

Appendix: Quantifying and Reducing Imbalance in Networks

Yoosof Mashayekhi

Bo Kang

yoosof.mashayekhi@ugent.be

bo.kang@ugent.be

Ghent University

Ghent, Belgium

Jefrey Lijffijt

Tijl De Bie

jefrey.lijffijt@ugent.be

tijl.debie@ugent.be

Ghent University

Ghent, Belgium

A NETWORK IMBALANCE IS EQUIVALENT TO EMD

Earth Mover's Distance (EMD) is a measure of the distance between two distributions. In this case, EMD computes the minimum amount of work (the distance of movement of nodes in the embedding space) to match the embedding of the source nodes and the target nodes. We show that EMD computes the same value as the network imbalance ψ . To show that, we first introduce EMD.

EMD was invented to solve certain kinds of transportation problems [1]. The transportation problem is a particular type of linear programming where the objective is to minimize the cost of transporting any commodity from one group of sources to another group of destinations. Formally:

DEFINITION 1 (EARTH MOVER'S DISTANCE [1]). Assume two distributions represented by signatures $P = \{(\mathbf{p}_1, w_{p1}), \dots, (\mathbf{p}_m, w_{pm})\}$ and $Q = \{(\mathbf{q}_1, w_{q1}), \dots, (\mathbf{q}_r, w_{qr})\}$, where \mathbf{p}_i and \mathbf{q}_j are cluster representatives and w_{pi} and w_{qj} are weights of clusters. Let $D = [d_{ij}]$ be the distance matrix between P and Q . EMD is a linear program whose goal is to find a flow $F = [f_{ij}]$ between two distributions P and Q that minimizes the overall cost:

$$C = \sum_{i=1}^m \sum_{j=1}^r f_{ij} d_{ij}$$

$$\text{s.t. } f_{ij} \geq 0, \quad 1 \leq i \leq m, \quad 1 \leq j \leq r,$$

$$\sum_{j=1}^r f_{ij} \leq w_{pi}, \quad 1 \leq i \leq m, \quad (1)$$

$$\sum_{i=1}^m f_{ij} \leq w_{qj}, \quad 1 \leq j \leq r, \quad (2)$$

$$\sum_{i=1}^m \sum_{j=1}^r f_{ij} = \min \left\{ \sum_{i=1}^m w_{pi}, \sum_{j=1}^r w_{qj} \right\}. \quad (3)$$

The optimal flow F is found by solving this linear optimization problem. D_{EMD} is defined as the work normalized by the total flow:

$$D_{EMD}(D, P, Q) = \frac{\sum_{i=1}^m \sum_{j=1}^r f_{ij} d_{ij}}{\sum_{i=1}^m \sum_{j=1}^r f_{ij}}.$$

EMD, which is a cross-bin distance measure, computes the distance between two distributions (source nodes and target nodes

in our case). Note that bin-bin distance measures such as Kullback–Leibler divergence do not take the distance between values into account.

We use D_{EMD} to compute imbalance in a network and we refer to it as ψ_{EMD} . We assign equal weights to each node in S and also equal weights to each node in T in a way that the sum of the weight of nodes in S and T are equal. Formally:

DEFINITION 2 (IMBALANCE MEASURE ψ_{EMD}). Given a network $G = (V, E)$, two disjoint sets of nodes, namely source nodes $\emptyset \subset S \subset V$ and target nodes $\emptyset \subset T \subset V$, and the cost of matching each pair of nodes $D = [d_{ij}]$, we define two signatures $C_S = \{(i, w_i = \frac{1}{|S|}) | i \in S\}$ and $C_T = \{(j, w_j = \frac{1}{|T|}) | j \in T\}$. We define imbalance measure ψ_{EMD} :

$$\psi_{EMD}(D, S, T) = D_{EMD}(D, C_S, C_T).$$

PROPOSITION 1 (EQUIVALENCE OF ψ_{EMD} AND ψ). ψ_{EMD} equals to ψ .

PROOF. We show that ψ_{EMD} solves the same problem as ψ . Since $w_i = \frac{1}{|S|}, i \in S$ and $w_j = \frac{1}{|T|}, j \in T$ (Definition 2), $\sum_{i \in S} w_i = \sum_{j \in T} w_j = 1$. Since constraint 3 becomes $\sum_{i \in S} \sum_{j \in T} f_{ij} = \sum_{i \in S} w_i = \sum_{j \in T} w_j = 1$, constraints 1 and 2 also change to $\sum_{j \in T} f_{ij} = w_i, i \in S$ and $\sum_{i \in S} f_{ij} = w_j, j \in T$. Hence, we can rewrite ψ_{EMD} as the linear program:

$$C = \sum_{i \in S} \sum_{j \in T} f_{ij} d_{ij}$$

$$\text{s.t. } f_{ij} \geq 0 \quad \forall (i, j) \in S \times T,$$

$$\sum_{j \in T} f_{ij} = w_i \quad \forall i \in S, \quad \sum_{i \in S} f_{ij} = w_j \quad \forall j \in T,$$

which solves the same problem as ψ . \square

B GRAB ALGORITHM

Algorithm ?? shows the greedy link selection in Grab. Algorithm 2 the complete generic method Grab to select k links to add to the network.

REFERENCES

- [1] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. 1998. A metric for distributions with applications to image databases. In *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*. IEEE, 59–66.

Algorithm 2: Graph Balancing (GraB)

Input: NE (network embedding method), A (adjacency matrix), S (source nodes), T (target nodes), U (auxiliary nodes), k (number of links to add), b (batch size), l (batch size coefficient)

Output: $links$ (k links)

Function GraB(NE, A, S, T, U, k, b, l):

```

     $links = []$ ,  $links\_idx = 1$ 
     $X = NE.Embeddings(A)$ 
    while  $links\_idx \leq k$  do
         $candidate\_list = SelectCandidateLinks(NE, A, S, T,$ 
             $U, b \cdot l)$ 
         $X\_new, A\_new = NE.ReEmbed(X, A,$ 
             $candidate\_list)$  // Re-embed the network after adding
            the  $candidate\_list$  links
         $current\_batch\_list = [], idx = 1$ 
        foreach  $(i, j)$  in  $candidate\_list$  do
            if  $\delta(x_i, X, S, T) < \delta(x\_new_i, X\_new, S, T)$ 
                then
                     $links[links\_idx] = (i, j)$ 
                     $links\_idx += 1$ 
                     $current\_batch\_list[idx] = (i, j)$ 
                     $idx += 1$ 
                    if  $idx > b$  then
                        break
                    end
            end
        end
         $X, A = NE.ReEmbed(X, A, current\_batch\_list)$ 
        // Re-embed the network after adding the
         $current\_batch\_list$  links
    end
    return  $links$ 

```