

机器学习的数学理论课程项目作业：动态 Kmeans

王英梹

211840213

wangbottlecap@gmail.com

June 19, 2025

Abstract

本报告是 25 Spring 数学学院课程《机器学习的数学理论》项目作业。针对给定的数据集，我们首先实现了经典 Kmeans, Kmeans++, Hartigan & Wong 算法, Elkan 加速 Kmeans, miniBatch Kmeans, Bisecting Kmeans 等多种算法并给出它们的主要思想和算法概述。随后，我们引入动态 kmeans 的概念，考虑在动态情形下对数据进行动态聚类，并给出对应的分析。动态聚类的所有代码作为附录在报告的最后给出。

Contents

1	Classical Kmeans	3
1.1	问题背景与符号表示	3
1.2	目标函数的分块优化思想	3
1.2.1	给定中心后的最优分配	3
1.2.2	给定分配后的最优中心	3
1.3	经典 KMeans 算法迭代流程	4
2	Kmeans++	4
2.1	核心思想	4
2.2	初始化算法的数学步骤	5
2.3	算法描述	5
3	Hartigan & Wong (1979)	5
3.1	逐点搬移的思想	5
3.2	算法步骤	7
4	Elkan 加速版	7
4.1	主要思想	7
4.2	算法框架	9
5	Mini-Batch KMeans	9
5.1	主要思想	9
5.2	算法框架	11
6	Bisecting KMeans	11
6.1	主要思想	11
6.2	算法框架	12
7	动态聚类	13
7.1	问题描述	13
7.2	数学建模	13
7.3	算法流程	14
7.4	复杂度分析	14
7.4.1	增量指派 Complexity of Assignment	14
7.4.2	中心更新 Complexity of Update	14
7.4.3	簇拆分 Complexity of Split	14
7.5	整体复杂度 Overall Complexity	15
8	程序结果	15

1 Classical Kmeans

1.1 问题背景与符号表示

令数据集记为

$$\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d,$$

其中 $\mathbf{x}_i \in \mathbb{R}^d$ 表示第 i 个样本， n 表示总样本数， d 表示特征维度。我们希望将 \mathcal{X} 划分为 k 个不相交的子集（簇）

$$S_1, S_2, \dots, S_k, \quad \bigcup_{i=1}^k S_i = \mathcal{X}, \quad S_i \cap S_j = \emptyset \ (i \neq j),$$

并为每个簇 S_i 找到一个中心 $\mu_i \in \mathbb{R}^d$ ，以最小化以下目标函数（簇内平方误差和）：

$$J(S_1, \dots, S_k; \mu_1, \dots, \mu_k) = \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mu_i\|^2.$$

1.2 目标函数的分块优化思想

KMeans 算法通过分块坐标下降（Block Coordinate Descent）策略，将优化拆分为两个交替步骤：

1. 给定当前中心 μ_1, \dots, μ_k ，求解最优的簇分配 S_1, \dots, S_k 。
2. 给定当前簇分配 S_1, \dots, S_k ，求解最优的聚类中心 μ_1, \dots, μ_k 。

通过在这两个步骤间迭代更新，我们不断降低目标函数，最终在有限步内收敛到某个局部最优解。

1.2.1 给定中心后的最优分配

设当前中心为 $\{\mu_1, \dots, \mu_k\}$ ，则最优的分配策略是将每个样本分到离它最近的中心所在的簇：

$$\mathbf{x} \in S_i \iff i = \arg \min_{j \in \{1, \dots, k\}} \|\mathbf{x} - \mu_j\|.$$

1.2.2 给定分配后的最优中心

设当前的簇分配结果为 S_1, \dots, S_k ，则目标函数可表示为

$$J = \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mu_i\|^2 = \sum_{i=1}^k J_i,$$

其中 $J_i = \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mu_i\|^2$. 要分别最小化各 J_i ，只需对各簇独立求解，即

$$\mu_i^* = \arg \min_{\mu_i} \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mu_i\|^2.$$

可以证明此时的最优解是簇内所有样本的算术平均值：

$$\mu_i^* = \frac{1}{|S_i|} \sum_{\mathbf{x} \in S_i} \mathbf{x}.$$

1.3 经典 KMeans 算法迭代流程

基于上述分块最优思想，KMeans 算法流程可概括如下：

1. **初始化**：随机选择 k 个样本作为初始中心，记为 $\mu_1^{(0)}, \dots, \mu_k^{(0)}$ 。

2. **迭代**：对于第 t 次迭代：

(a) 分配样本 (Assignment Step)：

$$S_i^{(t)} = \left\{ \mathbf{x} \mid i = \arg \min_j \|\mathbf{x} - \mu_j^{(t-1)}\| \right\}, \quad i = 1, \dots, k.$$

(b) 更新中心 (Update Step)：

$$\mu_i^{(t)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x} \in S_i^{(t)}} \mathbf{x}, \quad i = 1, \dots, k.$$

3. **收敛或终止条件**：若中心更新的变化量足够小（小于某个阈值），或达到最大迭代次数，停止迭代并输出结果。

2 Kmeans++

2.1 核心理想

与经典 KMeans 相同，令数据集

$$\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d,$$

其中 $\mathbf{x}_i \in \mathbb{R}^d$ 表示第 i 个样本。我们希望选取 k 个“初始中心”，以改善后续 KMeans 聚类的质量和收敛速度。KMeans++ 的核心思想是：

- 第一个中心等概率地从数据集中随机挑选；
- 后续中心依照与已选中心距离平方成正比的概率进行抽样，距离越远的点被选为新中心的概率越大。

这一策略能够让初始中心彼此分散，从而加速收敛并减少陷入不良局部最优的风险。

2.2 初始化算法的数学步骤

为选取第 k 个中心（假设已经选出 $\mu_1, \mu_2, \dots, \mu_{k-1}$ ）时，执行以下步骤：

1. 计算每个点到最近已选中心的距离平方：对每个 $\mathbf{x}_i \in \mathcal{X}$ ，定义

$$D(\mathbf{x}_i) = \min_{1 \leq j \leq k-1} \|\mathbf{x}_i - \mu_j\|^2.$$

2. 基于距离平方构建抽样分布：

$$P(\mathbf{x}_i) = \frac{D(\mathbf{x}_i)}{\sum_{r=1}^n D(\mathbf{x}_r)}, \quad i = 1, 2, \dots, n.$$

3. 按此分布随机抽样：随机生成一个 \mathbf{x}_q 使 $\Pr[q = i] = P(\mathbf{x}_i)$ ，并令 $\mu_k = \mathbf{x}_q$.

当所有 k 个中心均选出后，便得到一组初始中心。

2.3 算法描述

Algorithm 1 KMeans++ Initialization

Require: Data set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$, number of centers k .

Ensure: k initial centers μ_1, \dots, μ_k .

- 1: Randomly pick one sample \mathbf{x}_p from \mathcal{X} , set $\mu_1 \leftarrow \mathbf{x}_p$.
 - 2: **for** $t = 2 \rightarrow k$ **do**
 - 3: For each i , $D(\mathbf{x}_i) \leftarrow \min_{1 \leq j \leq t-1} \|\mathbf{x}_i - \mu_j\|^2$.
 - 4: $P(\mathbf{x}_i) \leftarrow \frac{D(\mathbf{x}_i)}{\sum_{r=1}^n D(\mathbf{x}_r)}$.
 - 5: Sample \mathbf{x}_q from \mathcal{X} with probability $P(\mathbf{x}_i)$.
 - 6: $\mu_t \leftarrow \mathbf{x}_q$
 - 7: **end for**
 - 8: **return** μ_1, \dots, μ_k
-

3 Hartigan & Wong (1979)

3.1 逐点搬移的思想

设当前数据集为

$$\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d,$$

被划分成 k 个簇 S_1, S_2, \dots, S_k ，对应中心为 $\mu_1, \mu_2, \dots, \mu_k$ 。我们考虑将单个样本 \mathbf{x} 从其当前所在簇 A 移动到另一个簇 B 时，对目标函数

$$J = \sum_{i=1}^k \sum_{\mathbf{z} \in S_i} \|\mathbf{z} - \mu_i\|^2$$

所带来的变化量。具体地，令 S_A, S_B 分别表示簇 A 和簇 B 的成员集合：

$$S'_A = S_A \setminus \{\mathbf{x}\}, \quad S'_B = S_B \cup \{\mathbf{x}\}.$$

相应的新中心记为

$$\boldsymbol{\mu}'_A = \frac{1}{|S_A| - 1} \sum_{\mathbf{y} \in S_A \setminus \{\mathbf{x}\}} \mathbf{y}, \quad \boldsymbol{\mu}'_B = \frac{1}{|S_B| + 1} \left(\sum_{\mathbf{z} \in S_B} \mathbf{z} + \mathbf{x} \right).$$

若将样本 \mathbf{x} 从簇 A 搬移到簇 B 后导致目标函数值严格下降，即

$$\Delta = \left(\sum_{\mathbf{z} \in S'_A} \|\mathbf{z} - \boldsymbol{\mu}'_A\|^2 + \sum_{\mathbf{z} \in S'_B} \|\mathbf{z} - \boldsymbol{\mu}'_B\|^2 \right) - \left(\sum_{\mathbf{z} \in S_A} \|\mathbf{z} - \boldsymbol{\mu}_A\|^2 + \sum_{\mathbf{z} \in S_B} \|\mathbf{z} - \boldsymbol{\mu}_B\|^2 \right) < 0,$$

则执行此“逐点搬移”可使 J 单调下降，从而向局部最优解靠近。

3.2 算法步骤

Algorithm 2 Simple Hartigan & Wong KMeans (Point Reassignment)

Require: Data set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, number of clusters k , max iterations M .

Ensure: Final cluster centers μ_1, \dots, μ_k and assignments S_1, \dots, S_k .

```
1: Initialize centers  $\mu_1, \dots, \mu_k$  using some method (e.g., KMeans++).
2: Assign each  $\mathbf{x}_i$  to the nearest center.
3: for  $t = 1 \rightarrow M$  do
4:    $changed \leftarrow \text{False}$ 
5:   for each sample  $\mathbf{x}$  in  $\mathcal{X}$  do
6:     Let  $A$  be the cluster containing  $\mathbf{x}$ , and  $\mu_A$  be its center.
7:      $\delta_{\text{best}} \leftarrow 0$ ,  $B_{\text{best}} \leftarrow A$ 
8:     for each cluster  $B \neq A$  do
9:       Compute the new centers if  $\mathbf{x}$  were moved from  $A$  to  $B$ :
           $\mu'_A$  and  $\mu'_B$ .
10:      Evaluate  $\Delta = J_{\text{after}} - J_{\text{before}}$  for the local change (see text).
11:      if  $\Delta < \delta_{\text{best}}$  then
12:         $\delta_{\text{best}} \leftarrow \Delta$ 
13:         $B_{\text{best}} \leftarrow B$ 
14:      end if
15:    end for
16:    if  $B_{\text{best}} \neq A$  then
17:      Move  $\mathbf{x}$  from  $A$  to  $B_{\text{best}}$ , update  $\mu_A$  and  $\mu_{B_{\text{best}}}$ .
18:       $changed \leftarrow \text{True}$ 
19:    end if
20:  end for
21:  if not  $changed$  then
22:    break
23:  end if
24: end for
25: return final centers  $\{\mu_1, \dots, \mu_k\}$  and cluster sets  $\{S_1, \dots, S_k\}$ .
```

4 Elkan 加速版

4.1 主要思想

与经典 KMeans 相同，假设数据集为

$$\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d,$$

聚类中心为 $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$, 目标函数为

$$J = \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mathbf{c}_i\|^2.$$

Elkan 加速版的关键在于利用对偶边界（上下界）和三角不等式减少不必要的距离计算。大体上，为每个样本 \mathbf{x} 维护：

- **上界 $u(\mathbf{x})$** ：估计 \mathbf{x} 与其当前分配中心 $\mathbf{c}_{\text{assigned}(\mathbf{x})}$ 的距离，

$$u(\mathbf{x}) \geq \|\mathbf{x} - \mathbf{c}_{\text{assigned}(\mathbf{x})}\|.$$

- **下界 $l_j(\mathbf{x})$** ：对任意中心 \mathbf{c}_j 的距离，

$$l_j(\mathbf{x}) \leq \|\mathbf{x} - \mathbf{c}_j\|.$$

记中心间的距离为

$$d_{ij} = \|\mathbf{c}_i - \mathbf{c}_j\|.$$

若对样本 \mathbf{x} 有

$$u(\mathbf{x}) < \frac{d_{ij}}{2},$$

则可推断 $\|\mathbf{x} - \mathbf{c}_j\| > \|\mathbf{x} - \mathbf{c}_i\|$, 因此中心 j 不会是 \mathbf{x} 更好的分配对象，无需计算真实距离。三角不等式保障了此结论的正确性，从而实现大规模剪枝。此外，当中心 \mathbf{c}_i 在两次迭代间移动量 m_i 较小，则可依据

$$u(\mathbf{x}) \leftarrow u(\mathbf{x}) + m_i$$

等方式修正上界，而依旧保留安全的剪枝判断。此即 Elkan 加速算法的根本逻辑：**以有限的上/下界更新与三角不等式将大量“无关距离”排除在外**，从而减少时间复杂度。

4.2 算法框架

Algorithm 3 Elkan Accelerated KMeans

Require: Data set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, number of clusters k , max iterations M .

Ensure: Cluster centers $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ and partition S_1, \dots, S_k .

- 1: Initialize centers $\mathbf{c}_1, \dots, \mathbf{c}_k$ (e.g., KMeans++), and assign each sample to the nearest center.
 - 2: For each \mathbf{x}_i , set $u(\mathbf{x}_i) \leftarrow \|\mathbf{x}_i - \mathbf{c}_{\text{assigned}(\mathbf{x}_i)}\|$, and initialize $l_j(\mathbf{x}_i)$ for $j \neq \text{assigned}(\mathbf{x}_i)$.
 - 3: **for** $t = 1 \rightarrow M$ **do**
 - 4: **Update bounds:** compute center movements $m_i = \|\mathbf{c}_i^{(\text{old})} - \mathbf{c}_i^{(\text{new})}\|$ and correct $u(\mathbf{x})$ or $l_j(\mathbf{x})$ accordingly.
 - 5: **Assignment step:**
 - 6: **for** each sample \mathbf{x}_i **do**
 - 7: Let $i_0 = \text{assigned}(\mathbf{x}_i)$.
 - 8: **for** each center $j \neq i_0$ **do**
 - 9: **Pruning check:** if $u(\mathbf{x}_i) < d_{i_0 j}/2$, skip distance computation for center j .
 - 10: Otherwise, compute $\|\mathbf{x}_i - \mathbf{c}_j\|$ and update $u(\mathbf{x}_i), l_j(\mathbf{x}_i)$ if needed.
 - 11: **end for**
 - 12: Re-assign \mathbf{x}_i to the closest center found if it differs from i_0 .
 - 13: **end for**
 - 14: **Update step:** recalculate each center $\mathbf{c}_j \leftarrow \frac{1}{|S_j|} \sum_{\mathbf{x} \in S_j} \mathbf{x}$.
 - 15: **Check convergence:** if center changes are small enough or no re-assignments, break.
 - 16: **end for**
 - 17: **return** $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ and final partition $\{S_1, \dots, S_k\}$.
-

5 Mini-Batch KMeans

5.1 主要思想

与经典 KMeans 相同，假设数据集

$$\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d,$$

我们希望寻找 k 个中心 μ_1, \dots, μ_k 并将 \mathcal{X} 划分为 k 个不相交子集 S_1, \dots, S_k ，使得以下目标函数（簇内平方误差和）最小化：

$$\min_{S_1, \dots, S_k; \mu_1, \dots, \mu_k} \sum_{j=1}^k \sum_{\mathbf{x} \in S_j} \|\mathbf{x} - \mu_j\|^2.$$

然而，在大规模数据场景或流式数据场景中，每次迭代都处理全部样本 $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ 成本非常高。为此，Mini-Batch KMeans 在每次迭代只抽取一个 **小批量**（mini-batch），记为

$$B^{(t)} = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_m}\} \subset \mathcal{X}, \quad |B^{(t)}| = m \ll n,$$

并执行类似 KMeans 的分配与更新步骤，但仅基于 $B^{(t)}$ 更新当前的中心 $\{\mu_j\}$ 。

令在 $B^{(t)}$ 中分配到簇 j 的那些样本记为

$$B_j^{(t)} = \{\mathbf{x} \in B^{(t)} \mid \text{assigned to } \mu_j\},$$

则 μ_j 的更新常采取如下形式：

$$\mu_j^{(t+1)} = \mu_j^{(t)} + \eta \left(\bar{\mathbf{x}}_{\text{batch}, j} - \mu_j^{(t)} \right),$$

其中

$$\bar{\mathbf{x}}_{\text{batch}, j} = \frac{1}{|B_j^{(t)}|} \sum_{\mathbf{x} \in B_j^{(t)}} \mathbf{x}, \quad \eta \in (0, 1]$$

为更新步长（又可视为学习率）。这使得中心向本批次样本的均值方向移动，而无需处理全部 n 个样本。

通过在每次迭代使用较小规模的批数据 $B^{(t)}$ ，算法单次迭代的复杂度降至 $O(mk)$ ，明显低于经典 KMeans 的 $O(nk)$ ；同时在多轮迭代后，聚类解可近似逼近全量 KMeans 的结果。Mini-Batch KMeans 在大规模数据或在线学习环境中非常实用，能在保证一定精度下显著降低计算开销。

5.2 算法框架

Algorithm 4 Mini-Batch KMeans

Require: Data set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$, number of clusters k , mini-batch size m , max iterations M .

Ensure: Final centers $\{\mu_1, \dots, \mu_k\}$ and assignments S_1, \dots, S_k .

```

1: Initialize centers  $\mu_1, \dots, \mu_k$  (e.g., random or KMeans++).
2: for  $t = 1 \rightarrow M$  do
3:   Sample a mini-batch  $B^{(t)} = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_m}\} \subset \mathcal{X}$ .
4:   Assignment step (on mini-batch):
5:   for each  $\mathbf{x} \in B^{(t)}$  do
6:     Find  $j^* = \arg \min_j \|\mathbf{x} - \mu_j\|^2$ .
7:     Assign  $\mathbf{x}$  to cluster  $j^*$ .
8:   end for
9:   Update step (on mini-batch):
10:  for  $j = 1 \rightarrow k$  do
11:    Let  $B_j^{(t)}$  be the subset of  $B^{(t)}$  assigned to center  $j$ .
12:    if  $B_j^{(t)} \neq \emptyset$  then
13:       $\bar{\mathbf{x}}_j^{(t)} \leftarrow \frac{1}{|B_j^{(t)}|} \sum_{\mathbf{x} \in B_j^{(t)}} \mathbf{x}$ .
14:       $\mu_j \leftarrow \mu_j + \eta (\bar{\mathbf{x}}_j^{(t)} - \mu_j)$ ,
15:      where  $\eta$  is a learning rate in  $(0, 1]$ .
16:    end if
17:  end for
18:  (Optional) Check convergence if center movement is small or no changes.
19: end for
20: Final assignment: reassign every  $\mathbf{x}_i \in \mathcal{X}$  to its nearest center  $\mu_j$ .
21: return  $\{\mu_1, \dots, \mu_k\}$  and  $S_1, \dots, S_k$ .

```

6 Bisecting KMeans

6.1 主要思想

Bisecting KMeans (又称二分 KMeans) 将 KMeans 与层次聚类思想结合, 通过 **自顶向下** (top-down) 的分裂方式得到 k 个簇。初始时将所有数据视为一个大簇, 然后反复挑选其中“最需要细分”的簇, 做 2-means 将其分成两个簇, 直至达到 k 个簇。

具体而言, 令 \mathcal{X} 初始在一个簇 $S_{\text{root}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ 中, 中心

$$\mu_{\text{root}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i.$$

若当前聚类数小于 k , 则:

1. 选出一个簇 S_{\max} (可根据簇内误差最大的原则), 其中包含若干样本 $\{\mathbf{x}_{a_1}, \dots, \mathbf{x}_{a_r}\}$ 。
2. 在该簇上运行 **2-means** ($k = 2$ 的 KMeans), 解为

$$\min_{S_a, S_b, \mu_a, \mu_b} \left(\sum_{\mathbf{x} \in S_a} \|\mathbf{x} - \mu_a\|^2 + \sum_{\mathbf{x} \in S_b} \|\mathbf{x} - \mu_b\|^2 \right),$$

使 $S_a \cup S_b = S_{\max}$, $S_a \cap S_b = \emptyset$ 。

3. 用 S_a 与 S_b 替换 S_{\max} , 整体簇数加 1。

通过这种不断“二分”最大的簇, Bisecting KMeans 最终可得到 k 个簇。其核心思想是:

逐次选出最需要被拆分的簇, 用 2-means 精细地对它进行划分。

在每次二分时, 仍遵循 KMeans 的优化目标, 即最小化簇内误差和, 只不过将问题限制在某个子簇 S_{\max} 上。这种自顶向下的层次聚类与 KMeans 结合的方式在实践中具有较好的可解释性和灵活度, 同时每次只在局部簇做细分, 能在一定程度上减少计算量。虽然与一次性 k -means 一样只能得到局部最优, 但常能更好地针对层次结构做可视化与分析。

6.2 算法框架

Algorithm 5 Bisecting KMeans (二分聚类)

Require: Data set $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$, number of clusters k , max iterations for each 2-means M_2 .

Ensure: k final clusters S_1, \dots, S_k with centers $\{\mu_1, \dots, \mu_k\}$.

- 1: Initialize a single cluster $S_{\text{root}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$.
 - 2: Let the set of active clusters $C \leftarrow \{S_{\text{root}}\}$.
 - 3: **while** $|C| < k$ **do**
 - 4: Select the cluster $S_{\max} \in C$ with the largest SSE (or largest cardinality).
 - 5: Remove S_{\max} from C .
 - 6: **Run 2-means** on S_{\max} (e.g., use Lloyd or Mini-Batch with $k = 2$, M_2 iterations):
 - 7: Split S_{\max} into two subsets S_a and S_b with centers μ_a and μ_b .
 - 8: Add S_a and S_b to C .
 - 9: **end while**
 - 10: Let $\{S_1, \dots, S_k\} = C$, and compute each center as $\mu_j \leftarrow \frac{1}{|S_j|} \sum_{\mathbf{x} \in S_j} \mathbf{x}$.
 - 11: **return** S_1, \dots, S_k and $\{\mu_1, \dots, \mu_k\}$.
-

7 动态聚类

7.1 问题描述

随着海量数据以流式方式到达，本节所述算法旨在对数据进行动态聚类。假设数据流被划分为等长批次，每批含 $B = 100$ 条样本。记在时刻 t 到达的新批次为 $\Delta\mathcal{D}_t$ ，总体数据集为

$$\mathcal{D}_t = \bigcup_{s=1}^t \Delta\mathcal{D}_s.$$

我们要求：

- 每批到达后立即更新聚类结构；
- 保证每个簇包含的数据条数在 $[N_{\min}, N_{\max}] = [60, 150]$ 范围内；
- 若某簇大小超过 N_{\max} ，则将其拆分为两个子簇；
- (可扩展) 若删除过期数据后某簇小于 N_{\min} ，则与最近邻簇合并。

7.2 数学建模

令第 t 时刻的聚类中心集合为

$$C_t = \{c_{t,1}, \dots, c_{t,K_t}\}, \quad c_{t,i} \in \mathbb{R}^d,$$

样本点集合被分配至簇 $\mathcal{D}_{t,i} = \{x \mid \text{assign}(x) = i\}$ 。我们仍以 KMeans 惯性 (Inertia) 作为损失函数：

$$\mathcal{L}_t = \sum_{i=1}^{K_t} \sum_{x \in \mathcal{D}_{t,i}} \|x - c_{t,i}\|^2.$$

动态更新包含以下运算：

1. **增量指派**：对新到达点 $x \in \Delta\mathcal{D}_t$,

$$\text{assign}(x) = \arg \min_{1 \leq i \leq K_{t-1}} \|x - c_{t-1,i}\|.$$

2. **中心更新**：重新计算

$$c_{t,i} = \frac{1}{|\mathcal{D}_{t,i}|} \sum_{x \in \mathcal{D}_{t,i}} x.$$

3. **簇拆分**：对所有 i 满足 $|\mathcal{D}_{t,i}| > N_{\max}$ ，在该簇内部再执行一次 $k = 2$ 的 KMeans，得子中心 $\{c'_{i,1}, c'_{i,2}\}$ ，用 $c'_{i,1}$ 替换原中心并新增 $c'_{i,2}$ 。

Algorithm 6 Dynamic KMeans Clustering Algorithm

Require: Data stream $\{\Delta\mathcal{D}_t\}_{t=1}^T$, batch size B , thresholds N_{\min}, N_{\max}

Ensure: Final cluster centers C_T

```
1:  $t \leftarrow 1$ 
2: Perform standard KMeans on  $\Delta\mathcal{D}_1$  to obtain  $C_1, A_1$ 
3: for  $t = 2, 3, \dots, T$  do
4:   Read new batch  $\Delta\mathcal{D}_t$ 
5:   for all  $x \in \Delta\mathcal{D}_t$  do
6:      $\text{assign}(x) \leftarrow \arg \min_i \|x - c_{t-1,i}\|$ 
7:   end for
8:   Update centers:  $c_{t,i} \leftarrow \frac{1}{|\mathcal{D}_{t,i}|} \sum_{x \in \mathcal{D}_{t,i}} x$ 
9:   for all clusters  $i$  with  $|\mathcal{D}_{t,i}| > N_{\max}$  do
10:    Run 2-cluster KMeans on  $\mathcal{D}_{t,i}$  to get  $\{c'_{i,1}, c'_{i,2}\}$ 
11:    Replace  $c_{t,i} \leftarrow c'_{i,1}$  and add new center  $c'_{i,2}$ 
12:   end for
13: end for
```

7.3 算法流程

7.4 复杂度分析

该算法包含三个主要步骤：增量指派、中心更新、簇拆分，现逐一给出详细复杂度推导。

7.4.1 增量指派 Complexity of Assignment

在第 t 批处理中，新数据量为 B ，当前簇数为 K_{t-1} ，数据维度为 d 。对每个点计算到 K_{t-1} 个中心的欧氏距离，时间复杂度为

$$O(B \times K_{t-1} \times d).$$

由于 K_{t-1} 保持在总数据条数 n_t 的 1% 左右，即 $K_{t-1} \approx 0.01 n_t / B$ ，故该步骤总体接近线性。

7.4.2 中心更新 Complexity of Update

中心更新需遍历所有簇及其成员，总数据量为 $n_t = \sum_i |\mathcal{D}_{t,i}|$ ，对每簇计算均值耗时 $O(|\mathcal{D}_{t,i}| \times d)$ ，合计

$$\sum_{i=1}^{K_t} O(|\mathcal{D}_{t,i}| \times d) = O(n_t \times d).$$

7.4.3 簇拆分 Complexity of Split

最坏情况下，每个簇都可能超过阈值并被拆分，拆分过程对簇内点再执行一次 $k = 2$ KMeans。对簇 i ，若其规模为 m_i ，一次 2-簇 KMeans 时间为 $O(m_i \times d)$ （假设迭代次数有限常数次）。合

并所有簇的拆分成本为

$$\sum_{i=1}^{K_t} O(m_i \times d) = O(n_t \times d).$$

7.5 整体复杂度 Overall Complexity

综合上述三步，每批次迭代的总复杂度为

$$O(BK_{t-1}d + n_t d + n_t d) = O((BK_{t-1} + 2n_t)d).$$

鉴于 $B \ll n_t$ 且 $K_{t-1} = O(n_t/B)$ ，可简化为

$$O(n_t d),$$

即对流式数据的聚类更新具有近线性时间复杂度，适用于大规模在线场景。

8 程序结果

在给定迭代步数和给定收敛时间的收敛曲线中，Minibatch 和 Bisecting 聚类算法都比经典 kmeans 有了明显的提速。同时，由于算法自身设计，Bisecting 的聚类速度相对更稳定，MiniBatch 算法更加依赖挑选 Batch 的策略。

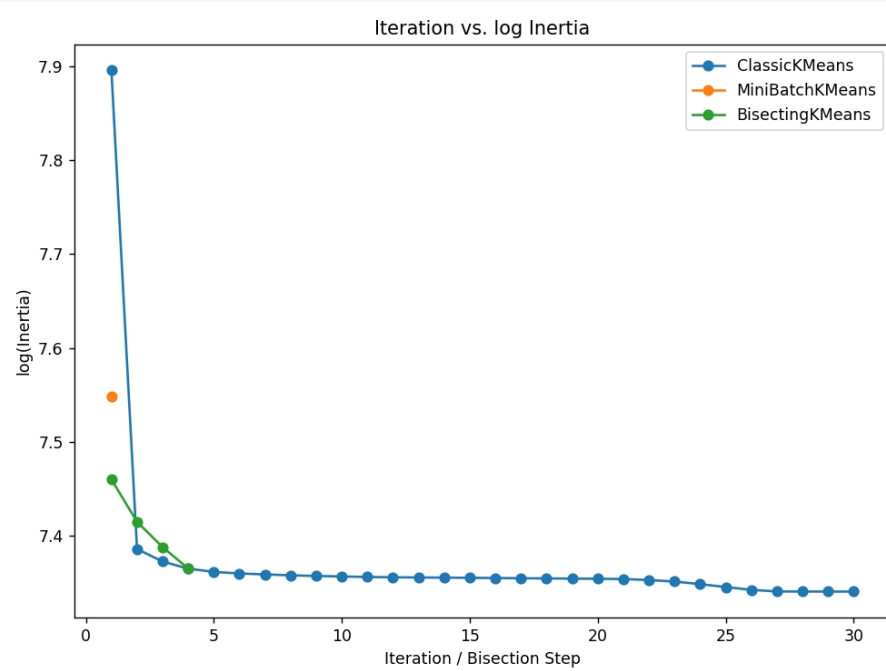


Figure 1: IterVSlogLoss

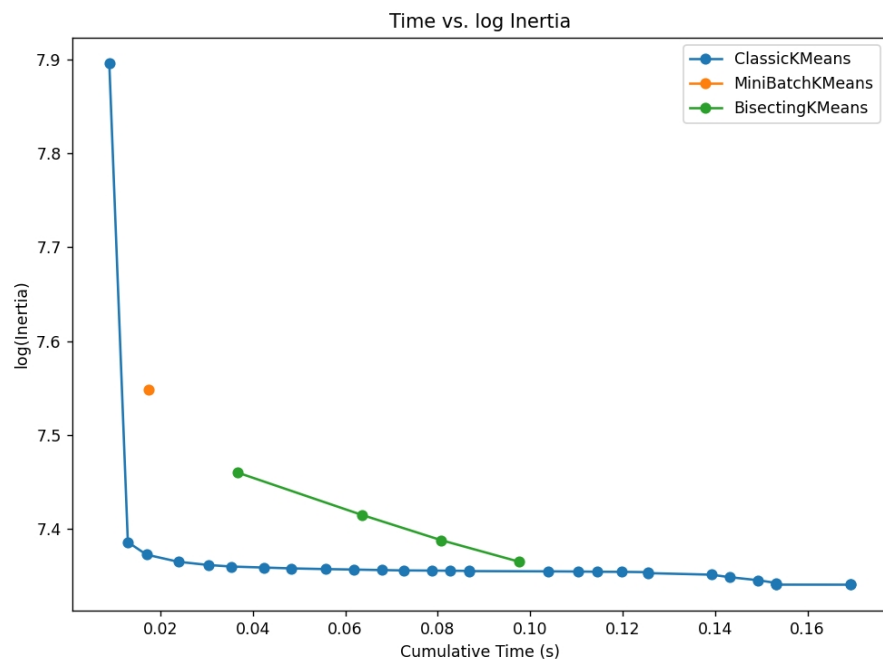


Figure 2: TimeVSlogLoss

聚类效果如下：

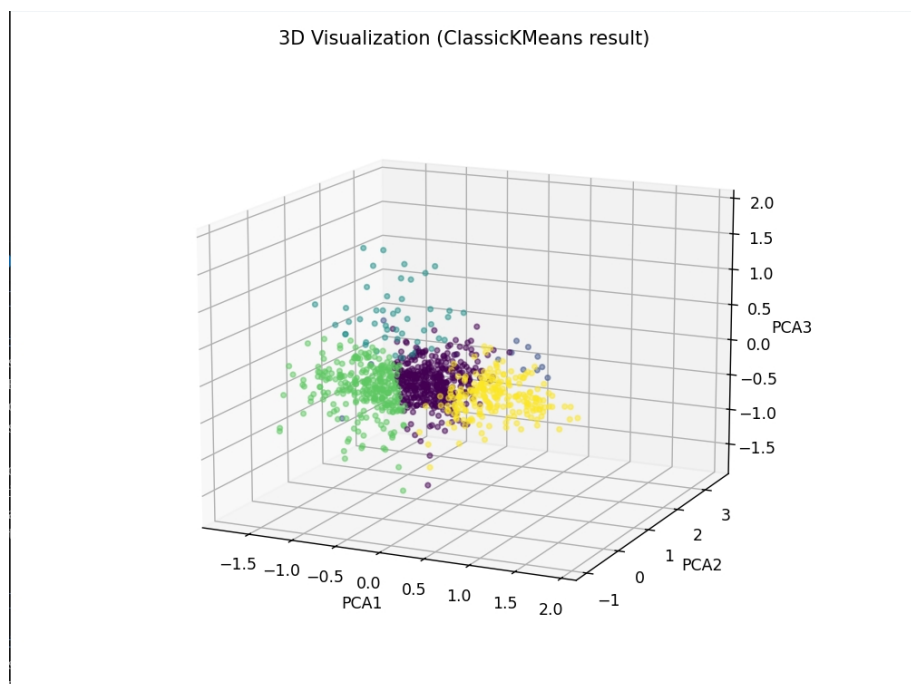


Figure 3: TimeVSlogLoss

随后，针对动态聚类问题，我们计算得出在相同的中心数量下，动态方法和静态方法平均最近距离为 0.3699。随着引入数据的增大，动态方法的 Loss 函数也在递增，这种递增方式无

法避免，因而需要更加具有代表性的 loss 函数设计以表达在动态过程中聚类的效果；时间方面，动态方法引入时间较长，这一方面是因为动态方法需要多次调用 `kmeans` 和裂分裂函数，另外一方面是因为动态方法也会对原始数据进行重复读取。当使用数据量较大（如百万条以上）时，动态方法因其快速迭代初始解和分批读取数据对内存需求较低等特性会相较本次实验出现不同的表现。

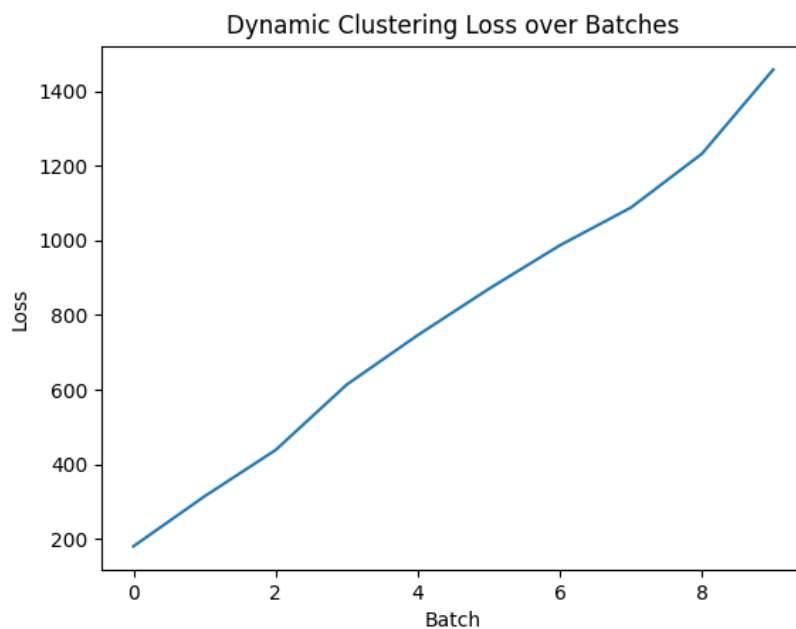


Figure 4: Dynamic Method: Loss

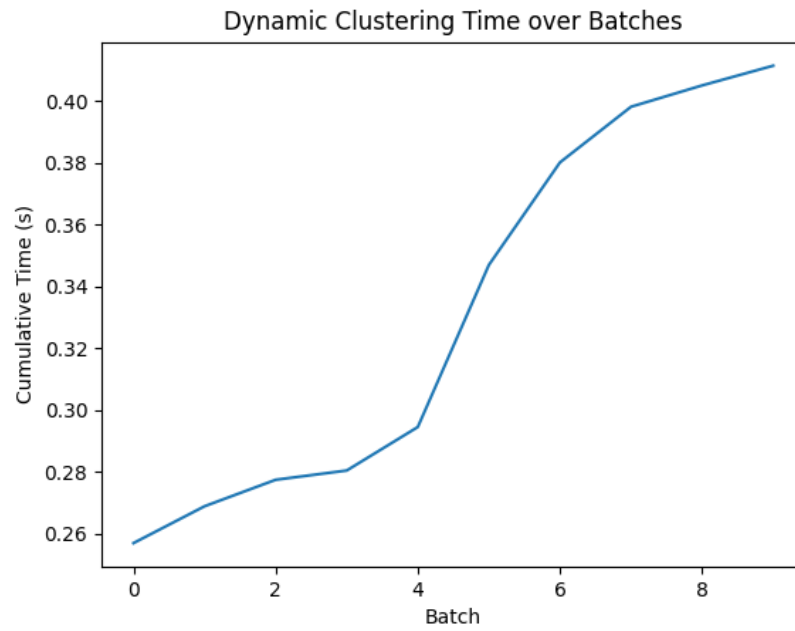


Figure 5: Dynamic Method: Time

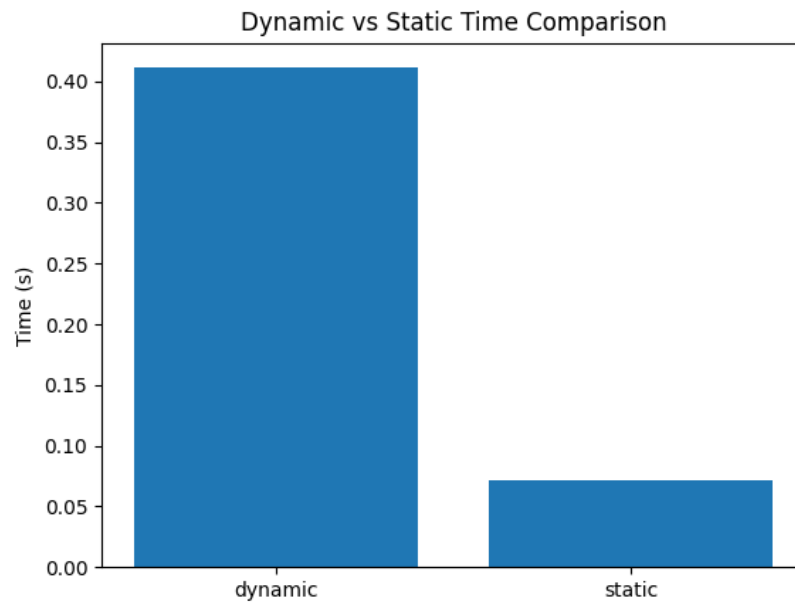


Figure 6: Time Compare