

Social Sampling

Marcel Mohler
mohlerm@student.ethz.ch

March 22, 2015

Abstract

This report provides the reader with an overview of the concept of *social sampling*. Social sampling is a class of methods where informants in a poll answer with a summary of their friends cumulative responses.

1 Introduction

Polling is a widely used method to achieve a public opinion on a certain subject gained from a sample. For example, one could think of polls before *Bundesrat* elections to estimate the outcome or polls to receive student's opinions on their visited lectures. Following the first example, pollers would like to ask only a subset of the population while still get significant results which party might win the vote. This is commonly known as sampling. The two fundamental quantities of interest in polling are *number of samples* (how many citizens to ask) and the *error rate* (error caused by observing a sample instead of the whole voters). They are related in a way, that if we wish to approximate a fraction within an additive error of ϵ , roughly $O(1/\epsilon^2)$ samples are both necessary and sufficient **[[TODO: cite]]**.

In practice we want to have a low error but still only require as little samples as possible. Since the sampling costs are usually a significant burden, research has proposed several alternatives to reduce the sample size. For instance instead of sampling uniformly one could choose their samples with a built in bias. However this method is vulnerable to introduce a so called *systematic bias* which can lead to systematic errors in the outcome.

Another, recently popular approach is to use "expectation polling", where voters are asked about their expectation about the outcome of the poll in contrast to the more classical "intent polling" (where voters are asked about their polling intent). In [3] David Rothschild and Justin Wolvers explored the value of expectation polling. Considering the focus of press and pollsters on intent polling they call it conventional wisdom that these type of polls are more accurate. However, executed on the example of Presidential Electoral College races they provide robust evidence that expectation based polling

yields to more accurate prediction of election outcomes. Unfortunately they fall short on providing any theoretical guarantees on either the sampling bias or sampling error.

In this paper we present an extended idea of "expectation polling" called "social sampling", suggested and described by Anirban Dasgupta, Ravi Kumar and D Sivakumar in [1]. Their idea is to reduce the sampling size by asking a member of a social network about the opinion of their friends. Even though this means that the actual structure of the network plays a role when analysing the error, they assume a saving in number of samples by a factor of around d , the average amount of friends of a member in a social network.

What follows are various social sampling strategies with the aim to observe sampling bias and sampling error. We will also provide precise characterizations of these errors and how they behave applied to large, real-world graphs.

2 Algorithms

Note that graph $G = (V, E)$ describes a social network with V denoting the set of nodes and E the set of edges. Further, $|V| = n$, $|E| = m$ and u, v are nodes of G . The set of neighbors of u is denoted by $N(u)$ and the degree of node u by $d_u = |N(u)|$. We will always assume $d_u \geq 1 \forall u$. Additionally, $f : V \rightarrow \{0, 1\}$ is a binary function which takes a node as input and outputs whether this node satisfies a certain property or not and $\bar{f} = \frac{1}{n} |\{u \mid f(u) = 1\}|$ is the fraction we wish to estimate.

Our goal is to estimate the portion of nodes $v \in C$ such that $f(v) = 1$. We introduce the concept of a set of algorithms called sampler which is defined as follows:

Definition 2.1 (sampler) *A sampler is a randomized algorithm with input r (sample size), ϵ (sampler accuracy), δ (sampler error) and f and outputs $\hat{f} = \frac{1}{r} \sum_{u \in S} f(u)$ with probability $1 - \delta$ and $|\hat{f} - \bar{f}| < \epsilon$*

*We will start by looking at an intuitive approach for estimating \bar{f} where we sample a set of nodes and poll each node u to test if $f(u) = 1$. The algorithm will return the fraction of nodes that satisfy the condition and we call this approach the **Naive** estimator [2].*

Definition 2.2 (Naive estimator) **Definition 2.3 (Ideal estimator)**

3 Results

References

- [1] A. Dasgupta, R. Kumar, and D. Sivakumar. Social sampling. In Proceedings of the 18th ACM SIGKDD international conference on Knowl-

edge discovery and data mining, *pages 235–243. ACM, 2012.*

[2] *O. Goldreich. A sample of samplers: A computational perspective on sampling. def, 1:2n, 1997.*

[3] *D. Rothschild. Forecasting elections comparing prediction markets, polls, and their biases. Public Opinion Quarterly, 73(5):895–916, 2009.*