

1. Background

Wasserstein Distance

- Powerful tool to measure the distance between distributions.
- High computational cost. (e.g., linear programming, Sinkhorn algorithm)

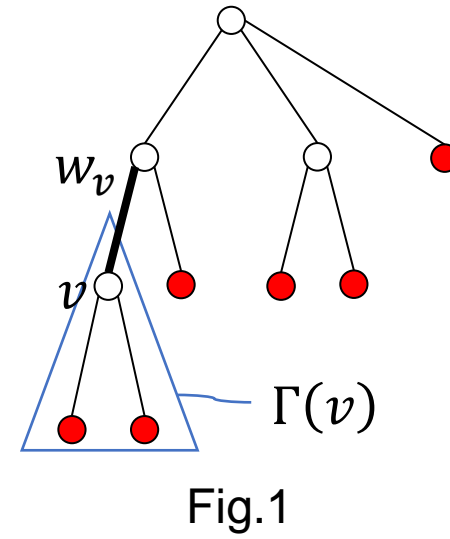
$$W(\mu_i, \mu_j) = \inf_{\gamma \in \Pi(\mu_i, \mu_j)} \int d(x_i, x_j) \gamma(dx, dy)$$

Tree-Wasserstein Distance

- The Wasserstein distance on a tree.
- Closed form solution, which can be computed in linear time.

$$W_{d_T}(\mu_i, \mu_j) = \sum_{v \in V} w_v \left| \sum_{x \in \Gamma(v)} \mu_i(x) - \mu_j(x) \right|_1$$

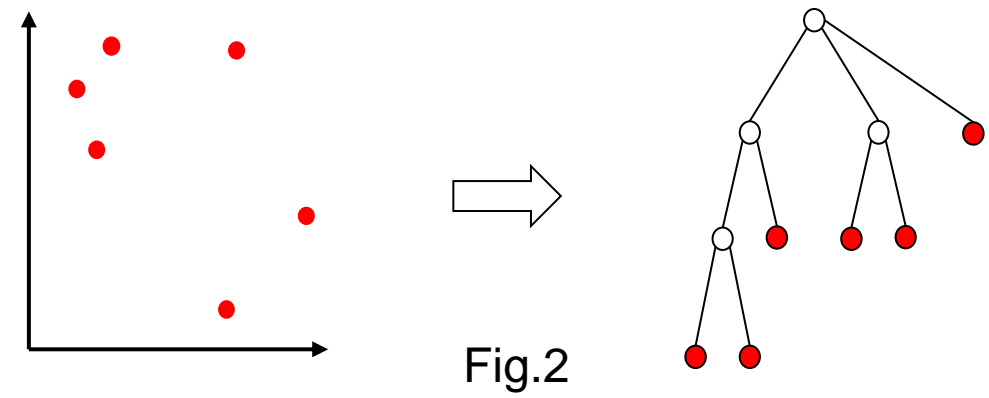
- w_v : the edge length between v and its parent node.
- $\Gamma(v)$: a set of nodes contained in the subtree rooted at v .



Methods to Construct a Tree.

- Quadtree [Indyk & Thaper 03]
- Clustering based method [Le+ 19]

These methods are unsupervised.



2. Contribution Summary

- We propose the **Supervised tree-Wasserstein (STW) distance** to construct a tree that can represent task-specific distances using the label information of documents.
- We propose the **Soft variant of the tree-Wasserstein distance** that is differential w.r.t. parent-child relationships in a tree.
- The STW distance outperforms other unsupervised tree-based methods in document classification tasks.
- Since the STW distance is GPU suitable, it can compute the tree-Wasserstein distance more efficiently.

3. Problem Setting

- Input :**
- $Z = \{z_1, z_2, \dots, z_{N_{\text{leaf}}}\}$: a set of words.
 - μ_i : probability measure of a document i . (i.e., normalized bag-of-words)
 - $y_i \in \mathcal{N}$: a label of document i .
 - $D = \{(\mu_i, y_i)\}_{i=1}^M$: a training dataset.
 - N_{in} : the number of internal nodes. (hyper parameter)

Then, we denote as follows :

- $V = \{v_1, v_2, \dots, v_{N_{\text{in}}+N_{\text{leaf}}}\}$: a set of nodes (v_1 is a root).

Goal :

To obtain the tree metric that can represent task specific distance.

4. Soft Tree-Wasserstein Distance

Difficulty :

- Optimization w.r.t. a tree structure is discrete optimization. (e.g., $\Gamma(v)$)

Soft Tree-Wasserstein Distance :

- Differential w.r.t. parent-child relationships.
- $P_{\text{sub}}(x | v)$: probability that $x \in \Gamma(v)$. (i.e., x is contained the subtree rooted at v .)

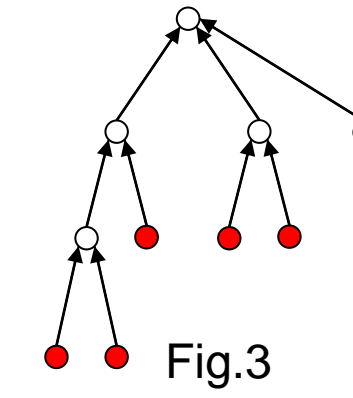
$$W_{d_T}^{\text{soft}}(\mu_i, \mu_j) = \sum_{v \in V} w_v \left| \sum_{x \in V} P_{\text{sub}}(x | v) (\mu_i(x) - \mu_j(x)) \right|_\alpha$$

Theorem 2 :

If the tree metric is given and $\|\cdot\|_\alpha$ approaches $\|\cdot\|_1$, then the soft tree-Wasserstein distance converges to the tree-Wasserstein distance.

Parent-Child Relationships :

- These relationships can be represented by a directed tree.
- i.e., an adjacency matrix D_{par} .



Theorem 1 : conditions of an adjacency matrix

Let $D_{\text{par}} \in \{0, 1\}^{|\mathcal{V}| \times |\mathcal{V}|}$ be an adjacency matrix of the directed graph G . If D_{par} satisfies the followings:

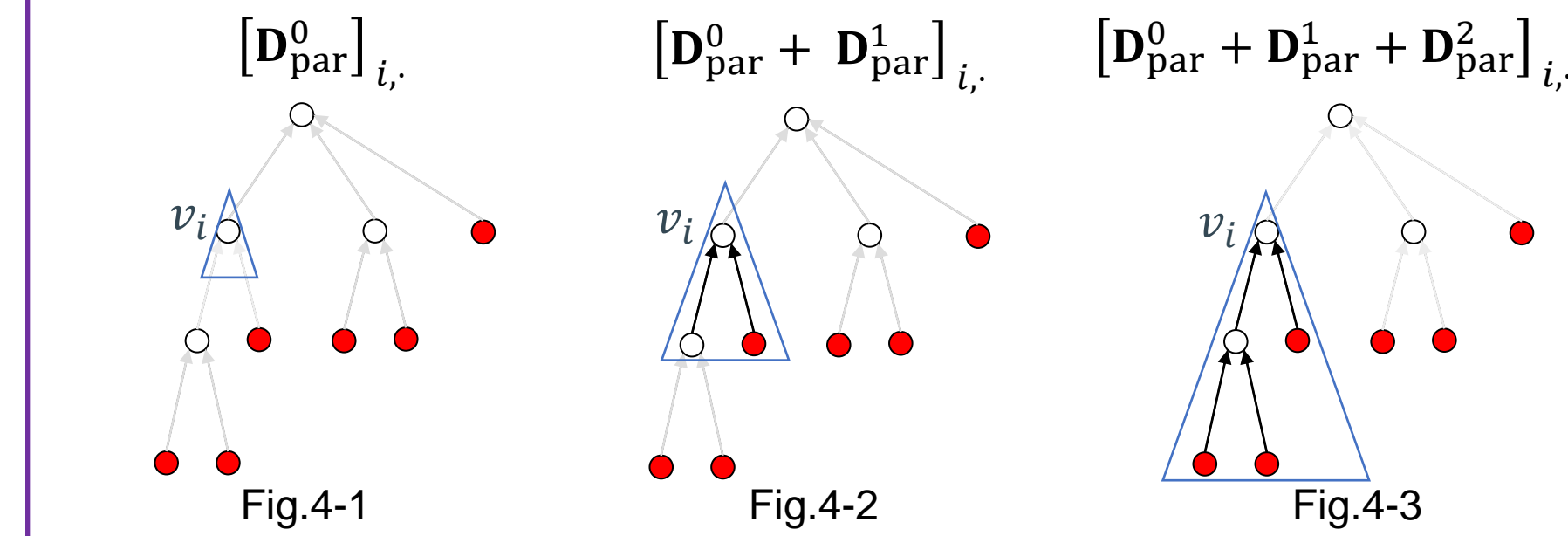
- D_{par} is a strictly upper triangular matrix.
- $D_{\text{par}}^T \mathbf{1} = (0, 1, \dots, 1)^T$.

Then G is a directed tree.

Formulation of $P_{\text{sub}}(x | v)$:

- $[D_{\text{par}}^k]_{i,j}$: the probability that there exists the path from v_j to v_i with exactly k steps.

Example



$$P_{\text{sub}}(v_j | v_i) = \left[\sum_{k=0}^{\infty} D_{\text{par}}^k \right]_{i,j} = \left[(\mathbf{I} - D_{\text{par}})^{-1} \right]_{i,j}$$

Matrix Form Formulation :

- \mathbf{a}_i and \mathbf{a}_j are normalized bag-of-words.
- Batch processing

$$W_{d_T}(\mathbf{a}_i, \mathbf{a}_j) = \left| \mathbf{w}_v \circ (\mathbf{I} - D_{\text{par}})^{-1} \begin{pmatrix} 0 \\ \mathbf{a}_i - \mathbf{a}_j \end{pmatrix} \right|_1$$

5. Supervised Tree-Wasserstein (STW) Distance

Loss Function :

$$L(D_{\text{par}}, \mathbf{w}_v) = \frac{1}{|D_p|} \sum_{(i,j) \in D_p} W_{d_T}^{\text{soft}}(\mu_i, \mu_j) - \frac{1}{|D_n|} \sum_{(i,j) \in D_n} \min\{W_{d_T}^{\text{soft}}(\mu_i, \mu_j), m\}$$

where $D_p = \{(i, j) | y_i = y_j\}$, $D_n = \{(i, j) | y_i \neq y_j\}$ and $\mathbf{w}_v = (w_{v_1}, \dots, w_{v_N})$.

- We can minimize the loss by using the stochastic gradient descent.
- After the optimization, we select the most probable parent node for each node, and then construct a tree.

Techniques :

- We fix the edge length w_v to 1.
- We fix the tree structure of internal nodes.

6. Experimental Results

Comparison Methods :

- Word Mover's Distance (WMD)
 - Supervised WMD
 - Quadtree
 - TSW
 - Flowtree
- Unsupervised tree-Wasserstein distance based methods

Document Classification Accuracy :

- On four datasets, STW outperforms unsupervised tree-based methods.

Table 1 : kNN test error rate.

	TWITTER	AMAZON	CLASSIC	BBCSPORT	OHSUMED	REUTERS
WMD	28.7 ± 0.6	7.4 ± 0.3	2.8 ± 0.1	4.6 ± 0.7	44.5	3.5
S-WMD	27.5 ± 0.5	5.8 ± 0.1	3.2 ± 0.2	2.1 ± 0.5	34.3	3.2
QUADTREE	30.4 ± 0.8	10.7 ± 0.3	4.1 ± 0.4	4.5 ± 0.5	44.0	5.2
FLOWTREE	29.8 ± 0.9	9.9 ± 0.3	5.6 ± 0.6	4.7 ± 1.1	44.4	4.7
TSW-1	30.2 ± 1.3	14.5 ± 0.6	5.5 ± 0.5	12.4 ± 1.9	58.4	7.5
TSW-5	29.5 ± 1.1	9.2 ± 0.1	4.1 ± 0.4	11.9 ± 1.3	51.7	5.8
TSW-10	29.3 ± 1.0	8.9 ± 0.5	4.1 ± 0.6	11.4 ± 0.9	51.1	5.4
STW	28.9 ± 0.7	10.1 ± 0.7	4.4 ± 0.7	3.4 ± 0.8	40.2	4.4

Running Time :

- The Tree-Wasserstein distance is faster than the Wasserstein distance.

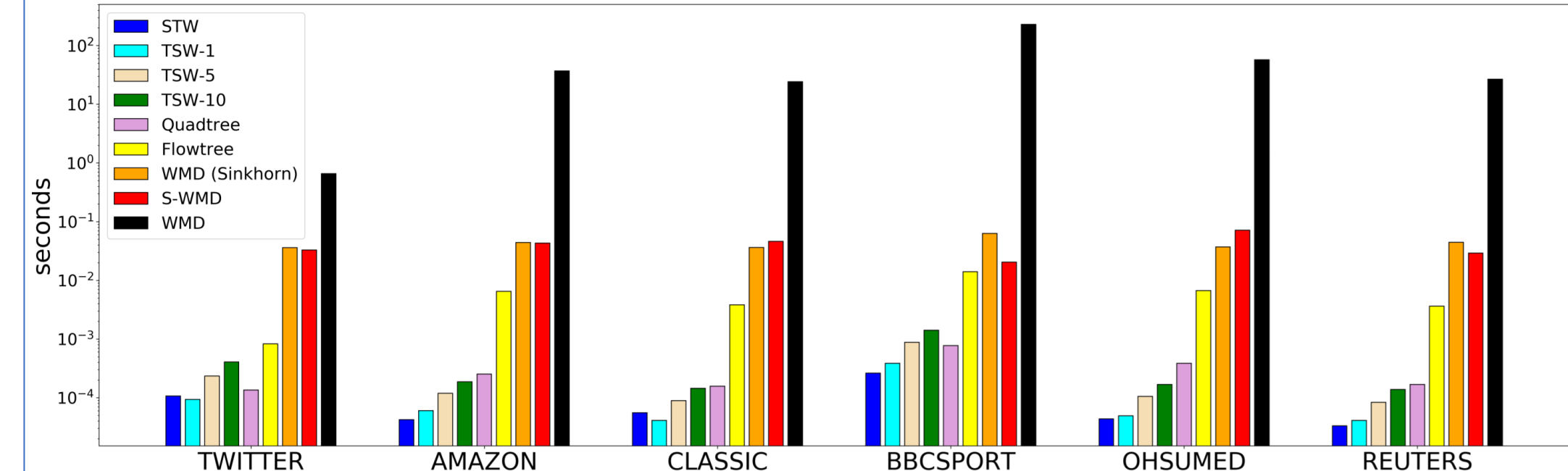


Fig.5 : Average time consumption for comparing 500 documents with one document.

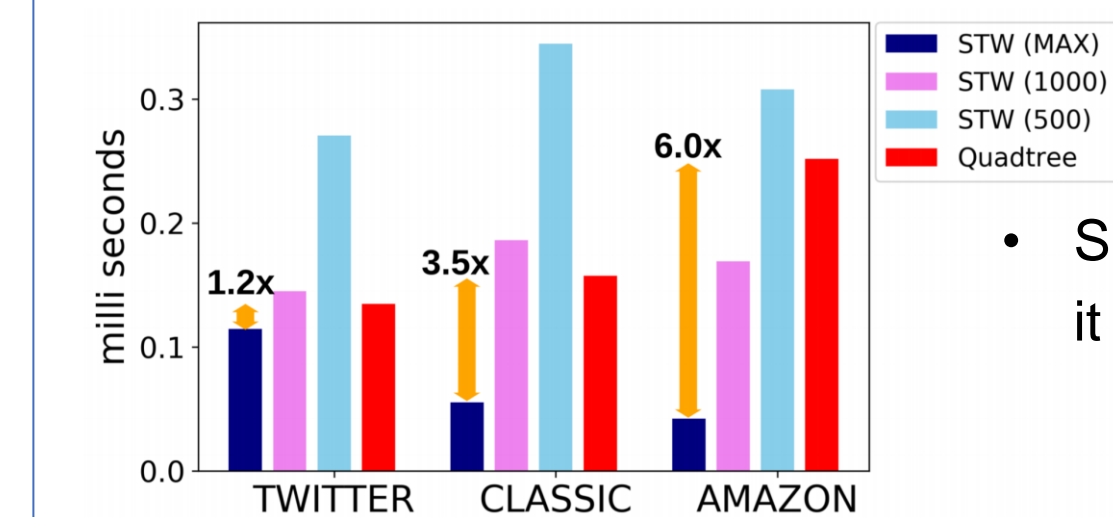


Fig. 6 : Running time varying the batch size.

- Since STW is suitable for batch processing, it can be computed faster.

Our code is available :

<https://github.com/yukiTakezawa/STW>