Proposal 0000 Experiment 00000 References 00000

# Wasserstein k-means with sparse simplex projection

#### Takumi Fukunaga<sup>1</sup> Hiroyuki Kasai <sup>2,3</sup>

<sup>1</sup>Dept. of Computer Science and Engineering,School of Fundamental Science and Engineering,Waseda University,Japan

<sup>2</sup>Dept. of Communications and Computer Engineering, School of Fundamental Science and Engineering, Waseda University, Japan

<sup>3</sup>Dept. of Computer Science and Communications Engineering, Graduate School of Fundamental Science and Engineering, Waseda University, Japan

#### December 9,2020

Introduction	
<b>●</b> 00	

Experiment 00000 References 00000

# Introduction

Wasserstein k-means

- Propose
  - Improvement of computational cost of Wasserstein k-means
- Fast and efficient approaches
  - [Cuturi and Doucet, 2014, Bonneel et al., 2015, Anderes et al., 2016]
- Contribution of our methods
  - Sparsifying data on probability simplex and shrinking them by removing the zero elements



Introduction	
000	

Experiment 00000 References 00000

# Introduction

Clustering Algorithm and its issue

- ▶ k-means [Lloyd, 1982]
  - High computational cost per iteration  $\mathcal{O}(qk)$ 
    - Efficient approaches [Arthur and Vassilvitskii, 2007, Kanungo et al., 2002]
- ▶ Wasserstein *k*-means [Ye et al., 2017]
  - Adopting Wasserstein distance and Wasserstein barycenter
    - High computational cost to caluculate its distance *O*(n<sup>3</sup> log n)[Cuturi, 2013, Rubner et al., 2000]

assignment step

$$s_i = \operatorname*{arg min}_{j=1,\ldots,k} d(\boldsymbol{x}_i, \boldsymbol{c}_j), \forall i \in [q]$$

### update step

$$c_j = \text{mean}(\{x | s_i = j\}) \text{ or } \text{barycenter}(\{x | s_i = j\}), \forall j \in [k]$$

Introduction	
000	

Experiment 00000 References 00000

### Introduction Optimal Transport

- Caluculate the minimum transport cost [Peyre and Cuturi, 2019]
- ▶ When n = m, its cost is called Wasserstein distance of order p
- Using this distance, calculate Wasserstein barycenter [Benamou et al., 2015]

$$\begin{array}{lll} \mathbf{T}^* &=& \displaystyle \operatorname*{arg\ min}_{\mathbf{T}\in\mathcal{U}_{mn}} \langle \mathbf{T},\mathbf{C}\rangle \\ W_p(\boldsymbol{\mu},\boldsymbol{\nu}) &=& \displaystyle \operatorname*{min}_{\mathbf{T}\in\mathcal{U}_{mn}} \langle \mathbf{T},\mathbf{C}\rangle = \langle \mathbf{T}^*,\mathbf{C}\rangle \\ g(\boldsymbol{\mu}) &=& \displaystyle \frac{1}{n}\sum_i W_p(\boldsymbol{\mu},\boldsymbol{\nu}_i) \end{array}$$

- $\blacktriangleright$  u and  $\mu$  of points
- a and b are in probability simplex
- C is ground matrix
- Row and Colunmn marginal constraint

$$\mathcal{U}_{mn} = \{\mathbf{T} \in \mathbb{R}^{m imes n}_+ : \mathbf{T} \mathbf{1}_n = \boldsymbol{a}, \mathbf{T}^T \mathbf{1}_m = \boldsymbol{b}\}$$

Proposal ●000 Experiment 00000 References 00000

# Proposal

Motivation

- Reduce the size of data and centroid
- Adopt two following approaches
- 1. Sparsify datas
  - Make data sparser than the original ones
  - Maintain degradation of the clustering quality as small as possible
- 2. Shrink datas
  - No degradation
  - Key operator to reduce the computational complexities

Introduction	
000	

Proposal O●OO Experiment 00000 References 00000

# Proposal

Basic idea - Sparse simplex projection

Sparse simplex projection GSHP [Kyrillidis et al., 2013]

$$\hat{\boldsymbol{\beta}} = \operatorname{Proj}^{\gamma(t)}(\boldsymbol{\beta}) = \left\{ \begin{array}{ll} \hat{\boldsymbol{\beta}}_{|\mathcal{S}^{\star}} &= \mathcal{P}_{\Delta_{\kappa}}(\boldsymbol{\beta}_{|\mathcal{S}^{\star}}) \\ \hat{\boldsymbol{\beta}}_{|(\mathcal{S}^{\star})^c} &= 0, \end{array} \right.$$

S is the subset of N = {1,...,n}
 a<sub>|S</sub> extracts the elements of S in a
 (P<sub>Δ<sub>κ</sub></sub>(β<sub>|S\*</sub>))<sub>v</sub> = [(β<sub>|S\*</sub>)<sub>v</sub> + τ]<sub>+</sub>, τ := <sup>1</sup>/<sub>κ</sub>(1 + Σ<sup>|S\*|</sup>β<sub>|S\*</sub>)
 S\* = supp(P<sub>Δ<sub>κ</sub></sub>)
 supp(a) = {i : a<sub>i</sub> ≠ 0}

Introduction	
000	

Proposal 00€0 Experiment 00000 References 00000

# Proposal

Basic idea - Shrinking datas

The zero elements don't have effect on transport matrix

Removing the zero elements

Define shrinking operator to vector and matrix

$$ilde{oldsymbol{
u}}_i = \operatorname{shrink}(\hat{oldsymbol{
u}}_i) = (\hat{oldsymbol{
u}}_i)_{|\mathcal{S}_{\mathsf{samp}}|} \in \mathbb{R}^{|\mathcal{S}_{\mathsf{samp}}|}$$

$$ilde{m{c}}_i \hspace{0.1 cm} = \hspace{0.1 cm} { ext{shrink}}(\hat{m{c}}_i) = (\hat{m{c}}_i)_{|\mathcal{S}_{ ext{cent}}|} \in \mathbb{R}^{|\mathcal{S}_{ ext{cent}}|}$$

$$\tilde{\mathbf{C}} = \mathsf{Shrink}(\mathbf{C}_{\boldsymbol{\nu}\boldsymbol{\mathcal{C}}}) = \mathbf{C}_{\mathsf{supp}(\hat{\boldsymbol{\nu}}_i),\mathsf{supp}(\hat{\boldsymbol{c}}_i)} \in \mathbb{R}^{|\mathcal{S}_{\mathsf{samp}}| \times |\mathcal{S}_{\mathsf{cent}}|}$$



Proposal 000● Experiment 00000 References 00000

# Proposal

Basic procedure

- 1. Update sparsity ratio  $\gamma(t)$
- 2. Project  $\boldsymbol{\nu}_i$  into  $\hat{\boldsymbol{\nu}_i}$  and shrink  $\hat{\boldsymbol{\nu}_i}$  into  $\tilde{\boldsymbol{\nu}_i}$
- 3. Project  $c_j$  into  $\hat{c_j}$  and shrink  $\hat{c_j}$  into  $ilde{c_i}$
- 4. Shrink ground cost matrix  ${\bf C}$  into  $\tilde{{\bf C}}$
- 5. Find closest centroids and update centroids
- 6. Unless cluster centroids stop changing, repeat step1
- Control parameter of sparse ratio  $\gamma(t)$

$$\gamma(t) := \begin{cases} \gamma_{\min} & (FIX) \\ 1 - \frac{(1 - \gamma_{\min})}{T_{\max}} t & (DEC) \\ \gamma_{\min} + \frac{(1 - \gamma_{\min})}{T_{\max}} t & (INC), \end{cases}$$

•  $T_{\max}$  is maximum iteration of k-means

Introduction
000

Experiment •0000 References 00000

# Experiment

Settings

- Use of Algorithm
  - Wasserstein barycenter [Cuturi and Doucet, 2014]
  - k-means with litekmeans
  - linprog of Mosek to solve LP [Andersen et al., 2000]
- Datasets
  - COIL-100 [Nene et al., 1996]
  - the USPS handwritten dataset

Proposal 0000 Experiment 00000 References 00000

### Experiment

#### 2-D histogram evaluation



Figure: Performance results of 2-D histogram data on the USPS dataset. KASAI Laboratory, WASEDA University. All Rights Reserved.

Proposal 0000 Experiment 00000 References 00000

### Experiment

#### Convergence Performance



Figure: Left : Convergence performance with different projection data using DEC algorithm of  $\gamma_{\min} = 0.5$  Right:Convergence performance comparison of different algorithm of  $\gamma(t)$  of  $\gamma_{\min} = 0.5$ 

Introduction	Proposal	Experiment
000	0000	00000

# Experiment

#### Comparoson on different sparsity



Figure: Performance comparison on different ratios on the USPS dataset.

References

Proposal 0000 Experiment 00000 References 00000

# Experiment

Conclusion

- We propose a faster Wasserstein k-means algorithm
- Our experiments demonstrate the effectiveness of our method
  - Reducing the computational complexity of Wasserstein distance
  - Keeping accuracy before sparsifying and shrinking

Introduction	
000	

Experiment 00000 References

# References I

- Anderes, E., Borgwardt, S., and Miller, J. (2016).
   Discrete wasserstein barycenters: optimal transport for discrete data.
   Mathematical Methods of Operations Research, 84:389–409.
- Andersen, E. D., Roos, C., and Terlaky, T. (2000). <u>High Performance Optimization</u>, volume 33, chapter The MOSEK interior point optimizer for linear programming: an implementation of the homogeneous algorithm. <u>Springer, Boston, MA.</u>
- Arthur, D. and Vassilvitskii, S. (2007).
   K-means++: The advantages of careful seeding.
   In Proc. 18th Ann. ACM-SIAM Symp. Discr. Alg. (SODA).

 Benamou, J. D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. (2015).
 Iterative bregman projections for regularized transportation problems. SIAM Journal on Scientific Computing, 37(2):1111–A1138.

Introduction
000

Experiment 00000 References

# References II

- Bonneel, N., Rabin, J., and Peyré, G. (2015).
   Sliced and radon wasserstein barycenters of measures.
   Journal of Mathematical Imaging and Vision, 51:22–45.
- Cuturi, M. (2013).
   Sinkhorn distances: Lightspeed computation of optimal transport. In <u>Annual Conference on Neural Information Processing Systems</u> (NIPS).
- Cuturi, M. and Doucet, A. (2014).
   Fast computation of wasserstein barycenters. In <u>ICML</u>.
- Kanungo, T., Mount, D. M., Netanyahu, N. S., Piatko, C. D., Silverman, R., and Wu, A. Y. (2002).
   An efficient k-means clustering algorithm: analysis and implementation.
   IEEE Trans. Pattern Anal. Mach. Intell., 24(7):881–892.

Introduction	
000	

Experiment 00000 References

# References III

- Kyrillidis, A., Becker, S., Cevher, V., and Koch, C. (2013).
   Sparse projections onto the simplex.
   In <u>ICML</u>.
- ► Lloyd, S. (1982).

Least squaares quantization in pcm. IEEE Trans. Inf. Theory, 28(2):129–137.

- Nene, S. A., Nayar, S. K., and Murase, H. (1996). Columbia university image library (COIL-20). Technical Report CUCS-005-96.
- Peyre, G. and Cuturi, M. (2019).
   Computational optimal transport.
   Foundations and Trends in Machine Learning, 11(5-6):355–607.
- Rubner, Y., Tomasi, C., and Guibas, L. J. (2000).
   The earth mover's distance as a metric for image retrieval.
   International Journal of Computer Vision, 40(2):99–121.

Introduction	
000	

Experiment 00000 References

### **References IV**

▶ Ye, Y., Wu, P., Wang, J. Z., and Li, J. (2017).

Fast discrete distribution clustering using wasserstein barycenter with sparse support.

IEEE Trans. Signal Process, 65(9):2317–2332.

Proposal 0000 Experiment 00000 References

# Thank you for listening.



# WASEDA University