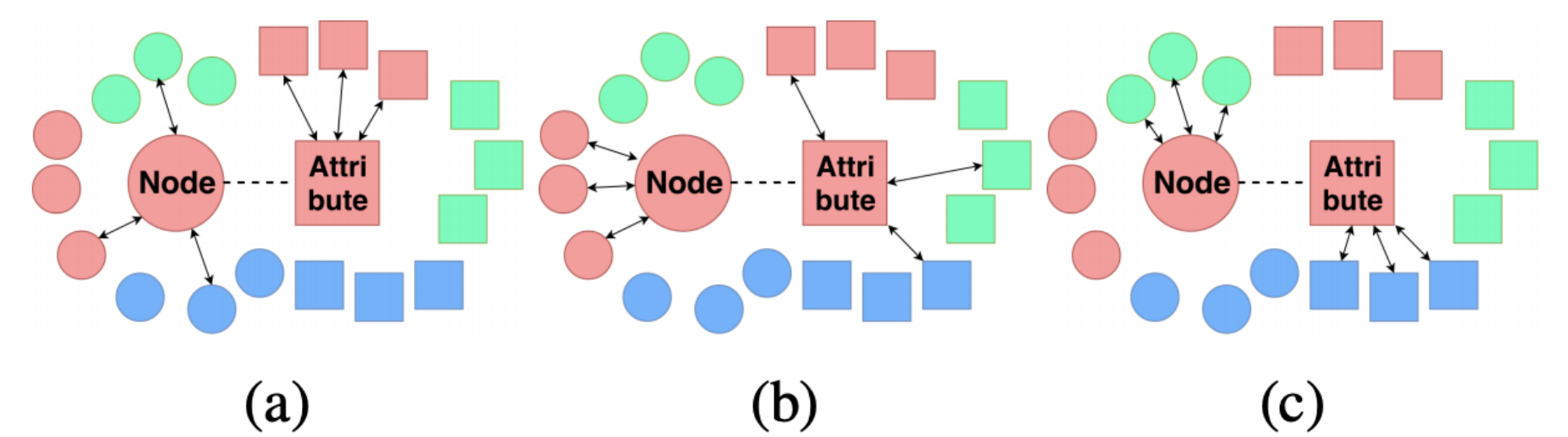


Introduction and Motivation

- In an attributed network, each node comes with some content within (also known as **node attributes**).
- Existing attributed network approaches work well when the network is **consistent between structure and attributes**.
- Real world networks often have anomalous (**outlier**) nodes, which affect the embeddings of other nodes in the network. Thus all the downstream network mining tasks may fail miserably in the presence of such outliers.
- We proposed an **unsupervised Outlier aware Network Embedding algorithm (ONE)** for attributed networks, which minimizes the effect of the outlier nodes, and hence generates robust network embeddings.
- To the best of our knowledge, this is the first work to propose a completely unsupervised algorithm for attributed network embedding integrated with outlier detection.**



- (a) **Structural Outlier**: The node has edges to nodes from different communities.
- (b) **Attribute Outlier**: The attributes of the node are similar to attributes of the nodes from different communities.
- (c) **Combined Outlier**: Node belongs to a community structurally but it has a different community in terms of attribute similarity.

Problem Formulation and Solution Approach

- An attributed information network is represented by a graph as $\mathcal{G} = (V, E, C)$
- $V = \{v_1, v_2, \dots, v_N\}$ is the set of nodes, $E \subset \{(v_i, v_j) | v_i, v_j \in V\}$ is the set of edges.
- $N \times N$ dimensional adjacency matrix of the graph \mathcal{G} is $A = (a_{ij})$, where $a_{ij} = w_{v_i, v_j}$ if $(v_i, v_j) \in E$, and $a_{ij} = 0$ otherwise.
- C is a $N \times D$ matrix where $C_i \in \mathbb{R}^D$ is the attribute vector associated with the node $v_i \in V$.
- C_{id} is the value of the attribute d for the node v_i . For example, if there is only textual content in each node, c_i can be the tf-idf vector for the content of the node v_i .
- For a given network \mathcal{G} , network embedding is a technique to learn a function $f: v_i \mapsto y_i \in \mathbb{R}^K$, i.e., it maps every vertex to a K dimensional vector, where $K < \min(N, D)$.
- The representations should preserve the underlying semantics of the network.

- Learning from the Link Structure:

$$\mathcal{L}_{str} = \sum_{i=1}^N \sum_{j=1}^N \log \left(\frac{1}{O_{1i}} \right) (A_{ij} - G_i \cdot H_j)^2 \quad (1)$$

- Learning from the Attributes:

$$\mathcal{L}_{attr} = \sum_{i=1}^N \sum_{d=1}^D \log \left(\frac{1}{O_{2i}} \right) (C_{id} - U_i \cdot V_d)^2 \quad (2)$$

- Connecting Structure and Attributes: We want to find a matrix W which minimizes $\|G - WU\|_F$.

$$\mathcal{L}_{dis} = \sum_{i=1}^N \sum_{k=1}^K \log \left(\frac{1}{O_{3i}} \right) (G_{ik} - U_i \cdot (W^T)_k)^2 \quad (3)$$

- If we restrict W to be an orthogonal matrix, then a **closed form solution** can be obtained from the solution concept of Procrustes problem:

$$W^* = \underset{W \in \mathcal{O}_K}{\operatorname{argmin}} \|G - UW^T\|_F \quad (4)$$

where $W^* = XY^T$ with $XY^T = \operatorname{SVD}(G^T U)$, \mathcal{O}_K is the set of all orthogonal matrices of dimension $K \times K$.

- Joint Loss Function:

$$\mathcal{L} = \mathcal{L}_{str} + \alpha \mathcal{L}_{attr} + \beta \mathcal{L}_{dis}$$

$$0 < O_{1i}, O_{2i}, O_{3i} \leq 1, \forall v_i \in V$$

$$\sum_{i=1}^N O_{1i} = \sum_{i=1}^N O_{2i} = \sum_{i=1}^N O_{3i} = \mu$$

$$W \in \mathcal{O}_K \iff W^T W = \mathcal{I}$$

Update Rules and Experimental Setup

$$G_{ik} = \frac{G_{ik}^{num1} + \beta \log \left(\frac{1}{O_{3i}} \right) (W_k \cdot U_i)}{\log \left(\frac{1}{O_{1i}} \right) (H_k \cdot H_k) + \beta \log \left(\frac{1}{O_{3i}} \right)}$$

$$G_{ik}^{num1} = \log \left(\frac{1}{O_{1i}} \right) \sum_{j=1}^N (A_{ij} - \sum_{k' \neq k} G_{ik'} H_{k'j}) H_{kj}$$

$$H_{kj} = \frac{\sum_{i=1}^N \log \left(\frac{1}{O_{1i}} \right) (A_{ij} - \sum_{k' \neq k} G_{ik'} H_{k'j}) G_{ik}}{\sum_{i=1}^N \log \left(\frac{1}{O_{1i}} \right) G_{ik}^2}$$

$$O_{1i} = \frac{\left(\sum_{j=1}^N (A_{ij} - G_i \cdot H_j)^2 \right) \cdot \mu}{\sum_{i=1}^N \sum_{j=1}^N (A_{ij} - G_i \cdot H_j)^2}$$

$$O_{2i} = \frac{\left(\sum_{d=1}^D (C_{id} - U_i \cdot V_d)^2 \right) \cdot \mu}{\sum_{i=1}^N \sum_{j=1}^D (C_{id} - U_i \cdot V_d)^2}$$

$$O_{3i} = \frac{\left(\sum_{k=1}^K (G_{ik} - W_i \cdot U_k)^2 \right) \cdot \mu}{\sum_{i=1}^N \sum_{k=1}^K (G_{ik} - W_i \cdot U_k)^2}$$

- To check the performance of the algorithms in the presence of outliers, we manually planted a total of 5% outliers (with equal numbers for each type) in each dataset.

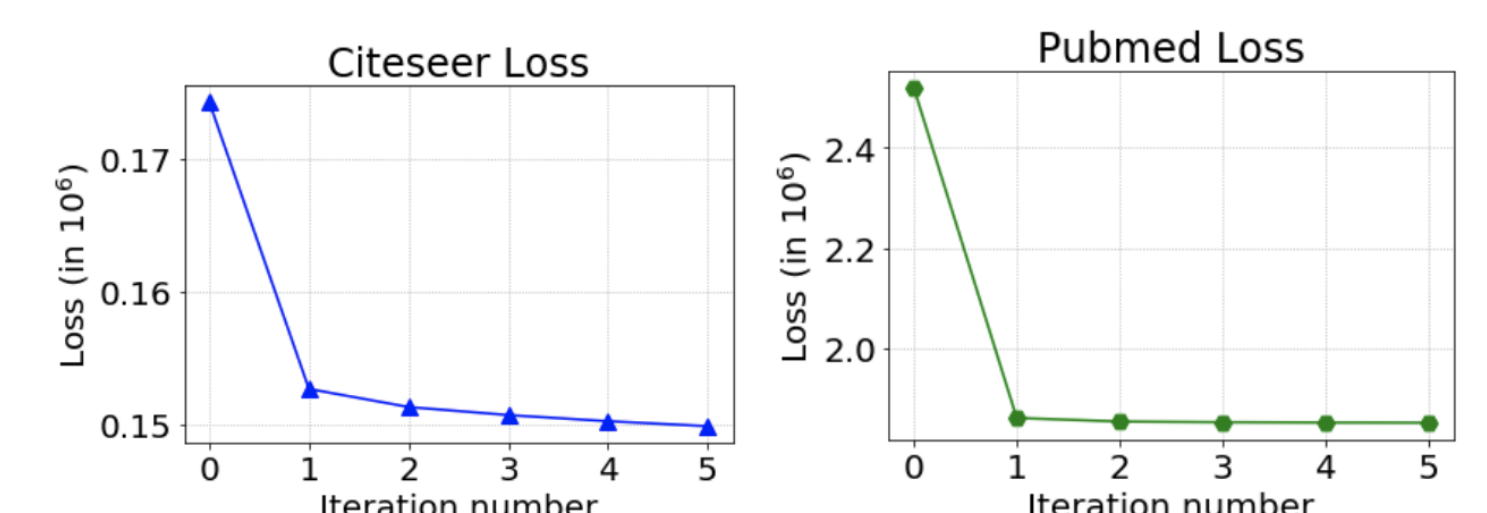


Figure 2: Values of Loss function over different iterations of ONE for Citeseer and Pubmed (seeded) datasets

Lemma 1 The joint cost function decreases after each iteration of the set of update rules.

Experimental Results

- Observations:

- Though SEANO uses 20% labeled data as the extra supervision to generate the embeddings, its accuracy is always close (or less) to ONE which is completely unsupervised in nature.
- ONE converges very fast on real datasets. But updating most of the variables in this framework takes $O(N)$ time, which leads to $O(N^2)$ runtime for ONE without any parallel processing.

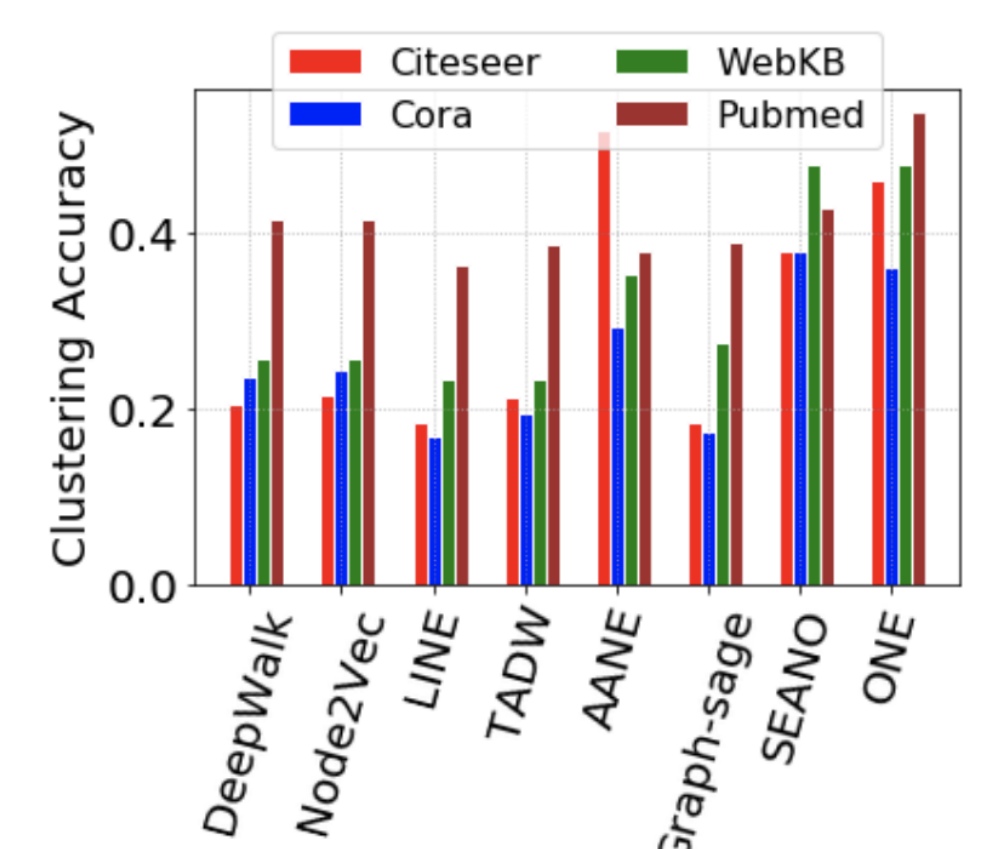


Figure 5: Clustering Accuracy of KMeans++ on the embeddings generated by different algorithms.

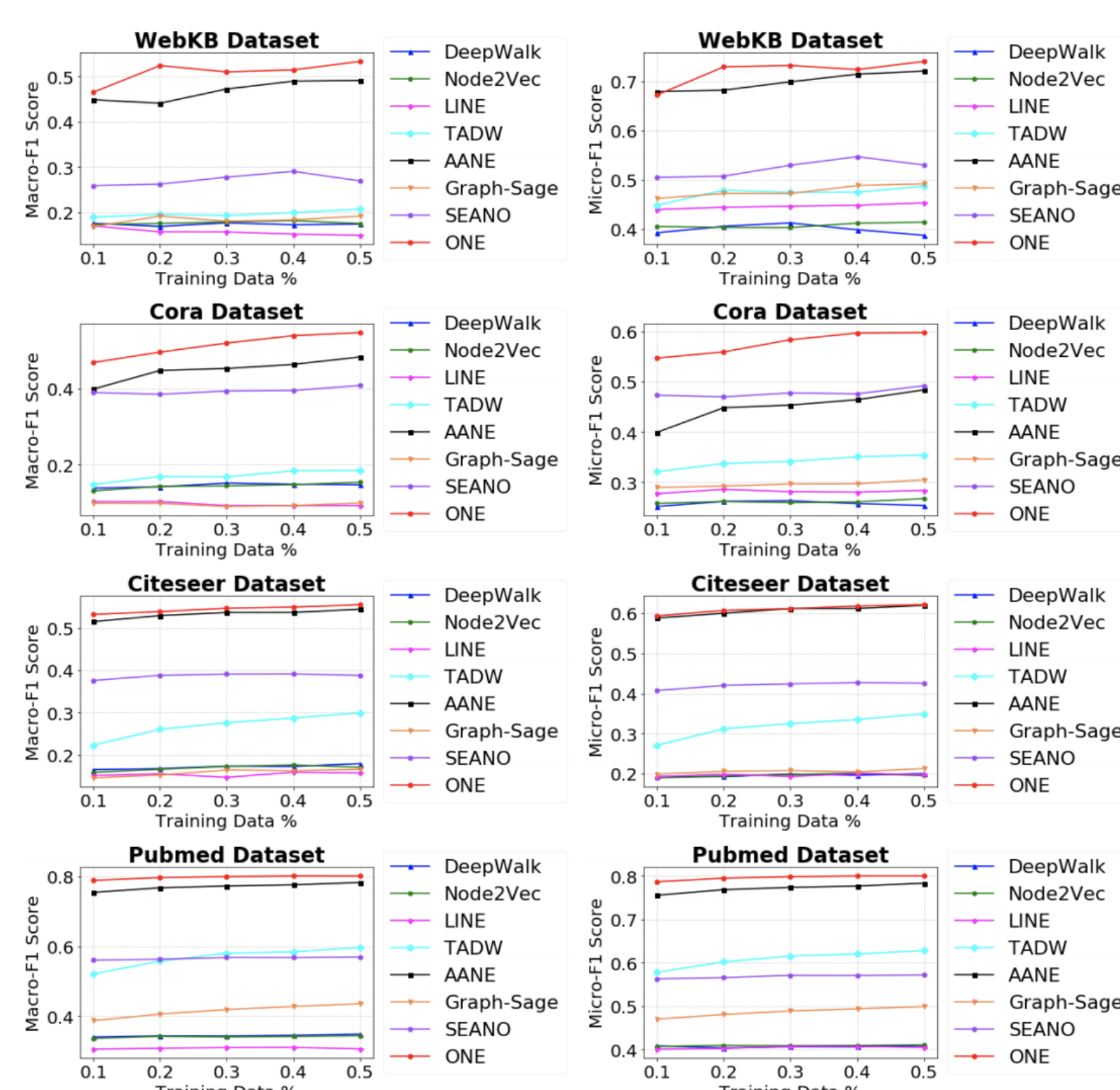


Figure 4: Performance of different embedding algorithms for Classification with Random Forest

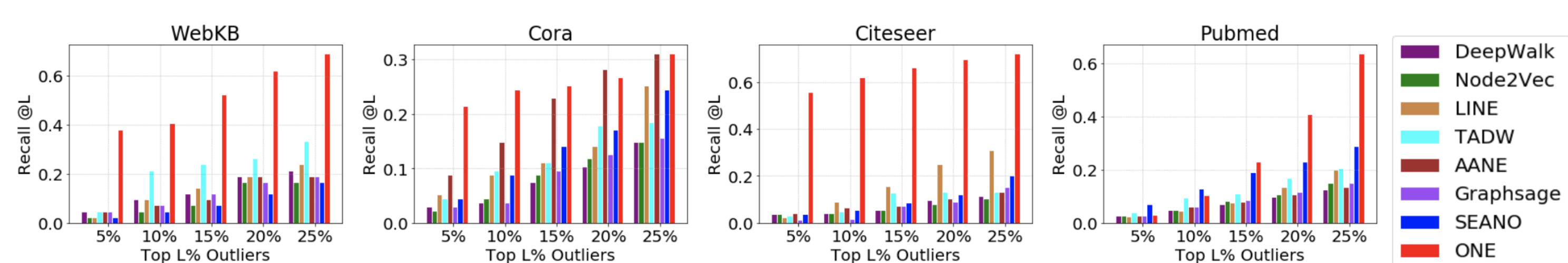


Figure 3: Outlier Recall at top L% from the ranked list of outliers for all the datasets. ONE, though it is an unsupervised algorithm, outperforms all the baseline algorithms in most of the cases. SEANO uses 20% labeled data for training.

