

The background of the slide features a repeating pattern of stylized, light pink flowers and leaves. The flowers have five petals and are arranged in a circular, wreath-like fashion around the central text.

Fast Discrete Distribution Clustering Using Wasserstein Barycenter with Sparse Support

IEEE TRANSACTIONS ON SIGNAL PROCESSING(2017)

Jianbo Ye; Panruo Wu; James Z. Wang; Jia Li
Reporter: Fengjiao Gong

April 14, 2022

Outline

1. Discrete Distribution Clustering(D2-clustering)
 - 1.1 Wasserstein Distance & Barycenter
 - 1.2 Clustering Criterion
 - 1.3 Algorithm
2. Alternating Direction Method of Multipliers(ADMM)
 - 2.1 Precursors
 - 2.2 Algorithm
3. Bregman ADMM(B-ADMM)
 - 3.1 Introduction
 - 3.2 Algorithm
4. Experiments

Discrete Distribution Clustering(D2-clustering)

Wasserstein Distance

Consider two distributions

$$P^{(a)} = \left\{ \left(w_i^{(a)}, x_i^{(a)} \right), i = 1, \dots, m_a \right\}$$

$$P^{(b)} = \left\{ \left(w_i^{(b)}, x_i^{(b)} \right), i = 1, \dots, m_b \right\}$$

Then

$$\begin{aligned} \left(W_p \left(P^{(a)}, P^{(b)} \right) \right)^p &:= \min_{\{\pi_{ij} \geq 0\}} \sum_{i \in \mathcal{I}_a, j \in \mathcal{I}_b} \pi_{ij} c \left(x_i^{(a)}, x_j^{(b)} \right) \\ \text{s.t. } \sum_{i=1}^{m_a} \pi_{ij} &= w_j^{(b)}, \forall j \in \mathcal{I}_b, \\ \sum_{j=1}^{m_b} \pi_{ij} &= w_i^{(a)}, \forall i \in \mathcal{I}_a \end{aligned} \tag{1}$$

Discrete Distribution Clustering(D2-clustering)

Wasserstein Distance

where

$$c\left(x_i^{(a)}, x_j^{(b)}\right)=\left\|x_i^{(a)}-x_j^{(b)}\right\|_p^p$$

$$\mathcal{I}_a=\left\{1, \ldots, m_a\right\} \quad \text { (index set) }$$

$$\mathcal{I}_b=\left\{1, \ldots, m_b\right\} \quad \text { (index set) }$$

$$\left\{\pi_{i, j}\right\} \quad \begin{array}{l} \text { matching weights} \\ \text { optimal coupling} \end{array}$$

Here

$$W=W_2=W_p|_{p=2}$$

Discrete Distribution Clustering(D2-clustering)

Wasserstein Barycenter

Suppose $\{P^{(1)}, \dots, P^{(N)}\}$ is a set of discrete distributions

$$\min_P \frac{1}{N} \sum_{k=1}^N W^2(P, P^{(k)}) \quad (2)$$

where

$$P : \{(w_1, x_1), \dots, (w_m, x_m)\}$$

$$\sum_{i=1}^m w_i = 1, w_i \geq 0$$

So P is centroid

Discrete Distribution Clustering(D2-clustering)

Clustering Criterion

Criterion:

- ▶ minimize total within-cluster variation under the Wasserstein distance
- ▶ find a set of centroid distributions $\{Q^{(i)}, i = 1, \dots, K\}$

$$\min_{\{Q^{(i)}\}} \sum_{k=1}^{\bar{N}} \min_{i=1, \dots, K} W^2(Q^{(i)}, P^{(k)})$$

Discrete Distribution Clustering(D2-clustering) Algorithm

Alternating Optimization

- ▶ **Assignment:** assign each instance to the nearest centroid
- ▶ **Update Centroids:** optimize centroids

Discrete Distribution Clustering(D2-clustering) Algorithm

Algorithm 1 D2 Clustering

```
1: procedure D2CLUSTERING( $\{P^{(k)}\}_{k=1}^M, K$ )
2:   Denote the label of each objects by  $l^{(k)}$ .
3:   Initialize  $K$  random centroid  $\{Q^{(i)}\}_{i=1}^K$ .
4:   repeat
5:     for  $k = 1, \dots, M$  do ▷ Assignment Step
6:        $l^{(k)} := \operatorname{argmin}_i W(Q^{(i)}, P^{(k)})$ ;
7:     for  $i = 1, \dots, K$  do ▷ Update Step
8:        $Q^{(i)} := \operatorname{argmin}_Q \sum_{l^{(k)}=i} W(Q, P^{(k)})$  (*)
9:   until the number of changes of  $\{l^{(k)}\}$  meets some
      stopping criterion
10:  return  $\{l^{(k)}\}_{k=1}^M$  and  $\{Q^{(i)}\}_{i=1}^K$ .
```

Figure 1: Outer Loop

Computational challenge \rightarrow optimal centroid *for each cluster at each iteration*

Discrete Distribution Clustering(D2-clustering)

Algorithm - Main question

$$\min_P \frac{1}{N} \sum_{k=1}^N W^2 \left(P, P^{(k)} \right) \quad (2)$$

Variables:

- ▶ weights in centroid $\{w_i \in \mathbb{R}^+\}$
- ▶ support points $\{x_i \in \mathbb{R}^d\}$
- ▶ optimal coupling between P and $P^{(k)}$ $\left\{ \pi_{i,j}^{(k)} \right\}$

Discrete Distribution Clustering(D2-clustering)

Algorithm - Main question

$$\min_P \frac{1}{N} \sum_{k=1}^N W^2(P, P^{(k)}) \quad (2)$$

Alternating Optimization:

- Fix \mathbf{w} and Π , cost function (2) is quadratic in terms of \mathbf{x} :

$$x_i := \frac{1}{Nw_i} \sum_{k=1}^N \sum_{j=1}^{m_k} \pi_{i,j}^{(k)} x_j^{(k)}, \quad i \in \mathcal{I}' = \{1, \dots, m\} \quad (3)$$

Matrix form

$$\mathbf{x} := \frac{1}{N} X \Pi^T \text{diag}(1./\mathbf{w})$$

Discrete Distribution Clustering(D2-clustering)

Algorithm - Main question

$$\min_P \frac{1}{N} \sum_{k=1}^N W^2(P, P^{(k)}) \quad (2)$$

Alternating Optimization:

- Fix \mathbf{x} , updating \mathbf{w} and Π is challenging (large LP):

$$\begin{aligned} \min_{\Pi \in \mathbb{R}_{m \times n}^+, \mathbf{w} \in \Delta_m} \sum_{k=1}^N \left\langle C(\mathbf{x}, \mathbf{x}^{(k)}), \Pi^{(k)} \right\rangle \\ \text{s.t.} \quad \mathbf{1} \cdot \left(\Pi^{(k)} \right)^T = \mathbf{w} \\ \mathbf{1} \cdot \Pi^{(k)} = \mathbf{w}^{(k)} \\ \forall k \in \mathcal{I}^c = \{1, \dots, N\} \end{aligned} \quad (4)$$

Scalability

Discrete Distribution Clustering(D2-clustering) Algorithm

Algorithm 2 Centroid Update with Full-batch LP

```
1: procedure CENTROID( $\{P^{(k)}\}_{k=1}^N$ )  
2:   repeat  
3:     Updates  $\{x_i\}$  from Eq. (3);  
4:     Updates  $\{w_i\}$  from solving full-batch LP (4);  
5:   until  $P$  converges  
6:   return  $P$ 
```

Figure 2: Inner Loop

Eq(4) is not end result but just one round of centroid update in outer loop

Scalable Methods: not accurate, but fast approximation

Alternating Direction Method of Multipliers(ADMM)

Precursors

- ▶ Dual Ascent
- ▶ Dual Decomposition
- ▶ Augmented Lagrangian & Method of Multipliers

Alternating Direction Method of Multipliers

- ▶ Algorithm
- ▶ Scaled Form

ADMM to solve Eq(4)

Alternating Direction Method of Multipliers(ADMM)

Precursors - Dual Ascent

Convex optimization problem:

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & Ax = b \end{array} \quad (5)$$

Lagrangian function:

$$L(x, y) = f(x) + y^T(Ax - b)$$

Strong duality:

same optimal values

Dual problem:

$$\text{maximize} \quad g(y) = \inf_x L(x, y)$$

Alternating Direction Method of Multipliers(ADMM)

Precursors - Dual Ascent

Lagrangian function

$$L(x, y) = f(x) + y^T(Ax - b)$$

x^* can be recovered by y^*

$$x^* = \arg \min_x L(x, y^*)$$

[reference] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers”, Foundations and Trends R in Machine Learning, vol. 3, no. 1, pp. 1–122, 2011.

Alternating Direction Method of Multipliers(ADMM)

Precursors - Dual Ascent

$$x^* = \arg \min_x L(x, y^*)$$

Dual Ascent Method:

$$\begin{aligned} x^{k+1} &:= \arg \min_x L(x, y^k) \\ y^{k+1} &:= y^k + \alpha^k (Ax^{k+1} - b) \\ \alpha^k &> 0 \end{aligned} \tag{6}$$

Alternating Direction Method of Multipliers(ADMM)

Precursors - Dual Decomposition

Suppose objective f is separable

$$\begin{aligned}f(\mathbf{x}) &= \sum_{i=1}^N f_i(x_i) \\ \mathbf{x} &= (x_1, \dots, x_N) \\ A &= [A_1, \dots, A_N]\end{aligned}$$

Then

$$L(x, y) = \sum_{i=1}^N L_i(x_i, y) = \sum_{i=1}^N f_i(x_i) + y^T(A_i x_i - \frac{1}{N}b)$$

Eq(6) can be solved in N separate problems in parallel

Alternating Direction Method of Multipliers(ADMM)

Precursors - Augmented Lagrangian & Method of Multipliers

Augmented Lagrangian

$$L_{\rho}(x, y) = f(x) + y^T(Ax - b) + \frac{\rho}{2}\|Ax - b\|_2^2, \quad \rho > 0 \text{ (penalty)}$$

Apply Dual Ascent:

$$\begin{aligned} x^{k+1} &:= \arg \min_x L_{\rho}(x, y^k) \\ y^{k+1} &:= y^k + \rho(Ax^{k+1} - b) \\ \rho &> 0 \end{aligned} \tag{7}$$

Method of Multipliers

Alternating Direction Method of Multipliers(ADMM) Algorithm

Problem form:

$$\begin{array}{ll}\text{minimize} & f(x) + g(z) \\ \text{subject to} & Ax + Bz = c\end{array}$$

suppose f and g are convex

$$\mathbf{x}_{\text{before}} \Rightarrow (x, z)$$

Alternating Direction Method of Multipliers(ADMM) Algorithm

Augmented Lagrangian

$$L_{\rho}(x, y, z) = f(x) + g(z) + y^T(Ax + Bz - c) + \frac{\rho}{2}\|Ax + Bz - c\|_2^2$$

Iteration:

$$\begin{aligned}x^{k+1} &:= \arg \min_x L_{\rho}(x, z^k, y^k) \\z^{k+1} &:= \arg \min_z L_{\rho}(x^{k+1}, z, y_k) \\y^{k+1} &:= y^k + \rho(Ax^{k+1} + Bz^{k+1} - c)\end{aligned}$$

Alternating Direction Method of Multipliers(ADMM)

Algorithm - Scaled Form

Denote $u = \frac{1}{\rho}y$

$$\begin{aligned}x^{k+1} &:= \arg \min_x (f(x) + \frac{\rho}{2} \|Ax + Bz^k - c + u^k\|_2^2) \\z^{k+1} &:= \arg \min_z (g(z) + \frac{\rho}{2} \|Ax^{k+1} + Bz - c + u^k\|_2^2) \\u^{k+1} &:= u^k + Ax^{k+1} + Bz^{k+1} - c\end{aligned}$$

Formulas shorter

ADMM to solve Eq(4)

Discrete Distribution Clustering(D2-clustering)

Algorithm - Main question

$$\min_P \frac{1}{N} \sum_{k=1}^N W^2(P, P^{(k)}) \quad (2)$$

Alternating Optimization:

- Fix \mathbf{x} , updating \mathbf{w} and Π is challenging (large LP):

$$\begin{aligned} \min_{\Pi \in \mathbb{R}_{m \times n}^+, \mathbf{w} \in \Delta_m} \sum_{k=1}^N \left\langle C(\mathbf{x}, \mathbf{x}^{(k)}), \Pi^{(k)} \right\rangle \\ \text{s.t.} \quad \mathbf{1} \cdot \left(\Pi^{(k)} \right)^T = \mathbf{w} \\ \mathbf{1} \cdot \Pi^{(k)} = \mathbf{w}^{(k)} \\ \forall k \in \mathcal{I}^c = \{1, \dots, N\} \end{aligned} \quad (4)$$

Scalability

Alternating Direction Method of Multipliers(ADMM)

► Scaled Augmented Lagrangian

$$L_{\rho}(\Pi, \mathbf{w}, \Lambda) = \sum_{k=1}^N \left\langle \mathbf{C} \left(\mathbf{x}, \mathbf{x}^{(k)} \right), \Pi^{(k)} \right\rangle + \rho \sum_{i \in \mathcal{I}' k \in \mathcal{I}^c} \lambda_{i,k} \left(\sum_{j=1}^{m_k} \pi_{i,j}^{(k)} - w_i \right) \\ + \frac{\rho}{2} \sum_{i \in \mathcal{I}' k \in \mathcal{I}^c} \left(\sum_{j=1}^{m_k} \pi_{i,j}^{(k)} - w_i \right)^2$$

► Iteration

$$\begin{aligned} \Pi^{n+1} &:= \underset{\Pi \in \Delta_{\Pi}}{\operatorname{argmin}} L_{\rho}(\Pi, \mathbf{w}^n, \Lambda^n), \\ \mathbf{w}^{n+1} &:= \underset{\mathbf{w} \in \Delta_m}{\operatorname{argmin}} L_{\rho}(\Pi^{n+1}, \mathbf{w}, \Lambda^n), \\ \lambda_{i,k}^{n+1} &:= \lambda_{i,k}^n + \sum_{j=1}^{m_k} \pi_{i,j}^{(k),n+1} - w_i^{n+1}, i \in \mathcal{I}', k \in \mathcal{I}^c. \end{aligned} \tag{8}$$

Alternating Direction Method of Multipliers(ADMM)

- Update $\pi_{i,j}^k$: equivalent form:

$$\begin{aligned} \min_{\pi_{i,j}^{(k)} \geq 0} & \left\langle C(\mathbf{x}, \mathbf{x}^{(k)}), \Pi^{(k)} \right\rangle + \frac{\rho}{2} \sum_{i=1}^m \left(\sum_{j=1}^{m_k} \pi_{i,j}^{(k)} - w_i^n + \lambda_{i,k}^n \right)^2 \\ \text{s.t. } & \mathbf{1} \cdot \Pi^{(k)} = \mathbf{w}^{(k)}, k \in \mathcal{I}^c. \end{aligned} \quad (9)$$

- Update \mathbf{w}_i : denote $\tilde{w}_i^{(k),n+1} = \sum_{j=1}^{m_k} \pi_{i,j}^{(k),n+1} + \lambda_{i,k}^n, i = 1, \dots, m$

$$\min_{\mathbf{w} \in \Delta_m} \sum_{i=1}^m \sum_{k=1}^N (\tilde{w}_i^{(k),n+1} - w_i)^2 \quad (10)$$

Alternating Direction Method of Multipliers(ADMM) Algorithm

Algorithm 3 Centroid Update with ADMM

```
1: procedure CENTROID( $\{P^{(k)}\}_{k=1}^N, P, \Pi$ )
2:   Initialize  $\Lambda^0 = 0$  and  $\Pi^0 := \Pi$ .
3:   repeat
4:     Updates  $\{x_i\}$  from Eq.(3);
5:     Reset dual coordinates  $\Lambda$  to zero;
6:     for  $iter = 1, \dots, T_{admm}$  do
7:       for  $k = 1, \dots, N$  do
8:         Update  $\{\pi_{i,j}\}^{(k)}$  based on QP
9:       Update  $\{w_i\}$  based on QP
10:      Update  $\Lambda$  based on Eq. (8);
11:   until  $P$  converges
12:   return  $P$ 
```

Figure 3: Centroid Update with ADMM

Computation limitation:

solve N of LP \Rightarrow solve N of QP

Bregman ADMM(B-ADMM)

Introduction

B-ADMM replace the quadratic penalty term by **Bregman divergence**:

$$L_{\rho}^{\phi}(\mathbf{x}, \mathbf{z}, \mathbf{y}) = f(\mathbf{x}) + g(\mathbf{z}) + \langle \mathbf{y}, \mathbf{Ax} + \mathbf{Bz} - \mathbf{c} \rangle + \rho B_{\phi}(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz})$$

where

$$B_{\phi}(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0$$

as $\phi : \Omega \rightarrow \mathbb{R}$ is a continuously differentiable and **strictly convex** function on the relative interior of a convex set Ω

Here $B_{\phi}(\mathbf{x}, \mathbf{y})$ is not necessarily convex w.r.t $y \Rightarrow$ *Not directly from ADMM*

Bregman ADMM(B-ADMM)

Introduction

B-ADMM Update:

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \langle \mathbf{y}_t, \mathbf{Ax} + \mathbf{Bz}_t - \mathbf{c} \rangle + \rho B_\phi(\mathbf{c} - \mathbf{Ax}, \mathbf{Bz}_t)$$

$$\mathbf{z}_{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathcal{Z}} g(\mathbf{z}) + \langle \mathbf{y}_t, \mathbf{Ax}_{t+1} + \mathbf{Bz} - \mathbf{c} \rangle + \rho B_\phi(\mathbf{Bz}, \mathbf{c} - \mathbf{Ax}_{t+1})$$

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \rho(\mathbf{Ax}_{t+1} + \mathbf{Bz}_{t+1} - \mathbf{c})$$

Ensure new updates do not violate the equality constraint significantly

Bregman ADMM(B-ADMM)

Application

Optimal Transport Problem:

$$\min \langle \mathbf{C}, \mathbf{X} \rangle \text{ s.t. } \mathbf{X}\mathbf{e} = \mathbf{a}, \mathbf{X}^T\mathbf{e} = \mathbf{b}, \mathbf{X} \geq 0 \quad (\textit{Linear Program})$$

where $\langle \mathbf{C}, \mathbf{X} \rangle = \text{Tr}(\mathbf{C}^T\mathbf{X})$, $\mathbf{C} \in \mathbb{R}^{m \times n}$ is cost matrix, $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{a} \in \mathbb{R}^{m \times 1}$, $\mathbf{b} \in \mathbb{R}^{n \times 1}$, \mathbf{e} is a column vector of ones.

Bregman ADMM(B-ADMM)

Application

Introduce variable \mathbf{Z} to split the constraints into two simplex such that

$$\Delta_{\mathbf{x}} = \{\mathbf{X} \mid \mathbf{X} \geq 0, \mathbf{X}\mathbf{e} = \mathbf{a}\}$$

$$\Delta_{\mathbf{z}} = \{\mathbf{Z} \mid \mathbf{Z} \geq 0, \mathbf{Z}^T \mathbf{e} = \mathbf{b}\}$$

Then

$$\min \langle \mathbf{C}, \mathbf{X} \rangle \text{ s.t. } \mathbf{X} \in \Delta_{\mathbf{x}}, \mathbf{Z} \in \Delta_{\mathbf{z}}, \mathbf{X} = \mathbf{Z}$$

Bregman ADMM(B-ADMM)

Application

B-ADMM Updates:

$$\begin{aligned}\mathbf{X}^{t+1} &= \operatorname{argmin}_{\mathbf{X} \in \Delta_{\mathbf{x}}} \langle \mathbf{C}, \mathbf{X} \rangle + \langle \mathbf{Y}^t, \mathbf{X} \rangle + \rho \mathbf{KL}(\mathbf{X}, \mathbf{Z}^t), \\ \mathbf{Z}^{t+1} &= \operatorname{argmin}_{\mathbf{Z} \in \Delta_{\mathbf{z}}} \langle \mathbf{Y}^t, -\mathbf{Z} \rangle + \rho \mathbf{KL}(\mathbf{Z}, \mathbf{X}^{t+1}), \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t + \rho (\mathbf{X}^{t+1} - \mathbf{Z}^{t+1}).\end{aligned}$$

Closed Form:

$$X_{ij}^{t+1} = \frac{Z_{ij}^t \exp\left(-\frac{C_{ij} + Y_{ij}^t}{\rho}\right)}{\sum_{j=1}^n Z_{ij}^t \exp\left(-\frac{C_{ij} + Y_{ij}^t}{\rho}\right)} a_i, \quad Z_{ij}^{t+1} = \frac{X_{ij}^{t+1} \exp\left(\frac{Y_{ij}^t}{\rho}\right)}{\sum_{i=1}^m X_{ij}^{t+1} \exp\left(\frac{Y_{ij}^t}{\rho}\right)} b_j$$

BADMM can be faster than ADMM by a factor of $O(n/\ln n)$

Bregman ADMM(B-ADMM) Algorithm

Consider two sets of variables

$$\Delta_{k,1} := \left\{ \pi_{ij}^{(k,1)} \geq 0 : \sum_{i=1}^m \pi_{ij}^{(k,1)} = w_j^{(k)}, j \in \mathcal{I}_k \right\}$$
$$\Delta_{k,2}(\mathbf{w}) := \left\{ \pi_{ij}^{(k,2)} \geq 0 : \sum_{j=1}^{m_k} \pi_{ij}^{(k,2)} = w_i, i \in \mathcal{I}' \right\}$$

Then $\Pi^{(k,1)} \in \Delta_{k,1}$ and $\Pi^{(k,2)} \in \Delta_{k,2}(\mathbf{w})$

Bregman ADMM(B-ADMM)

Algorithm

Notation

$$\begin{aligned}\bar{\Pi}^{(1)} &= \{\Pi^{(1,1)}, \Pi^{(2,1)}, \dots, \Pi^{(N,1)}\} \\ \bar{\Pi}^{(2)} &= \{\Pi^{(1,2)}, \Pi^{(2,2)}, \dots, \Pi^{(N,2)}\} \\ \bar{\Pi} &= \{\bar{\Pi}^{(1)}, \bar{\Pi}^{(2)}\} \\ \Lambda &= \{\Lambda^{(1)}, \dots, \Lambda^{(N)}\} \\ \Lambda^{(k)} &= \left(\lambda_{i,j}^{(k)}\right), i \in \mathcal{I}', j \in \mathcal{I}_k\end{aligned}$$

Formulation

$$\begin{aligned}\min_{\bar{\Pi}, \mathbf{w}} \quad & \sum_{k=1}^N \left\langle C\left(\mathbf{x}, \mathbf{x}^{(k)}\right), \Pi^{(k,1)} \right\rangle \\ \text{s.t. } \quad & \mathbf{w} \in \Delta_m \\ & \Pi^{(k,1)} \in \Delta_{k,1}, \quad \Pi^{(k,2)} \in \Delta_{k,2}(\mathbf{w}) \\ & \Pi^{(k,1)} = \Pi^{(k,2)}, \quad k = 1, \dots, N\end{aligned}$$

Bregman ADMM(B-ADMM)

Algorithm

B-ADMM Updates:

$$\begin{aligned}\bar{\Pi}^{(1),n+1} &:= \operatorname{argmin}_{\{\Pi^{(k,1)} \in \Delta_{k,1}\}} \sum_{k=1}^N (\langle C(\mathbf{x}, \mathbf{x}^{(k)}) , \Pi^{(k,1)} \rangle + \langle \Lambda^{(k),n}, \Pi^{(k,1)} \rangle + \rho \text{KL}(\Pi^{(k,1)}, \Pi^{(k,2),n})) \\ \bar{\Pi}^{(2),n+1}, \mathbf{w}^{n+1} &:= \operatorname{argmin}_{\substack{\{\Pi^{(k,2)} \in \Delta_{k,1}(\mathbf{w})\} \\ \mathbf{w} \in \Delta_m}} \sum_{k=1}^N (-\langle \Lambda^{(k),n}, \Pi^{(k,2)} \rangle + \rho \text{KL}(\Pi^{(k,2)}, \Pi^{(k,1),n+1})) \\ \Lambda^{n+1} &:= \Lambda^n + \rho (\bar{\Pi}^{(1),n+1} - \bar{\Pi}^{(2),n+1})\end{aligned}$$

Kullback-Leibler divergence(KL divergence):

$$B_{\phi}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i \log \frac{x_i}{y_i}$$

Bregman ADMM(B-ADMM)

Algorithm

Algorithm 4 Centroid Update with B-ADMM

```
1: procedure CENTROID( $\{P^{(k)}\}_{k=1}^N, P, \Pi$ ).  
2:    $\Lambda := 0; \bar{\Pi}^{(2),0} := \Pi$ .  
3:   repeat  
4:     Update  $x$  from Eq.(3) per  $\tau$  loops;  
5:     for  $k = 1, \dots, N$  do  
6:       Update  $\Pi^{(k,1)}$   
7:       Update  $\{\tilde{\pi}_{i,j}^{(k,1)}\}$   
8:       Update  $w$   
9:       for  $k = 1, \dots, N$  do  
10:        Update  $\Pi^{(k,2)}$   
11:         $\Lambda^{(k)} := \Lambda^{(k)} + \rho(\Pi^{(k,1)} - \Pi^{(k,2)})$ ;  
12:   until  $P$  converges  
13:   return  $P$ 
```

Figure 4: Centroid Update with B-ADMM

where $\tilde{\pi}_{ij}^{(k,1),n+1} := \pi_{ij}^{(k,1),n+1} \exp \left[\frac{1}{\rho} \lambda_{ij}^{(k),n} \right] + \text{eps}$, eps is floating-point tolerance.

Experiments

Data	\bar{N}	d	m	K
synthetic	2,560,000	≥ 16	≥ 32	256
image color	5,000	3	8	10
image texture	-	-	-	-
USPS digits	11,000	2	80	360
BBC news abstract	2,225	300	16	15
Wiki events abstract	1,983	400	16	100
20newsgroups GV	18,774	300	64	40
20newsgroups WV	-	400	100	-

Figure 5: DATASETS IN THE EXPERIMENTS

\bar{N} : DATA SIZE

d : DIMENSION OF THE SUPPORT VECTORS

m : NUMBER OF SUPPORT VECTORS IN A CENTROID

K : MAXIMUM NUMBER OF CLUSTERS TESTED

ENTRY WITH SAME VALUE AS IN PREVIOUS ROW IS “-”

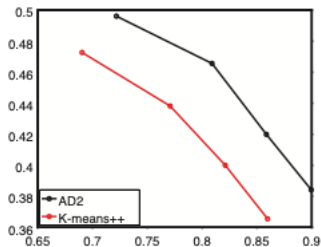
Experiments

# processors	32	64	128	256	512
SSE (%)	93.9	93.4	92.9	84.8	84.1
WSE on \bar{N} (%)	99	94.8	95.7	93.3	93.2
WSE on m (%)	96.6	89.4	83.5	79.0	-

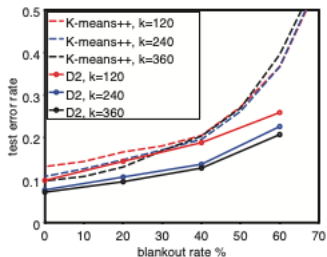
Figure 6: SCALING EFFICIENCY OF AD2-CLUSTERING IN PARALLEL IMPLEMENTATION

- ▶ Strong scaling efficiency(SSE):
speed-up gained from using more and more processors when the problem is fixed in size
- ▶ Weak scaling efficiency(WSE):
how stable the real computation time can be when proportionally more processors are used as the size of the problem grows

Experiments



(a) Homogeneity vs. completeness



(b) Test error rate vs. blankout rate

Figure 7: Comparisons between Kmeans++ and AD2-clustering on USPS dataset.

Thanks!