

# Extended Abstract: Hierarchical Data Analysis

Wendy Di, Stefan Wild

June 6, 2017

## 1 Notations

**Goal:** Given a halo with  $N_p$  particles, find its MBP.

- $D$ : dimension of the space,  $D = 2$  for now
- $\mathbf{X} = [X_i] \in \mathbb{R}^{N_p \times 2}$ : collection of particles
- $\mathbf{m} = [m_i] \in \mathbb{R}^{N_p}$ : collection of each particle's mass, normally  $m_i = 1, \forall i$
- $d(x, y)$ : Euclidean distance between points  $x$  and  $y$
- $\tilde{m} \in R$ : mass of super-particle as a specific collection of particles.

---

**Algorithm 1** Naive

---

$$1: MBP = \min_i \sum_{j \neq i} \frac{m_j}{d(X_i, X_j)}$$

---

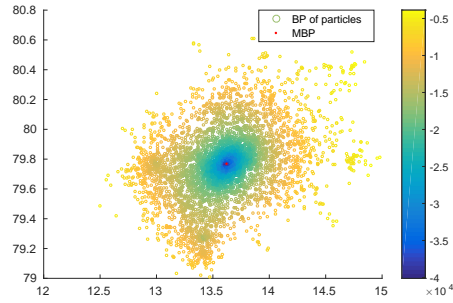


Figure 1: MBP Illustration

---

**Algorithm 2** Mixed Particle/Super-particle Hierarchy

---

```

1: procedure  $MBP = mixed_kmeans(X, m)$ 
2:    $[IDX, c] = kmeans(X, N_c)$ , where  $N_c$  is the number of clusters,  $IDX \in \mathbb{N}^{N_p}$  is the index function to indicate which cluster particle  $X_i$  belongs to,  $c_i \in \mathbb{R}^2$ ,  $i = 1, \dots, N_c$  is the centroid of each cluster.
3:    $\bar{m}_i = |\{j | IDX(j) = i\}|$ ,  $\forall i = 1 : N_c$ 
4:    $SBP_i = \sum_{\{j | IDX(i) \neq j\}} \frac{\bar{m}_j}{d(X_i, c_j)} + \sum_{\{j | IDX(j) = IDX(i), i \neq j\}} \frac{m_j}{d(X_i, X_j)}$ ,  $\forall i$ 
5:    $MBP = \min_i SBP_i$ 
6: end procedure

```

---



---

**Algorithm 3** Super-particle Hierarchy

---

```

1: procedure  $MBP = sp_kmeans(X, m)$ 
2:    $[IDX, c] = kmeans(X, N_c)$ 
3:    $\bar{m}_i = |\{j | IDX(j) = i\}|$ ,  $\forall i = 1 : N_c$ 
4:    $MBP = \min_i \sum_{\{j | IDX(i) \neq j\}} \frac{\bar{m}_j}{d(X_i, c_j)}$ 
5:    $n_c = N_c$ 
6:    $k = 1$ 
7:   while  $k \leq n_k$  do
8:      $MBP_{old} = MBP$ 
9:      $v = \{j | IDX(j) = IDX(MBP)\}$ 
10:     $[\tilde{IDX}, \tilde{c}] = kmeans(X(i), \tilde{N}_c)$ 
11:     $c = \{c_1, \dots, c_{MBP-1}\} \cup \tilde{c}_1 \cup \{c_{MBP+1}, \dots, c_{N_c}\} \cup \{\tilde{c}_2, \dots, \tilde{c}_{\tilde{N}_c}\}$ 
12:     $IDX(v) = \tilde{IDX} + kN_c - 1$ 
13:     $n_c = n_c - 1 + \tilde{N}_c$ 
14:     $\bar{m}_i = \frac{1}{|\{j | IDX(j) = i\}|}$ ,  $\forall i = 1 : n_c$ 
15:     $MBP = \min_i \sum_{\{j | IDX(i) \neq j\}} \frac{\bar{m}_j}{d(X_i, c_j)}$ 
16:    if  $MBP = MBP_{old}$  then
17:      Stop
18:    end if
19:     $k = k + 1$ 
20:  end while
21: end procedure

```

---

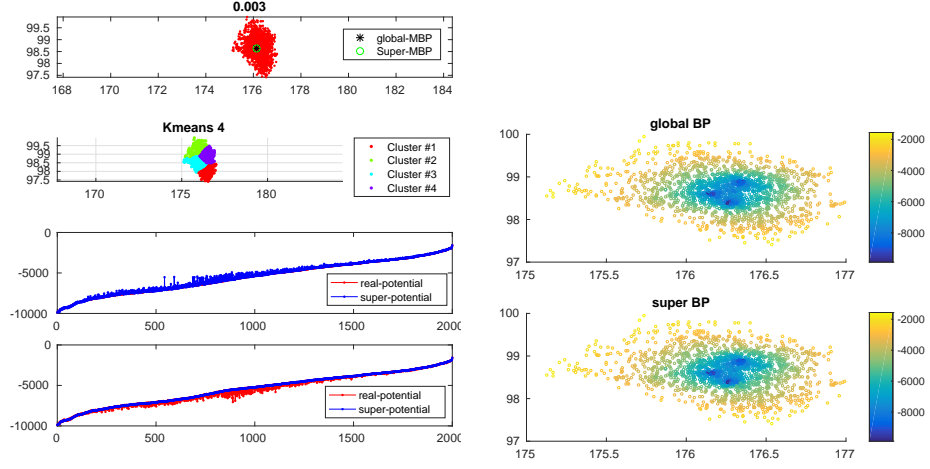


Figure 2: Result of Algorithm 2

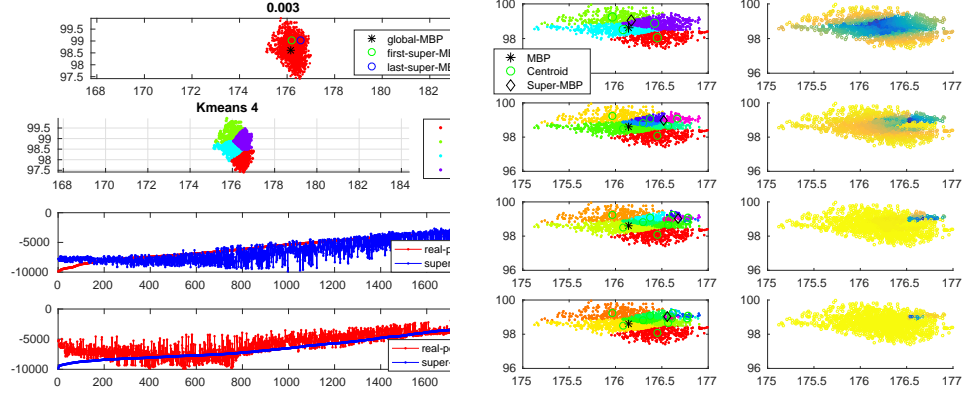


Figure 3: Result of Algorithm 3

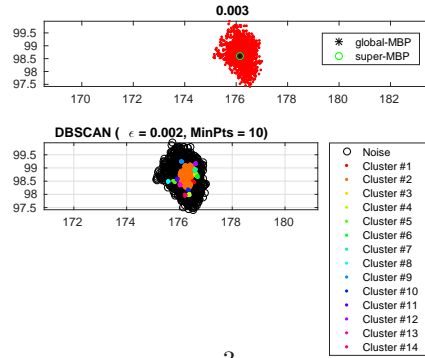


Figure 5: Result of Algorithm 4

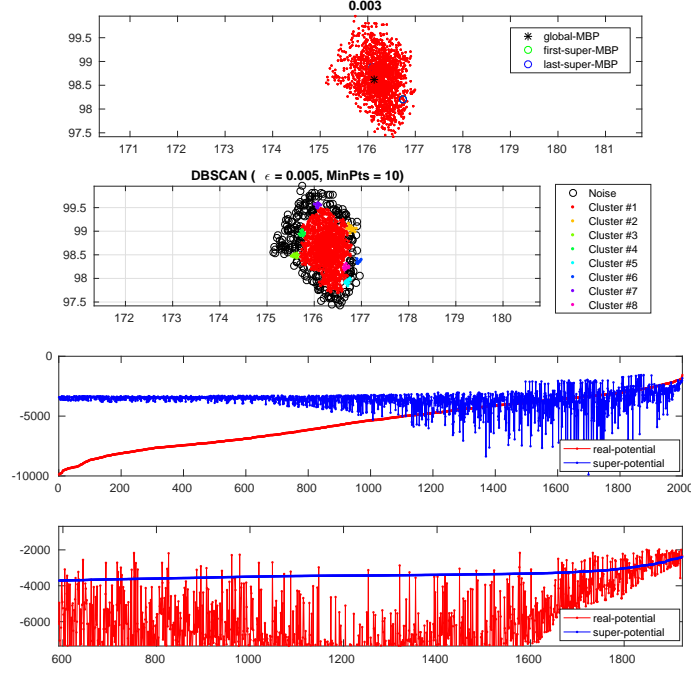


Figure 4: Result of replacing Kmeans by DBSCAN in Algorithm 3

---

**Algorithm 4** Locate MBP from a subset of particles which forms the most dense super-particle via DBSCAN

---

- 1: **procedure**  $MBP = mbp_{dbscan\_max}(X, m)$
  - 2:  $[IDX, c] = dbscan(X, \varepsilon)$ , where  $\varepsilon$  is the linkage length provided for DBSCAN.
  - 3:  $\bar{m}_i = |\{j | IDX(j) = i\}|, \forall i = 1 : N_c$ , where  $N_c$  is the resulted number of clusters given  $\varepsilon$
  - 4:  $i_s = \max_i \bar{m}_i$
  - 5:  $V = \{j | IDX(j) = i_s\}$
  - 6:  $SBP_i = \sum_{\{j | IDX(i) \neq j\}} \frac{\bar{m}_j}{d(X_i, c_j)} + \sum_{\{j | IDX(j) = IDX(i), i \neq j\}} \frac{m_j}{d(X_i, X_j)}, \forall i \in V$
  - 7:  $MBP = \min_{i \in V} SBP_i$
  - 8: **end procedure**
-