

# Supervised Tree-Wasserstein Distance

*Yuki Takezawa<sup>1 2</sup>, Ryoma Sato<sup>1 2</sup>, Makoto Yamada<sup>1 2</sup>*

*<sup>1</sup>Kyoto University, <sup>2</sup>RIKEN AIP*

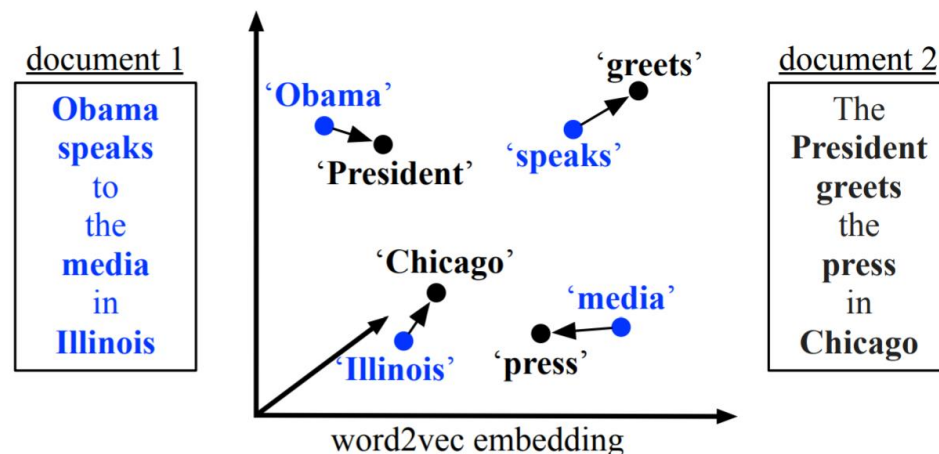
# Motivation

## The computational cost of the Wasserstein distance.

### Wasserstein distance

$$W_d(\mu_i, \mu_j) = \inf_{\gamma \in \Pi(\mu_i, \mu_j)} \int_{\Omega \times \Omega} d(x, y) \gamma(dx, dy)$$

- Linear Programming : cubic time.
- Sinkhorn Algorithm : quadratic time.



Word Mover's Distance [Kusner 15]

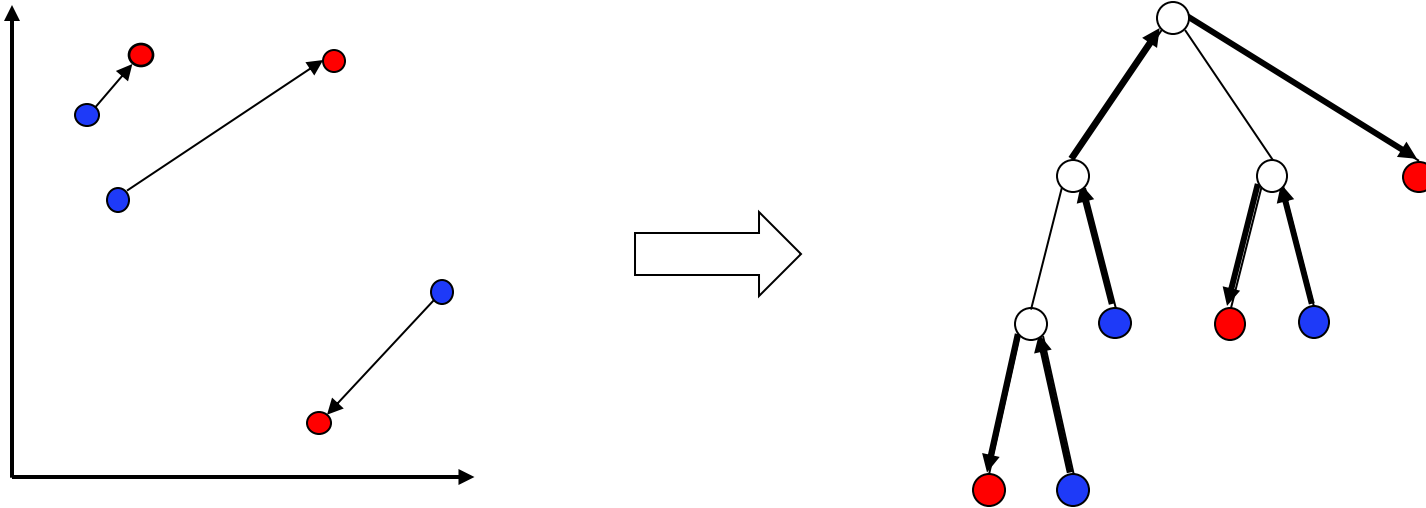
# Motivation

## Tree-Wasserstein Distance

---

The Wasserstein distance on a tree metric.

- The closed-form solution, which can be computed in linear time.



Quadtree [Indyk 03], clustering based method [Le 19]

- Unsupervised methods

# Motivation

## Supervised Tree-Wasserstein Distance

---

The distance between documents depends on the task.

- Topic, Author etc.
- Supervised WMD [Huang 16]
  - High computational cost.

We propose the **Supervised Tree-Wasserstein (STW) Distance**.

- Task-specific distance.
- Fast computation for the Wasserstein distance.
  - Tree-Wasserstein distance
  - GPU suitable. (e.g., batch processing)

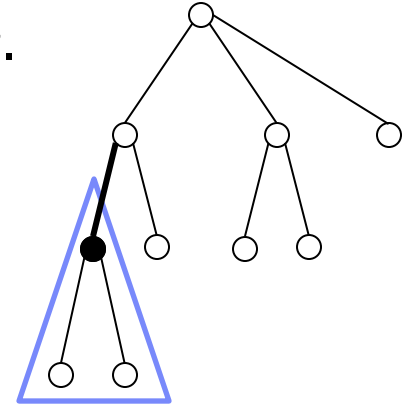
# Proposed Method

## Soft Tree-Wasserstein Distance

### Tree-Wasserstein Distance

- $\Gamma(v)$  : the set of nodes contained in the subtree rooted at  $v$ .

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



### Soft Tree-Wasserstein Distance

- $P_{\text{sub}}(x|v)$  : the probability that node  $x$  is contained in the subtree rooted at  $v$ .

$$\rightarrow \sum_{v \in V} w_v \left| \sum_{x \in V} P_{\text{sub}}(x|v) (\mu_i(x) - \mu_j(x)) \right|$$

# Proposed Method

## How to formulate $P_{\text{sub}}(x|v)$ ?

### Theorem 1 : conditions of an adjacency matrix

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

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

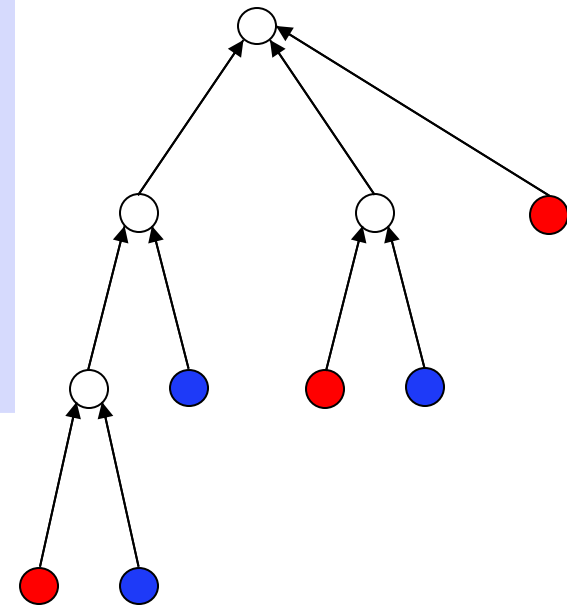
Then  $G$  is a directed tree.

Probability that  $v_i$  is a parent of  $v_j$ .

- $[\mathbf{D}_{\text{par}}]_{i,j} \in [0, 1]$

Probability that  $v_i$  is contained in the subtree rooted at  $v_j$ .

- $[\sum_{k=1}^{\infty} \mathbf{D}_{\text{par}}^k]_{i,j} = [(\mathbf{I} - \mathbf{D}_{\text{par}})^{-1}]_{i,j}$



# Proposed Method

## Supervised Tree-Wasserstein (STW) Distance

---

### Soft Tree-Wasserstein Distance

- Differentiable w.r.t. parent-child relationships. (i.e.,  $\mathbf{D}_{\text{par}}$ )

$$\begin{aligned} W_{d_T}^{\text{soft}}(\mu_i, \mu_j) &= \sum_{v \in V} w_v \left| \sum_{x \in V} P_{\text{sub}}(x|v) (\mu_i(v) - \mu_j(v)) \right|_{\alpha} \\ &= \left\| \mathbf{w}_v \circ (\mathbf{I} - \mathbf{D}_{\text{par}})^{-1} (\mathbf{a}_i - \mathbf{a}_j) \right\|_{\alpha} \end{aligned}$$

### Objective function:

$$\underset{\mathbf{D}_{\text{par}}}{\text{minimize}} \quad \sum_{i,j:\text{same label}} W_{d_T}^{\text{soft}}(\mu_i, \mu_j) - \sum_{i,j:\text{different label}} \min\{W_{d_T}^{\text{soft}}(\mu_i, \mu_j), m\}$$

s. t.  $\mathbf{D}_{\text{par}}$  satisfies the conditions of Th. 1.

# Experiment

## Document Classification Task

Table 1: The  $k$ NN test error for real datasets. The most accurate method based on the tree-Wasserstein distance is shown in blue.

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

Tree-Wasserstein distance



# Experiment

## Time Consumption

Tree-Wasserstein distance

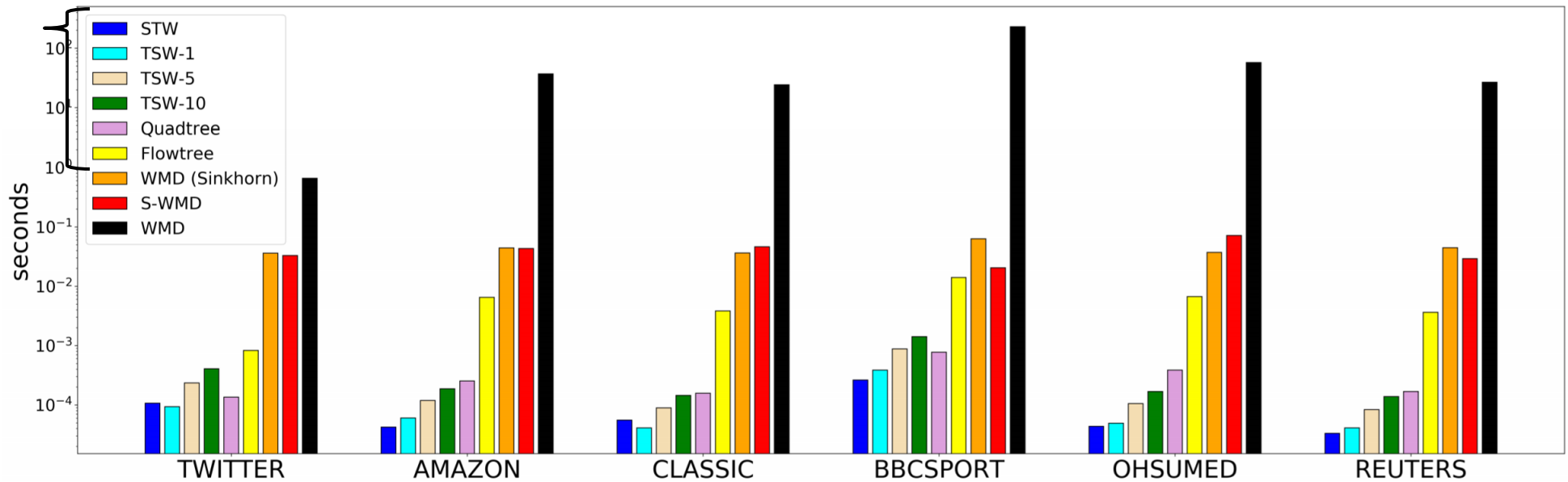


Figure 3. Average time consumption for comparing 500 documents with one document. For the STW distance and the TSW distance, the batch size is set to the number of documents contained in the training dataset. For WMD (Sinkhorn) and S-WMD, the batch size is set to 500 due to the memory size limitations. To obtain the average time consumption, we sample 100 documents as queries and measure the time consumption.