

# Higher Order Group Prediction

## ABSTRACT

Aim is to predict higher order groups given a previous group interaction history.

## General Terms

Hyperedge, Hypergraph

## Keywords

Hyperedge Prediction, Hypergraphs, Convex Coptimization

## 1. PROBLEM STATEMENT

Our input is the history of previous collaborations,  $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{(T-1)}\}$ . Here, each  $\mathbf{y}_t = \{v_1, v_2, \dots, v_k\}$  represents a particular hyperedge (group or collaboration) that occurs at time instance  $t$  with the given vertices each represented by  $v_i$  (total number of vertices being  $n$ ). Our aim is to predict the possible new hyperedge  $\mathbf{y}_{\text{new}}$  that occurs at  $t = T$ .

## 2. APPROACH 1

We treat this problem as an optimization task of finding the vector  $\mathbf{y}_{\text{new}}$  which is nearest to all the given previous collaboration vectors  $\mathbf{y}_i$ ,  $i \in \{1, 2, \dots, (T-1)\}$  and is also sparse. We define the following regularization framework:

$$\arg \min_{\mathbf{y}_{\text{new}}} \text{Dist}(\mathbf{Y}, \mathbf{y}_{\text{new}}) + \lambda \|\mathbf{y}_{\text{new}}\|_1 \quad (1)$$

where  $\text{Dist}(\mathbf{Y}, \mathbf{y}_{\text{new}})$  represents the distance of different kinds like euclidean distance  $\|\mathbf{Y} - \mathbf{1}\mathbf{y}_{\text{new}}^T\|_2$ , hamming distance etc. and  $\mathbf{1}$  is a vector of all ones. The second term is lasso constraint for sparsity. The parameter  $\lambda$  is used for tuning the extent of regularization.

## 3. APPROACH 2: PREDICTING SET OF GROUPS

We treat this problem as an optimization task of finding the stack of vectors  $\mathbf{Y}_{\text{new}} = \{\mathbf{y}_{\text{new}}^1, \dots, \mathbf{y}_{\text{new}}^P\}$  which is

nearest to all the given previous collaboration vectors  $\mathbf{y}_i$ ,  $i \in \{1, 2, \dots, (T-1)\}$  and is also sparse. We define the following regularization framework:

$$\arg \min_{\mathbf{Y}_{\text{new}}} \sum_{p=1}^P \sum_{t=1}^{T-1} \text{Dist}(\mathbf{y}_t, \mathbf{y}_{\text{new}}^p) + \lambda \sum_{p=1}^P \|\mathbf{y}_{\text{new}}^p\|_1 \quad (2)$$

where  $(\mathbf{y}_t, \mathbf{y}_{\text{new}}^p)$  represents the distance of different kinds like euclidean distance  $\|\mathbf{y}_t - \mathbf{y}_{\text{new}}^p\|_2$ , hamming distance etc. The second term is lasso constraint for sparsity. The parameter  $\lambda$  is used for tuning the extent of regularization.

## 4. APPROACH 3

**Approach 1** is static in the sense that it does not take into account the time dimension. Moreover, **Approach 1** assumes an oversimplified and over demanding in the sense that it wants to find out new groups which are similar to the whole stack of groups observed in past. Rather, a more realistic approach should consider vectors (groups or hyperedges) that are similar to some subset of previously observed groups. Moreover, it is safe to assume that these subsets are the various communities or clusters of related hyperedges observed in past. We therefore, represent the previous collaboration vectors  $\mathbf{y}_{ic}$ ,  $i \in \{1, 2, \dots, (T-1)\}$  and  $c \in \{1, 2, \dots, N_c\}$  ( $N_c$  being the total number of communities observed in past). We define the following regularization framework:

$$\arg \min_{\mathbf{y}_{\text{new}}^p} \sum_{t=1}^{T-1} \gamma^{-t} \text{Dist}(\mathbf{y}_{tc}, \mathbf{y}_{\text{new}}^p) + \lambda \|\mathbf{y}_{\text{new}}^p\|_1 \quad (3)$$

where  $\text{Dist}(\mathbf{y}_{tc}, \mathbf{y}_{\text{new}}^p)$  represents the distance of different kinds like euclidean distance  $\|\mathbf{y}_{tc} - \mathbf{1}(\mathbf{y}_{\text{new}}^p)^T\|_2$ , hamming distance etc,  $p$  is the index of the new vectors predicted for each community, and  $\mathbf{1}$  is a vector of all ones. The second term is lasso constraint for sparsity. The parameter  $\lambda$  is used for tuning the extent of regularization and parameter  $\gamma$  penalizes more if the predicted hyperedge is not similar to hyperedges in recent past. Community  $c$  for  $\mathbf{y}_{tc}$  is decided using any clustering methods like KNN.

## 5. ISSUES

Major issues with the above approaches in general is that  $\mathbf{y}_{\text{new}}$  can be predicted same as one of the training  $\mathbf{y}_{tc}$  i.e. there is no penalization for exact similarity.

- Either we penalize it directly, but how ?,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WOODSTOCK '97 El Paso, Texas USA

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

- Or we make constraints such that  $\mathbf{Dist}(\mathbf{y}_{tc}, \mathbf{y}_{new}^p) > \delta$  where  $\delta$  is the minimum diversity from the existing hyperedges in community  $c$ .
- Or we can measure distance from the cluster centers rather than the  $\mathbf{y}_{tc}$  directly. More generally we can take distance from a dictionary representing  $\mathbf{Y}$ .

## 6. ADDITIONAL AUTHORS

## 7. REFERENCES

- [1] M. Bowman, S. K. Debray, and L. L. Peterson. Reasoning about naming systems. *ACM Trans. Program. Lang. Syst.*, 15(5):795–825, November 1993.
- [2] J. Braams. Babel, a multilingual style-option system for use with latex’s standard document styles. *TUGboat*, 12(2):291–301, June 1991.
- [3] M. Clark. Post congress tristesse. In *TeX90 Conference Proceedings*, pages 84–89. TeX Users Group, March 1991.
- [4] M. Herlihy. A methodology for implementing highly concurrent data objects. *ACM Trans. Program. Lang. Syst.*, 15(5):745–770, November 1993.
- [5] L. Lamport. *LaTeX User’s Guide and Document Reference Manual*. Addison-Wesley Publishing Company, Reading, Massachusetts, 1986.
- [6] S. Salas and E. Hille. *Calculus: One and Several Variable*. John Wiley and Sons, New York, 1978.