# DynCluster

*Release 0.0.1*

**Yingxiao Wang**

**Jun 19, 2025**

# CONTENTS:

Add your content using `reStructuredText` syntax. See the [reStructuredText](#) documentation for details.

# SOURCE CODE

## 1.1 `main.py`

Listing 1: main.py

```python
"""
This project is a course project for Prof. Qinwu Xu's "Mathematical Theory of Machine␣
↪Learning" (Spring 2025), School of Mathematics, Nanjing University. It mainly␣
↪focuses on the problem of dynamic clustering.

The core process is as follows:
→ New data stream arrives
→ Perform clustering on the new data
→ Update cluster centers and adjust previous clustering results
→ If any cluster becomes too large, it is split
→ Adjust the clustering result accordingly
→ Wait for the next batch of data

Author: Yingxiao Wang
Email: wangbottlecap@gmail.com
"""

import logger
from constants.constants import *
from backgrounds import DataReader
from optimizer import DynamicKmeansController, StaticKmeansController
from processer import ComparisonProcessor
from plotter import ComparisonPlotter

def main():
    """
```

```python
25
26      """
27      LOGGER = logger.setup_logging()
28      LOGGER.info(f"This project is a course project for Prof. Qinwu Xu's Mathematical
   →Theory of Machine Learning (Spring 2025), School of Mathematics, Nanjing University.
   → It mainly focuses on the problem of dynamic clustering.")
29
30      rd = DataReader()
31      rd.read_fvecs()
32      print(rd.data)
33
34      dyn = DynamicKmeansController(data=rd.data)
35      dyn.run()
36
37      sta = StaticKmeansController(data=rd.data, k=dyn.getCentroidsNum())
38      sta.run()
39
40      processer = ComparisonProcessor(dyn, sta)
41      plotter = ComparisonPlotter(processer)
42      plotter.print_center_distance()
43      plotter.plotAll()
44
45      return
46
47
48  if __name__ == "__main__":
49      main()
```

## 1.2 `optimizer.py`

Listing 2: optimizer.py

```python
1  import numpy as np
2  import time
3  from sklearn.cluster import KMeans
4
5  import logging
6  LOGGER = logging.getLogger(__name__)
7
```

```python
8   class DynamicKmeansController():
9       """
10      Dynamic KMeans clustering controller for streaming data.
11
12      Core process:
13      → New data batch arrives
14      → Assign new points to existing clusters (or initialize)
15      → Update cluster centers
16      → Split clusters exceeding max_size into two
17      → Compute and print current loss
18      → Repeat for next batch
19
20      Attributes:
21          data (np.ndarray): full dataset for indexing.
22          batch_size (int): number of points per incoming batch.
23          min_size (int): minimum cluster size before merge (future extension).
24          max_size (int): maximum cluster size before split.
25          centroids (list[np.ndarray]): current cluster centers.
26          cluster_data_indices (dict[int, list[int]]): mapping cluster IDs to data␣
    ↪indices.
27          assignments (dict[int, int]): mapping data index to its cluster ID.
28          loss_history (list[float]): recorded loss after each batch.
29          time_history (list[float]): cumulative processing time after each batch.
30      """
31
32      def __init__(self, data:np.ndarray, batch_size=100, min_size=60, max_size=150,␣
    ↪ratio = 0.01,threshold = 0.01):
33          self.data = data                                    # 保存整个数据集
34          self.batch_size = batch_size                        # 每次处理的批大小
35          # self.threshold = threshold
36          # self.ratio = ratio
37          self.min_size = min_size                            # 最小簇大小（留作合并条件）
38          self.max_size = max_size                            # 最大簇大小（超过需分裂）
39          self.centroids = []                                 # 存储当前簇中心列表
40          self.cluster_data_indices = {}                      # 存储簇 ID 到样本索引的映射
41          self.assignments = {}                               # 存储样本索引到簇 ID 的映射
42          self.loss_history = []                              # 记录每批的损失值（inertia）
43          self.time_history = []                              # 记录累计时间
44          self._current_time = 0.0                            # 内部跟踪累计运行时间
45
```

```python
46      def getCentroidsNum(self): return len(self.centroids)

47

48      def run(self):
49          """Run dynamic clustering over the entire dataset in streaming batches."""
50          n_points = self.data.shape[0]                        # 数据总样本数
51          # 按批次遍历数据
52          for batch_start in range(0, n_points, self.batch_size):
53              # 计算本批样本全局索引范围
54              batch_indices = list(range(batch_start, min(batch_start + self.batch_size,
    → n_points)))
55              batch = self.data[batch_indices]              # 提取本批数据
56              t0 = time.time()                              # 计时开始

57

58              # 增量聚类、更新中心、拆分过大簇
59              self._add_batch(batch, batch_indices)
60              self._update_clusters()
61              self._split_large_clusters()

62

63              batch_loss = self._compute_loss()             # 计算本批后总损失
64              self._current_time += time.time() - t0        # 更新累计时间

65

66              self.loss_history.append(batch_loss)          # 记录损失
67              self.time_history.append(self._current_time)  # 记录时间
68              print(f"Batch {batch_start//self.batch_size + 1}: Loss = {batch_loss:.4f},
    → Time = {self._current_time:.4f}s")

69

70      def _add_batch(self, batch, indices):
71          """Assign a new batch of points to clusters or initialize clustering."""
72          if not self.centroids:                            # 如果还没有簇中心（首次调用）
73              # 计算初始簇数：确保至少 1 个簇
74              k_init = max(1, len(batch) // ((self.min_size + self.max_size) // 2))
75              # k_init = max(1, int(self.batch_size * self.ratio))
76              # 使用 KMeans 对首批数据做静态聚类初始化
77              km = KMeans(n_clusters=k_init, init='k-means++', random_state=42).
    →fit(batch)
78              self.centroids = list(km.cluster_centers_)    # 保存初始簇中心
79              # 记录每个样本的簇标签
80              for idx, lbl in zip(indices, km.labels_):
81                  self.assignments[idx] = lbl
82              # 构建簇到样本索引的映射
```

```python
        for cid in range(k_init):
            self.cluster_data_indices[cid] = [idx for idx in indices if self.
→assignments[idx] == cid]
    else:
        # 对后续批次，遍历每个新点
        for idx, point in zip(indices, batch):
            # 计算该点到所有簇中心的欧氏距离
            distances = [np.linalg.norm(point - c) for c in self.centroids]
            lbl = int(np.argmin(distances))          # 选择最近的簇
            self.assignments[idx] = lbl              # 更新样本–簇映射
            # 将样本索引追加到对应簇的列表中
            self.cluster_data_indices.setdefault(lbl, []).append(idx)


def _update_clusters(self):
    """Recompute centroids based on current cluster memberships."""
    # 遍历所有簇
    for cid, members in self.cluster_data_indices.items():
        # 提取该簇对应的所有样本点
        pts = np.array([self.data[i] for i in members])
        if pts.size > 0:
            # 重新计算簇中心为样本均值
            self.centroids[cid] = pts.mean(axis=0)


def _split_large_clusters(self):
    """Split any cluster larger than max_size into two subclusters."""
    # 下一个可用簇 ID
    next_cid = max(self.cluster_data_indices.keys(), default=-1) + 1
    # 找出需要分裂的簇 ID 列表
    oversized = [cid for cid, members in self.cluster_data_indices.items() if
→len(members) > self.max_size]
    for cid in oversized:
        members = self.cluster_data_indices.pop(cid)  # 暂时移除该簇
        pts = np.array([self.data[i] for i in members])
        # 对过大簇内样本再次做 2 簇 KMeans 拆分
        km2 = KMeans(n_clusters=2, init='k-means++', random_state=42).fit(pts)
        centers2 = km2.cluster_centers_
        labels2 = km2.labels_

        # 更新原簇中心及成员
        self.centroids[cid] = centers2[0]
```

```python
121            self.cluster_data_indices[cid] = [members[i] for i, l in
       enumerate(labels2) if l == 0]
122
123            # 添加新簇中心及成员
124            self.centroids.append(centers2[1])
125            self.cluster_data_indices[next_cid] = [members[i] for i, l in
       enumerate(labels2) if l == 1]
126
127            # 更新分裂后所有样本的簇归属
128            for idx, lbl in zip(members, labels2):
129                self.assignments[idx] = cid if lbl == 0 else next_cid
130
131            next_cid += 1                              # 更新下一个可用簇 ID
132
133    def _compute_loss(self):
134        """Compute total inertia (sum of squared distances) as loss."""
135        total_loss = 0.0
136        # 遍历每个簇计算簇内平方和
137        for cid, members in self.cluster_data_indices.items():
138            center = self.centroids[cid]              # 当前簇中心
139            pts = np.array([self.data[i] for i in members])
140            total_loss += np.sum((pts - center) ** 2)    # 累加平方误差
141        return total_loss
142
143
144
145 class StaticKmeansController:
146     """
147     Static KMeans clustering controller using faiss for full-dataset clustering.
148
149     Attributes:
150         data (np.ndarray): Full dataset of shape (N, D), dtype float32.
151         k (int): Number of clusters.
152         niter (int): Number of iterations for faiss KMeans.
153         seed (int): Random seed for reproducibility.
154         centroids (np.ndarray): Final cluster centers of shape (k, D).
155         assignments (dict[int, int]): Mapping from data index to assigned cluster ID.
156         cluster_data_indices (dict[int, list[int]]): Mapping from cluster ID to list
       of data indices.
157         loss (float): Final inertia (sum of squared distances to centroids).
```

```python
158             time_taken (float): Time spent on clustering (seconds).
159         """
160
161     def __init__(self, data: np.ndarray, k: int, niter: int = 100, seed: int = 42):
162         self.data = data.astype('float32')
163         self.k = k
164         self.niter = niter
165         self.seed = seed
166         self.centroids = None
167         self.assignments = {}
168         self.cluster_data_indices = {}
169         self.loss = None
170         self.time_taken = None
171
172     def run(self):
173         """
174         Run static KMeans clustering using faiss.Kmeans.
175         Records centroids, assignments, cluster_data_indices, loss, and time_taken.
176         """
177         try:
178             import faiss
179         except ImportError:
180             raise ImportError("faiss library not installed. Please install faiss to␣
    ↪use StaticKmeansController.")
181
182         # Initialize faiss KMeans
183         d = self.data.shape[1]
184         kmeans = faiss.Kmeans(d, self.k, niter=self.niter, verbose=False, seed=self.
    ↪seed)
185
186         # Train and time
187         start_time = time.time()
188         kmeans.train(self.data)
189         self.time_taken = time.time() - start_time
190
191         # Retrieve centroids
192         self.centroids = kmeans.centroids.copy()
193
194         # Assign each point to the nearest centroid
195         distances, labels = kmeans.index.search(self.data, 1)
```

```python
196         labels = labels.flatten()

197

198         # Build assignments and cluster-to-indices mapping
199         for idx, lbl in enumerate(labels):
200             self.assignments[idx] = int(lbl)
201             self.cluster_data_indices.setdefault(int(lbl), []).append(idx)

202

203         # Compute inertia (loss)
204         total_loss = 0.0
205         for idx, lbl in self.assignments.items():
206             center = self.centroids[lbl]
207             point = self.data[idx]
208             total_loss += np.sum((point - center) ** 2)
209         self.loss = float(total_loss)

210

211     def report(self):
212         """
213         Return a summary of clustering results.
214         """
215         return {
216             'Method': 'Static faiss.Kmeans',
217             'Time (s)': self.time_taken,
218             'Loss': self.loss,
219             'Num Clusters': self.k,
220             'Cluster Sizes': {cid: len(idxs) for cid, idxs in self.cluster_data_
    →indices.items()}
221         }
```

## 1.3 `background.py`

## 1.4 `logger.py`

Listing 3: logger.py

```python
1  import os
2  import logging

3

4  from constants.filepath_constants import RESULTS_DIR

5
```

```python
6   def setup_logging(fileName = 'app.log'):
7       # set up logging file, format and location, information level
8
9       # logFormat = "%(asctime)s - %(levelname)s - %(message)s"
10      logFormat = "%(message)s"
11      dateFormat = "%m/%d/%Y %H:%M:%S %p"
12
13      # fileHandler
14      path = os.path.join(RESULTS_DIR, fileName)
15      file_handler = logging.FileHandler(path, mode='w')
16
17      logging.basicConfig(
18          format=logFormat,
19          datefmt=dateFormat,
20          level=logging.DEBUG,
21          handlers= [
22              file_handler,
23              # stream_handler
24          ]
25      )
26
27      LOGGER = logging.getLogger()
28      return LOGGER
```

## 1.5 `processer.py`

Listing 4: processer.py

```python
1   import logging
2   from optimizer import DynamicKmeansController, StaticKmeansController
3   from util import *
4   LOGGER = logging.getLogger(__name__)
5
6   # 2. 数据处理类
7   class ComparisonProcessor:
8       """
9       对比动态 vs 静态 聚类结果，存储以下信息：
10        - center_distance: 两组簇中心平均最近距离
11        - dynamic_loss, static_loss
```

```python
12          - dynamic_time, static_time
13          - dynamic_history: loss & time 历史
14      """
15      def __init__(self, dynController:DynamicKmeansController,
    →staController:StaticKmeansController):
16          self.dyn = ComparisonProcessor.summarize_dynamic(dynController)
17          self.sta = ComparisonProcessor.summarize_static(staController)
18          self.results = {}
19          self._compute()
20
21      def _compute_center_distance(self):
22          # 对于每个动态中心，找到距离最近的静态中心，计算平均距离
23          dyn_c = self.dyn['centroids']
24          sta_c = self.sta['centroids']
25          dists = []
26          for c in dyn_c:
27              diff = sta_c - c
28              dist_to_sta = np.linalg.norm(diff, axis=1)
29              dists.append(np.min(dist_to_sta))
30          return float(np.mean(dists))
31
32      def _compute(self):
33          self.results['center_distance'] = self._compute_center_distance()
34          self.results['dynamic_loss'] = float(self.dyn['loss'])
35          self.results['static_loss'] = float(self.sta['loss'])
36          self.results['dynamic_time'] = float(self.dyn['time'])
37          self.results['static_time'] = float(self.sta['time'])
38          self.results['dynamic_loss_history'] = self.dyn['loss_history']
39          self.results['dynamic_time_history'] = self.dyn['time_history']
40
41
42      @staticmethod
43      def summarize_dynamic(controller:DynamicKmeansController):
44          """
45          输出动态聚类摘要信息
46          Returns a dict with:
47          - num_clusters: 最终簇数
48          - loss: 最终 loss
49          - time: 累计运行时间
50          - loss_history: 每批 loss 的列表
```

```
51              - time_history: 每批累计时间的列表
52              - centroids: 簇中心数组, shape=(k, D)
53              """
54              centroids = np.stack(controller.centroids, axis=0)
55              return {
56                  'method': 'dynamic',
57                  'num_clusters': centroids.shape[0],
58                  'loss': controller.loss_history[-1],
59                  'time': controller.time_history[-1],
60                  'loss_history': np.array(controller.loss_history),
61                  'time_history': np.array(controller.time_history),
62                  'centroids': centroids
63              }
64
65          @staticmethod
66          def summarize_static(controller:StaticKmeansController):
67              """
68              输出静态聚类摘要信息
69              Returns a dict with:
70              - num_clusters: 簇数 k
71              - loss: 最终 loss
72              - time: 聚类耗时
73              - centroids: 簇中心数组, shape=(k, D)
74              """
75              return {
76                  'method': 'static',
77                  'num_clusters': controller.k,
78                  'loss': controller.loss,
79                  'time': controller.time_taken,
80                  'centroids': controller.centroids
81              }
82
83
84
```

## 1.6 `plotter.py`

Listing 5: plotter.py

```
1   import matplotlib.pyplot as plt
2
3   import logging
4   from processer import ComparisonProcessor
5   from constants.filepath_constants import RESULTS_DIR
6
7   LOGGER = logging.getLogger(__name__)
8
9
10
11
12
13  # 3. 画图类
14  class ComparisonPlotter:
15      """
16      接收 ComparisonProcessor, 绘制:
17        - loss 对比 (bar)
18        - time 对比 (bar)
19        - 中心距离 (单个文本展示)
20        - 动态聚类 loss & time 历史曲线
21      """
22      def __init__(self, comp: ComparisonProcessor):
23          self.comp = comp
24
25      def plotAll(self):
26          self.plot_dynamic_history()
27          self.plot_loss_comparison()
28          self.plot_time_comparison()
29
30      def plot_loss_comparison(self):
31          fig = plt.figure()
32          labels = ['dynamic', 'static']
33          values = [self.comp.results['dynamic_loss'], self.comp.results['static_loss']]
34          plt.bar(labels, values)
35          plt.ylabel('Loss')
36          plt.title('Dynamic vs Static Loss Comparison')
37          plt.show()
```

```python
38
39      def plot_time_comparison(self):
40          fig = plt.figure()
41          labels = ['dynamic', 'static']
42          values = [self.comp.results['dynamic_time'], self.comp.results['static_time']]
43          plt.bar(labels, values)
44          plt.ylabel('Time (s)')
45          plt.title('Dynamic vs Static Time Comparison')
46          plt.show()
47
48      def print_center_distance(self):
49          dist = self.comp.results['center_distance']
50          print(f"Average nearest-center distance between dynamic and static: {dist:.4f}
    ↪")
51
52      def plot_dynamic_history(self):
53          # Loss history
54          fig1 = plt.figure()
55          plt.plot(self.comp.results['dynamic_loss_history'])
56          plt.xlabel('Batch')
57          plt.ylabel('Loss')
58          plt.title('Dynamic Clustering Loss over Batches')
59          plt.show()
60          # Time history
61          fig2 = plt.figure()
62          plt.plot(self.comp.results['dynamic_time_history'])
63          plt.xlabel('Batch')
64          plt.ylabel('Cumulative Time (s)')
65          plt.title('Dynamic Clustering Time over Batches')
66          plt.show()
```

## 1.7 `util.py`

Listing 6: util.py

```python
1   import numpy as np
```