\omega}{\sqrt{n}} + e^{-4\theta}} \right)^n}{(2 \cosh \theta)^n + (2 \sinh \theta)^n} \\
&\quad + \lim_{n \rightarrow \infty} \frac{e^{n\theta} \left(\cosh \frac{j\omega}{\sqrt{n}} - \sqrt{\sinh^2 \frac{j\omega}{\sqrt{n}} + e^{-4\theta}} \right)^n}{(2 \cosh \theta)^n + (2 \sinh \theta)^n} \\
&= \exp(-e^{2\theta} \omega^2 / 2).
\end{aligned} \tag{11}$$

In equation (11), we use Taylor's expansions of $\cosh \frac{j\omega}{\sqrt{n}}$ and $\sinh \frac{j\omega}{\sqrt{n}}$ to calculate the limit. As a result, $\sqrt{n} \overline{X}_n \xrightarrow{d} N(0, e^{2\theta})$. So, Theorem 3 tells that:

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \frac{1}{\pi} \operatorname{arccot} \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} e^\theta \right) > 0.$$

B. Complete Graph

In a complete graph, each pair of vertices is connected. The corresponding Ising MRF prior (12), also called the Curie-Weiss [11] prior, is defined slightly differently, in that the strength is weakened to θ/n , to ensure that the total strength from all neighbors of a vertex remains constant, i.e., does not grow with n , though the number of neighbors is $n-1$. Thus, with this standard modification, the prior is

$$p(\mathbf{x}) = \frac{\exp\left(\frac{\theta}{n}(\mathbf{1}^T \mathbf{x})^2\right)}{Z_n(\theta)}. \tag{12}$$

We are interested in the limiting distribution of $\sqrt{n} \overline{X}_n$, with the corresponding variable called X_{lim} . Intuitively, it should have the density $p(x_{\text{lim}}) \propto \exp(\theta x_{\text{lim}}^2) \exp(-x_{\text{lim}}^2/2)$, since in the i.i.d. case $\sqrt{n} \overline{X}_n \xrightarrow{d} N(0, 1)$, and the Curie-Weiss prior introduces a multiplier $\exp(\theta x_{\text{lim}}^2)$. In the case $\theta < \frac{1}{2}$, $\sqrt{n} \overline{X}_n$ converges to a normal distribution, and therefore, by Theorem 3, the detection error probability should be asymptotically positive; otherwise $\sqrt{n} \overline{X}_n$ diverges so the error probability should tend to 0. We shall prove all these formally in Proposition 3.

Proposition 3 *In the complete graph,*

(a) *if $\theta < \frac{1}{2}$, then*

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \frac{1}{\pi} \operatorname{arccot} \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} \frac{1}{\sqrt{1-2\theta}} \right) > 0.$$

(b) *if $\theta > \frac{1}{2}$, then $\lim_{n \rightarrow \infty} P_e^{(n)} = 0$. In fact, it tends to 0 exponentially fast: $\liminf_{n \rightarrow \infty} -\frac{1}{n} \log P_e^{(n)} > 0$.*

Proof: In this model, the partition function at field b , defined by $Z_n(\theta, b) = \sum_{\mathbf{x} \in \{-1,1\}^n} \exp\left(\frac{\theta}{n}(\mathbf{1}^T \mathbf{x})^2 + b \mathbf{1}^T \mathbf{x}\right)$, is [11]

$$Z_n(\theta, b) = \frac{2^n}{\sqrt{\pi}} \int_{-\infty}^{\infty} \exp(-t^2) \cosh^n \left(2\sqrt{\frac{\theta}{n}} t + b \right) dt. \tag{13}$$

For part (a), we calculate the characteristic function of $\sqrt{n} \overline{X}_n$ as

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \mathbb{E} \left[\exp(j\omega \sqrt{n} \overline{X}_n) \right] = \lim_{n \rightarrow \infty} \frac{Z_n(\theta, \frac{j\omega}{\sqrt{n}})}{Z_n(\theta)} \\
&= \frac{\int_{-\infty}^{\infty} \exp(-t^2) \lim_{n \rightarrow \infty} \cosh^n \left(2\sqrt{\frac{\theta}{n}} t + \frac{j\omega}{\sqrt{n}} \right) dt}{\int_{-\infty}^{\infty} \exp(-t^2) \lim_{n \rightarrow \infty} \cosh^n \left(2\sqrt{\frac{\theta}{n}} t \right) dt} \\
&= \frac{\int_{-\infty}^{\infty} \exp(-t^2) \exp\left(\frac{1}{2}(2\sqrt{\theta}t + j\omega)^2\right) dt}{\int_{-\infty}^{\infty} \exp(-t^2) \exp\left(\frac{1}{2}(2\sqrt{\theta}t)^2\right) dt} \\
&= \exp\left(-\frac{\omega^2}{2(1-2\theta)}\right).
\end{aligned}$$

In the second line, the integral and the limit are exchanged by monotone convergence theorem [9]. In the third line, we use the Gaussian integral formula: $\int_{-\infty}^{\infty} \exp(-x^2) dx = \sqrt{\pi}$, with appropriate changes of variables. As a result, $\sqrt{n} \overline{X}_n \xrightarrow{d} N\left(0, \frac{1}{1-2\theta}\right)$. So, Theorem 3 tells that

$$\lim_{n \rightarrow \infty} P_e^{(n)} = \frac{1}{\pi} \operatorname{arccot} \left(\frac{(1-2p)}{\sqrt{4p(1-p)}} \frac{1}{\sqrt{1-2\theta}} \right) > 0.$$

For part (b), we only need to prove that $\liminf_{n \rightarrow \infty} -\frac{1}{n} \log P_e^{(n)} > 0$. Define $C_p = \frac{(1-2p)^2}{8(1-p)^2}$. Start from Theorem 1:

$$\begin{aligned}
P_e^{(n)} &\leq \mathbb{E} \left[\exp(-C_p(\sqrt{n} \overline{X}_n)^2) \right] \\
&= \frac{1}{Z_n(\theta)} \sum_{\mathbf{x} \in \{-1,1\}^n} \exp\left(\frac{\theta}{n}(\mathbf{1}^T \mathbf{x})^2 - C_p(\sqrt{n} \overline{X}_n)^2\right) \\
&= \frac{Z_n(\theta - C_p)}{Z_n(\theta)} \\
&= \frac{\int_{-\infty}^{\infty} \exp(-t^2) \cosh^n \left(2\sqrt{\frac{\theta - C_p}{n}} t \right) dt}{\int_{-\infty}^{\infty} \exp(-t^2) \cosh^n \left(2\sqrt{\frac{\theta}{n}} t \right) dt} \\
&= \frac{\int_{-\infty}^{\infty} \exp(-ns^2 + n \log \cosh(2\sqrt{\theta - C_p}s)) ds}{\int_{-\infty}^{\infty} \exp(-ns^2 + n \log \cosh(2\sqrt{\theta}s)) ds}.
\end{aligned} \tag{14}$$

In the fourth line, we use formula (13), with $Z(\theta) \doteq Z(\theta, 0)$. In the fifth line, we change the variable to $s = t/\sqrt{n}$.

Define $f(\theta, s) = \log \cosh(2\sqrt{\theta}s) - s^2$. From equation (14)

we obtain that:

$$-\frac{1}{n} \log P_e^{(n)} \geq \frac{1}{n} \log \int_{-\infty}^{\infty} \exp(nf(\theta, s)) ds - \frac{1}{n} \log \int_{-\infty}^{\infty} \exp(nf(\theta - C_p, s)) ds.$$

By Laplace's approximation [12], listed below as Lemma 2, whose conditions are satisfied by $f(\theta, s)$ when $\theta > \frac{1}{2}$, $\lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{-\infty}^{\infty} \exp(nf(\theta, s)) ds = \max_s f(\theta, s)$, and similarly for the second term.

It can be observed (proof is omitted here) that $\max_s f(\theta, s)$ is 0 when $\theta \leq \frac{1}{2}$ and monotonically increases with θ when $\theta > \frac{1}{2}$, as illustrated in Fig. 2. Under the condition $\theta > \frac{1}{2}$, noticing that $0 < C_p \leq \frac{1}{8}$ by definition, we state that $\max_s f(\theta, s) > \max_s f(\theta - C_p, s)$. In conclusion,

$$\liminf_{n \rightarrow \infty} -\frac{1}{n} \log P_e^{(n)} \geq \max_s f(\theta, s) - \max_s f(\theta - C_p, s) > 0.$$

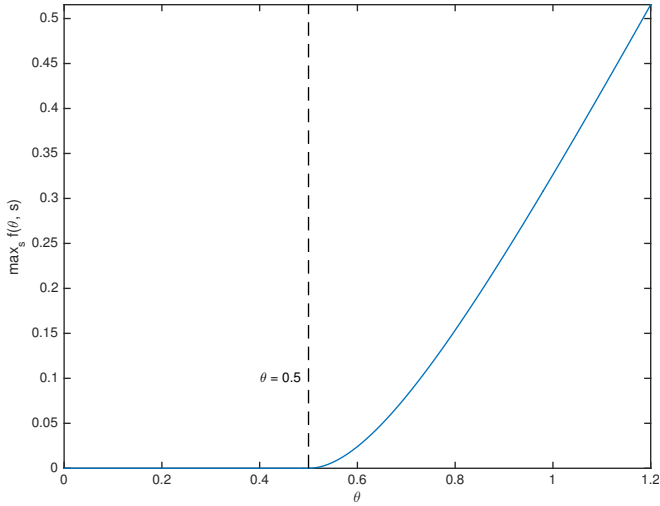


Fig. 2: $\max_s f(\theta, s)$ v.s. θ

Lemma 2 (Laplace's approximation) *Given that $g(s)$ is twice differentiable, with $s^* = \arg \max_s g(s)$, and $g''(s^*) < 0$, then*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{-\infty}^{\infty} \exp(ng(s)) ds = \max_s g(s). \quad (15)$$

V. NUMERICAL RESULTS

In this section, we will provide numerical results for two examples to illustrate the theoretical results presented.

First for the i.i.d. (empty graph) case, Fig. 3 illustrates the detection error probability $P_e^{(n)}$ versus user size n , for several cross-over probabilities p . The horizontal lines denote the limit in Proposition 1. It can be seen that $P_e^{(n)}$ tends to the positive constant given in Proposition 1.

Next for the complete graph case, Fig. 4 illustrates $P_e^{(n)}$ v.s. n for several cross-over probabilities p : Fig. 4a with $\theta = 0.3$ corresponds to part (a) of Proposition 3, where the detection error probability tends to the positive constant

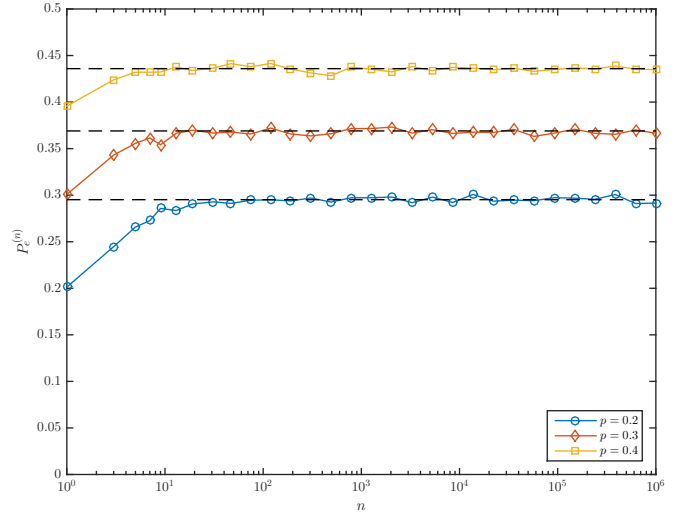


Fig. 3: $P_e^{(n)}$ v.s. n in i.i.d. (empty graph) case.

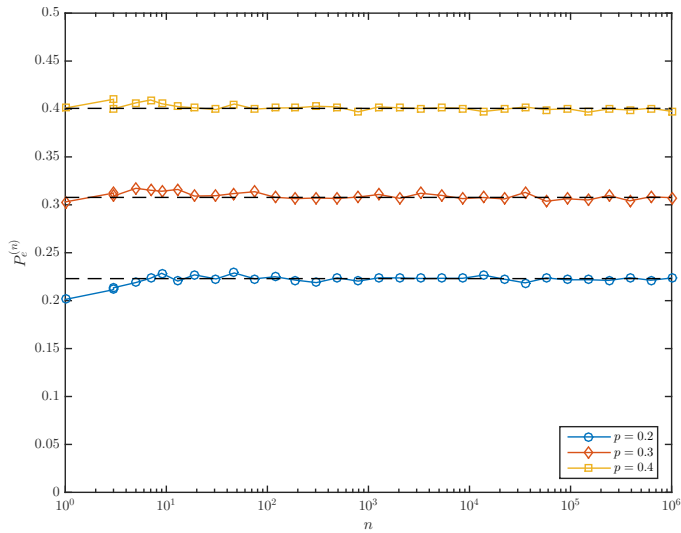
given there; Fig. 4b with $\theta = 0.7$ corresponds to part (b) of Proposition 3, where the error probability decays to 0 exponentially fast.

VI. CONCLUSION

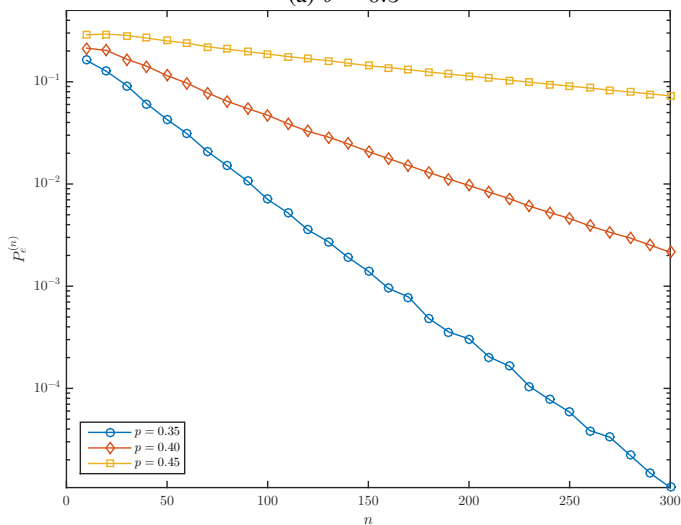
In this paper, we have analyzed the asymptotic performance of majority sentiment detection in online social networks, and revealed that the detection error probability of a majority sentiment detector is strongly related to the network connections. In the i.i.d. case, where users are not connected, the error probability never reduces to zero, regardless of how large the user base is. This result is interesting in its own right, since one would naively expect the error probability to decay to zero with increasing number of users. Furthermore, by modeling the underlying social network as an Ising Markov random field prior, we discovered an interesting phenomenon of phase transition: in the complete graph case, which is an example of a highly connected network, there exists a critical connection strength. If the strength is below this critical value, the error probability is asymptotically positive, while above this critical value, the error probability tends to zero as the number of users increases. This phase transition seen in the complete graph case is analogous to the critical temperature in statistical physics, which separates the paramagnetic phase, where atom magnetic spins are disordered, from the ferromagnetic phase, where the spins are predominantly in one direction resulting in a magnet. We remark that this phenomenon in the OSN model is not due to the type of detector used - but rather is due to the inherent similarity of opinions among users produced by the network model.

REFERENCES

- [1] Philip N Howard and Muzammil M Hussain. *Democracy's fourth wave?: digital media and the Arab Spring*. Oxford University Press on Demand, 2013.
- [2] Efthymios Kouloumpis, Theresa Wilson, and Johanna D Moore. Twitter sentiment analysis: The good the bad and the omg! *Icwsn*, 11:538–541, 2011.



(a) $\theta = 0.3$



(b) $\theta = 0.7$

Fig. 4: $P_e^{(n)}$ v.s. n in complete graph case.

- [3] Mike Thelwall, Kevan Buckley, and Georgios Paltoglou. Sentiment in twitter events. *Journal of the American Society for Information Science and Technology*, 62(2):406–418, 2011.
- [4] Walter Lippmann. *Public opinion*. Transaction Publishers, 1946.
- [5] Peter J Coughlin. *Probabilistic voting theory*. Cambridge University Press, 1992.
- [6] Bo Pang and Lillian Lee. Opinion mining and sentiment analysis. *Foundations and trends in information retrieval*, 2(1-2):1–135, 2008.
- [7] Rohit Negi, Vinay Uday Prabhu, and Miguel Rodrigues. Latent sentiment detection in online social networks: A communications-oriented view. In *Communications (ICC), 2014 IEEE International Conference on*, pages 3758–3763. IEEE, 2014.
- [8] K Binder. Ising model. *Hazewinkel, Michiel, Encyclopedia of Mathematics*. Springer, 1, 2001.
- [9] Richard M Dudley. *Real analysis and probability*, volume 74. Cambridge University Press, 2002.
- [10] P Pfeuty. An exact result for the 1d random ising model in a transverse field. *Physics Letters A*, 72(3):245–246, 1979.
- [11] Martin Kochmański, Tadeusz Paszkiewicz, and Sławomir Wolski. Curie–weiss magnet? a simple model of phase transition. *European Journal of Physics*, 34(6):1555, 2013.
- [12] George Pólya and Gabor Szegő. *Problems and Theorems in Analysis II: Theory of Functions. Zeros. Polynomials. Determinants. Number Theory. Geometry*. Springer Science & Business Media, 1997.