

Gossip-based Truth Discovery

Zhiying Xu

May 7, 2017

1 Objective Function

The objective function is based on [1]. However, we adopt a discretized truth vector here to avoid in complicated EM algorithm.

Suppose there are n sensors in an arbitrary connected sensor network and each sensor is asked to answer m questions. The questions for each sensor are same and each answer can only be 0 or 1. $\mathbf{X}^{n \times m} \in \{0, 1\}^{n \times m}$ is the observed matrix of the answer from the whole network. In \mathbf{X} , x_{ij} denote the answer of question j from sensor i , \mathbf{x}_i is the i th row of the \mathbf{X} , and \mathbf{x}_j is the j th column of the \mathbf{X} . $\mathbf{t}^m \in \{0, 1\}^m$ is the truth vector, the true answers of the m questions. $\mathbf{r}^n \in [0, 1]^n$ is the reliability vector, the possibility of telling truth of each sensor.

We try to discover the truth vector through observing the \mathbf{X} , that is,

$$\max p(\mathbf{X}|\mathbf{t}, \mathbf{r}).$$

Here \mathbf{t} and \mathbf{r} are related to each other, we further simplify the object function.

For question j , we denote sensors that observe 0 the set S_{j0} and those that observe 1 the set S_{j1} .

$$p(\mathbf{x}_j|\mathbf{r}, t_j = 0) = \prod_{i \in S_{j0}} r_i \prod_{j \in S_{j1}} (1 - r_i)$$

$$p(\mathbf{x}_j|\mathbf{r}, t_j = 1) = \prod_{i \in S_{j0}} (1 - r_i) \prod_{j \in S_{j1}} r_i$$

We decide the value of t_i by comparing $p(\mathbf{x}_j|t_i = 0)$ and $p(\mathbf{x}_j|t_i = 1)$. Taking all the m into consideration, we have

$$p(\mathbf{X}|\mathbf{t}, \mathbf{r}) = p(\mathbf{X}|\mathbf{r}) = \prod_{j=1}^m \max \left\{ \prod_{i \in S_{j0}} r_j \prod_{i \in S_{j1}} (1 - r_i), \prod_{i \in S_{j0}} (1 - r_i) \prod_{i \in S_{j1}} r_i \right\}$$

$$\ln p(\mathbf{X}|\mathbf{r}) = \sum_{j=1}^m \max \left\{ \sum_{i \in S_{j0}} \ln r_i + \sum_{i \in S_{j1}} \ln(1 - r_i), \sum_{i \in S_{j0}} \ln(1 - r_i) + \sum_{i \in S_{j1}} \ln r_i \right\}$$

$$= \sum_{j=1}^m \sum_{i=1}^n \frac{1}{2} \left(\sum_{i=1}^n \ln r_i + \sum_{i=1}^n \ln(1 - r_i) + \left| \sum_{i \in S_{j0}} \ln r_i + \sum_{i \in S_{j1}} \ln(1 - r_i) - \sum_{i \in S_{j0}} \ln(1 - r_i) - \sum_{i \in S_{j1}} \ln r_i \right| \right)$$

$$= \frac{1}{2} \sum_{j=1}^m \left(\sum_{i=1}^n \ln r_i (1 - r_i) + \left| \sum_{i=1}^n x_{ij} \ln \frac{r_i}{1 - r_i} \right| \right).$$

We get reliability vector \mathbf{r} from $\arg \max_{\mathbf{r}} \ln p(\mathbf{X}|\mathbf{r})$. Traditionally, to maximize $p(\mathbf{X}|\mathbf{t}, \mathbf{r})$, an EM algorithm is required. However, there is no need for iteration here. The reason is that future \mathbf{r} won't change after we updating \mathbf{t} based on current \mathbf{r} . Specifically, the result of r_i base on \mathbf{t} is $\frac{\|\mathbf{x}_i - \mathbf{t}\|_1}{m}$.

In this way, the object function can also be written as

$$\arg \max_{\mathbf{t}} \sum_{i=1}^n (d_i \ln d_i + (1 - d_i) \ln(1 - d_i)), \text{ where } d_i = \frac{1}{m} \|\mathbf{x}_i - \mathbf{t}\|_1.$$

2 Optimal Solution

Theorem 1. *The optimal solution depends on the rank of the observed matrix.*

Proof. content...

□

3 NP-hardness

Theorem 2. *Finding the truth vector is NP-hard.*

Proof. We prove the NP-hardness through reduction from exact 3-cover problem. The exact 3-cover asks, given set U and $S \subseteq \binom{U}{3}$, to decide if there exists $S' \subseteq S$, where S' is a partition of U .

Our proof can be divided into 3 parts. (1) We construct a graph based on U and S and define function on the nodes. (2) We prove minimizing the sum of function of all nodes is NP-hard. (2) We derive a matrix from the graph, and the value of object function of the matrix equals to the sum of function in all nodes.

(1) We construct a graph $G = (V, E)$ same as graph described in [3]. We create a vertex s_i per set S_i and a copy of gadgets per elements u_j . We link $u_{j,k}$ into $S_{j,k}$, where j_1, j_2, j_3 are the indices of the three sets containing u_j .

We transform the undirected graph into a directed graph. The total degree of node v is denoted as $d(v)$. The in-degree and out-degree of the node v is denoted $d_{in}(v)$ and $d_{out}(v)$. We define a function for v .

$$f(v) = -()$$

□

References

- [1] Wang, Dong, et al. "On truth discovery in social sensing: A maximum likelihood estimation approach." Information Processing in Sensor Networks (IPSN), 2012 ACM/IEEE 11th International Conference on. IEEE, 2012.
- [2] Boyd, Stephen, et al. "Randomized gossip algorithms." IEEE/ACM Transactions on Networking (TON) 14.SI (2006): 2508-2530.
- [3] Cardinal, Jean, Samuel Fiorini, and Gwenal Joret. "Minimum entropy orientations." Operations research letters 36.6 (2008): 680-683.