

# 并行计算的K-Means聚类算法实现

## 一，实验介绍

聚类是拥有相同属性的对象或记录的集合，属于无监督学习，K-Means聚类算法是其中较为简单的聚类算法之一，具有易理解，运算深度块的特点.

### 1.1 实验内容

通过本次课程我们将使用C++语言实现一个完整的面向对象的可并行K-Means算法.这里我们一起围绕着算法需求实现各种类，最终打造出一个健壮的程序.所以为了更好地完成这个实验，需要你有C++语言基础，会安装一些常用库，喜欢或愿意学习面向对象的编程思维.

### 1.2 实验知识点

- C++语言语法
- K-Means算法思路与实现
- 并行计算思路与实现
- boost库的常用技巧(Smart Pointers,Variant, tokenizer)

### 1.3 实验环境

- Xfce 终端（Xfce Terminal）：  
Linux 命令行终端，打开后会进入 Bash 环境，可以用来执行 Linux 命令和调用系统调用.
- GVim：非常好用的编辑器，不会使用的可以参考课程《Vim编辑器》.
- boost,MPICH2库

### 1.4 适合人群

本课程适合有C++语言基础，对聚类算法感兴趣并希望在动手能力上得到提升的同学.

### 1.5 代码获取

### 1.6 效果图

完成时间显示:

- 单进程

```
completed in 31.9997 seconds
number of processes: 1
```

- 8进程

```
completed in 7.35373 seconds
number of processes: 8
```

输出结果文件

```
Clustering Summary:
Number of clusters: 3
Size of Cluster 0: 4989
Size of Cluster 1: 5011
Size of Cluster 2: 5000

Number of given clusters: 3
Cross Tabulation:
Cluster ID   1   2   3
0             0   0  4989
1             50000  11
2             0   50000

Cluster Membership
Record ID, Cluster Index, Cluster Index Given
1, 1, 0
2, 1, 0
3, 1, 0
4, 1, 0
```

图1 输出文件图

## 1.7 项目结构与框架

项目的整个文件目录:

```
|— clusters
|   |— distance.hpp
|   |— record.hpp
|— datasets
|   |— attrinfo.hpp
|   |— dataset.hpp
|   |— dcattrinfo.hpp
|— mainalgorithm
|   |— kmean.hpp
|   |— kmeanmain.cpp
|— utilities
|   |— datasetreader.hpp
|   |— exceptions.hpp
|   |— null.hpp
|   |— types.hpp
```

这里简单介绍一下功能模块,在具体实践每一个类的时候会有详细UML图或流程图.

主要分为4个模块:数据集类,聚集类,实用工具类,算法类.

- 实用工具类:定义各种需要的数据类型;常用的异常处理;文件读取.
- 数据集类:将文件中的数据通过智能指针建立一个统一数据类,拥有丰富的属性和操作.
- 聚集类:在数据类基础上实现中心簇.
- 算法类:完成对聚集类的初始化,通过算法进行更新迭代,最终实现数据集的聚类并输出聚类结果.

## 二, 实验原理

这一章我们将配置好我们的实验环境并介绍一些基础知识.

### 2.1 依赖库安装

安装boost和mpich2

```
mpich2下载:
wget -c http://www.mpich.org/static/downloads/3.2.1/mpich-3.2.1.tar.gz

解压:
tar xvfz mpich-3.2.1.tar.gz

配置:
cd mpich-3.2.1
./configure

编译:
make

安装:
make install

boost下载:
wget -c https://dl.bintray.com/boostorg/release/1.68.0/source/boost_1_68_0.tar.gz

解压
tar xvfz boost_1_68_0.tar.gz
cd boost_1_68_0

编译:
sh bootstrap.sh

修改project-config.jam 文件
第19行添加一句:using mpi;

安装:
./bjam --with-programoptions --with-mpi install
```

检验boost是否安装成功,可以检测一下:  
运行源码,test/mpitest.cpp

```
mpic++ -o mpitest mpitest.cpp -L/usr/local/lib -lboost_mpi -lboost_serialization

mpirun -n 3 ./mpitest(3个进程)
```

若结果如下,有三个Process则证明安装成功!

```
Process 1: a msg from master

Process 2: a msg from master

Process 2:
Process 1:
Process 0: zero one two
Process 0: zero one two
Process 1: zero one two
Process 2: zero one two
```

## 2.2 boost的小技巧

### Smart Pointers

在Boost中，智能指针是存储指向动态分配对象的指针的对象.智能指针非常有用，因为它们确保正确销毁动态分配的对象，即使在异常情况下也是如此.事实上，智能指针被视为拥有指向的对象，因此负责在不再需要时删除对象.Boost智能指针库提供了六个智能指针类模板.表给出了这些类模板的描述.本实验中将大量使用智能指针.

类	描述
scoped_ptr	单个对象的简单唯一所有权,不可复制.

scoped_array	数组的简单唯一所有权.不可复制
shared_ptr	对象所有权在多个指针之间共享
shared_array	多个指针共享的数组所有权
weak_ptr	shared_ptr拥有的对象的非拥有观察者
intrusive_ptr	具有嵌入引用计数的对象的共享所有权.

表1 智能指针类型简介

## Variant versus Any

Boost Variant类模板是一个安全通用的联合容器，和std::vector不同储存单个类型的多个值，variant可以储存多个类型的单个值，本实验中将使用variant储存双精度和整数类型来表示不同类型的数据。

与variant一样，Boost any是另一个异构容器.虽然Boost anys有许多与Boost variant相同的功能.根据Boost库文档，Boost variant比Boost any具有以下优势：

- 1，variant保证其内容的类型是用户指定的有限类型集之一.
- 2，variant提供对其内容的编译时检查访问.
- 3，variant通过提供有效的，基于堆栈的存储方案，可以避免动态分配的开销.

同样Boost any也有一些优势：

- 1,any几乎允许任何类型的内容.
- 2,很少使用模板元编程技术.
- 3,any 对交换操作提供安全的不抛出异常保证.

## Tokenizer

Tokenizer提供了一种灵活而简单的方法通过分割符（如:" , "）将一个完整的string分隔开。

字符串为：“ A flexible,easy tokenizer “

如果通过","分割,则结果为：

[A flexible] [ easy tokenizer>]

以" " 为分隔符:

分割结果为：

[A] [flexible,] [easy] [tokenizer]

## 三，实验步骤

接下来将具体实践各个类,会给出每一个类的声明并解释其成员函数和数据成员以及相关类之间的继承关系和逻辑关系.涉及到重要的成员函数的实现会给出其定义代码,一些普通的成员函数的源码可以到下载的源文件中查看,里面也会有详细的注解。

### 3.1 数据集的构建

数据对于一个聚类算法来说非常重要,在这里我们将一个数据集描述为一个记录(record),一个记录由一些属性(Attribute)表征.因此自然而然将依次建立attributes,records,最后是数据集datasets。

在此之前我们需要了解一下我们在聚类中实际接触到的数据类型。  
这里有一个示例,心脏数据集。

```
//heart.data
70.0,1.0,4.0,130.0,322.0,0.0,2.0,109.0,0.0,2.4,2.0,3.0,3.0,2
67.0,0.0,3.0,115.0,564.0,0.0,2.0,160.0,0.0,1.6,2.0,0.0,7.0,1
57.0,1.0,2.0,124.0,261.0,0.0,0.0,141.0,0.0,0.3,1.0,0.0,7.0,2
64.0,1.0,4.0,128.0,263.0,0.0,0.0,105.0,1.0,0.2,2.0,1.0,7.0,1
74.0,0.0,2.0,120.0,269.0,0.0,2.0,121.0,1.0,0.2,1.0,1.0,3.0,1
65.0,1.0,4.0,120.0,177.0,0.0,0.0,140.0,0.0,0.4,1.0,0.0,7.0,1
.....
```

包含13个属性,age,sex,chest pain type(4 values),resting blood pressure.....

为了更好地表述不同数据相同属性的差异,我们需要对这些数据进行离散/连续处理,即对于有些数据我们认为它是连续的如:age,有些是离散的如:年龄.这样我们建立一个描述数据类型的文件:

```
//heart.names
schema file for heart.dat
///: schema
1, Continuous
2, Discrete
3, Discrete
4, Continuous
5, Continuous
6, Discrete
7, Discrete
8, Continuous
9, Discrete
10, Continuous
11, Discrete
12, Continuous
13, Discrete
14, Class
```

### 3.1.1 AttrValue类

AttrValue类有一个私有变量,有两个友元函数,一个公有成员函数.

\_value是一个variant类型变量,它可以存储一个双精度或无符号整形的数据,分类数据用无符号整形数据表示.

AttrValue类自身无法存储或获取数据.它的两个友元函数可以获取和修改数据\_value.

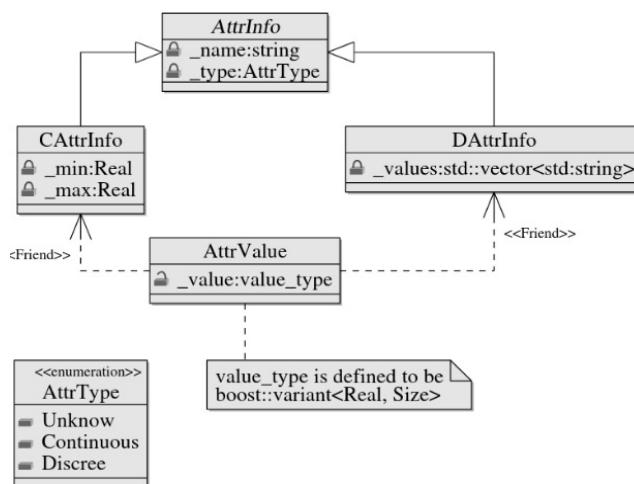


图2 数据类UML关系图

```

1 //source:datasets.attrinfo.hpp
2 class AttrValue
3 {
4     public:
5         friend class DAttrInfo; //友元函数可以访问_value
6         friend class CAttrInfo; //友元函数可以访问_value
7         typedef boost::variant<Real, Size> value_type; //可存储双精度和无符号整形数据
8         AttrValue();
9     private:
10         value_type _value;
11 };
12
13 inline AttrValue::AttrValue(): _value(Null<Size>()) {
14     } //构造函数, 将_value初始化为Null<Size>(定义在utilities/null.hpp中)

```

### 3.1.2 AttrInfo类

AttrInfo是一个基类, 包括了许多虚函数和纯虚函数. 这些函数都将在它的派生类中具体实现, 基类中仅进行声明和简单定义.

```

1 //source:datasets.attrinfo.hpp
2 //三种数据类型: 未知型, 连续型(双精度), 离散型(无符号整形)
3 enum AttrType
4 {
5     Unknow,
6     Continuous,
7     Discrete
8 };
9
10 class DAttrInfo;
11 class CAttrInfo;
12 class AttrInfo
13 {
14     public:
15         AttrInfo(const std::string &name, AttrType type); //每一栏的属性名(id, attr, label, ...)和该属性的数据类型
16         virtual ~AttrInfo() {} //虚析构函数
17         std::string &name(); //返回标签
18         AttrType type() const; //返回数据类型
19         virtual Real distance(const AttrValue&, const AttrValue&) const = 0;
20
21         virtual void set_d_val(AttrValue&, Size) const; //AttrValue赋值; 适用于DAttrInfo
22         virtual Size get_d_val(const AttrValue&) const; //获取_value
23         virtual void set_c_val(AttrValue&, Real) const; //AttrValue赋值; 适用于CAttrInfo
24         virtual Real get_c_val(const AttrValue&) const; //获取_value
25         virtual bool can_cast_to_d() const; //布尔值, 对于DAttrInfo类来说其返回值为true, 相反为false. 在基类的声明中
26         virtual bool can_cast_to_c() const;
27         virtual DAttrInfo& cast_to_d(); //返回DAttrInfo本身
28         virtual bool is_unknown(const AttrValue&) const = 0;
29         virtual void set_unknown(AttrValue&) const = 0;
30     private:
31         std::string _name;
32         AttrType _type;
33 };

```

### 3.1.3 CAttrInfo类和DAttrInfo类

CAttrInfo主要是用来表示连续型数据的一些属性和方法. 有两个数据成员: \_min和\_max. 表示最小值和最大值属性, 在初始化时都将设置为 Null<Size> . 这两个属性将在归一化的时候用到. CAttrInfo将会继承AttrInfo的一些函数, 并且重新定义.

```
1 //source:datasets/dcatrrinfo.hpp
2 class CAttrInfo: public AttrInfo
3 {
4     public:
5         CAttrInfo(const std::string& name); //构造函数
6         Real distance(const AttrValue&, const AttrValue&) const; //两个距离
7         void set_c_val(AttrValue &, Real) const;
8         void set_min(Real); //设置最小值
9         void set_max(Real); //设置最大值
10        Real get_min() const; //获取最小值
11        Real get_max() const; //获取最大值
12        Real get_c_val(const AttrValue&) const;
13        bool is_unknown(const AttrValue&) const;
14        bool can_cast_to_c() const;
15        void set_unknown(AttrValue&) const;
16    protected:
17        Real _min;
18        Real _max;
19 };
20 CAttrInfo::CAttrInfo(const std::string& name)
21     : AttrInfo(name, Continuous) {
22     _min = Null<Real>();
23     _max = Null<Real>();
24 }
```

DAttrInfo类有一个私有变量\_values,它是一个string类型的vector,用来存储一些离散的字符串.在DAttrInfo对象中所有的离散值都将由字符串转化为唯一的无符号整形.

```
1 //source:datasets/dcatrrinfo.hpp
2 class DAttrInfo: public AttrInfo //继承AttrInfo
3 {
4     public:
5         DAttrInfo(const std::string& name); //构造函数, 传入属性字符串
6         const std::string& int_to_str(Size i) const;
7         Size num_values() const; //获取长度
8         Size get_d_val(const AttrValue&) const; //接口定义
9         void set_d_val(AttrValue& , Size) const; //接口定义
10        Size add_value(const std::string&,
11                        bool bAllowDuplicate = true); //将一组离散值加入到_values中, 比如"X,X,Y,Z",
12                                                    //则values=[X,Y,Z], 对应的二进制数字为[0,0,1,2]
13                                                    //对于属性值, 则可以重复, 但对于id则具有唯一性, 不能重复
14        DAttrInfo& cast_to_d();
15        Real distance(const AttrValue&, const AttrValue&) const; //比较两个离散型变量的距离
16        bool is_unknown(const AttrValue& av) const; //值有缺省
17        bool can_cast_to_d() const;
18        void set_unknown(AttrValue&) const;
19    protected:
20        std::vector<std::string> _values;
21 };
```

add\_value 是一个将字符串转化为无符号整形数据的重要函数,返回值为该字符所表示的整形,并将为出现的字符添加进\_values.

Record	Attribute	AttrValue
1	"A"	0
2	"B"	1
3	"A"	0
4	"C"	2
5	"B"	1

Record	Attribute
0	"A"
1	"B"
2	"C"

表2 DAttrInfo的一个具体实例

通过上面表格中我们可以看到一组字符类型的数据被存储为该字符串所在的inex,如果该字符串第一次出现则为上一个字符串的index+1.这样相同的字符串都被转化为唯一的无符号整形.\_value这个辅助变量可以帮助实现这一功能.

```
1 //source:datasets/dcattribinfo.hpp
2 Size DAttrInfo::add_value(const std::string& s,
3     bool bAllowDuplicate) {
4     Size ind = Null<Size>();
5     //如果该字符串已经出现,则返回该字符串在_values中的index
6     for(Size i=0;i<_values.size();++i) {
7         if(_values[i] == s) {
8             ind = i;
9             break;
10        }
11    }
12    //如果未出现,则返回_values的大小-1.
13    //同时对于不允许重复字符串的数据,如ID,当出现重复字符串时则会错误提示.
14    if(ind == Null<Size>()) {
15        _values.push_back(s);
16        return _values.size()-1;
17    } else {
18        if(bAllowDuplicate) {
19            return ind;
20        } else {
21            FAIL("value "<<s<<" already exists");
22            return Null<Size>();
23        }
24    }
25 }
```

这里需要看一下distance这个函数的定义,它返回的是一个双精度类型数值.如果传入的两个数据类型为Unknow则返回为0.0,其中一个为Unknow则为1,对于两个双精度类型的数据返回其差值.

```
1 //source:datasets/dcattribinfo.hpp
2 Real CAttrInfo::distance(const AttrValue& av1,const AttrValue& av2) const {
3     if(is_unknown(av1) && is_unknown(av2)){
4         return 0.0;
5     }
6     if(is_unknown(av1) ^ is_unknown(av2)){
7         return 1.0;
8     }
9     return boost::get<Real>(av1._value) -
10         boost::get<Real>(av2._value);
11 }
```

对于离散型数据,两个离散数据之间的距离定义也会不同,这里主要是考虑到离散型数据都转化为相差为1的整形,所以只要两个DAttrInfo的值不同则距离就为1.0,所以在含有离散型和连续型数据的混合数据中连续型数据要进行归一化处理以满足量纲统一.



```

1 //source:datasets/dattrinfo.hpp
2 Real DAttrInfo::distance(const AttrValue& av1,
3                           const AttrValue& av2) const {
4     if(is_unknown(av1) && is_unknown(av2)) {
5         return 0.0; //如果两个值都有缺省,则距离为0
6     }
7     if(is_unknown(av1) ^ is_unknown(av2)) {
8         return 1.0; //如果有一个值缺省,距离为1
9     }
10    if(boost::get<Size>(av1._value) ==
11        boost::get<Size>(av2._value) ) {
12        return 0.0; //如果两个值相等,则无差距
13    } else {
14        return 1.0; //否则为最大距离1
15    }
16 }

```

### 3.1.4 Container类

Container类是一个基类模板,有一个vector的数据成员\_data.add函数可以将T类型的数据添加进入\_data,同样erase可以删除数据.[]是一个操作符重载,返回索引对应的数据.

```

1 //source:clusters/record.hpp
2 template <typename T>
3 class Container//基类模板
4 {
5     public:
6         typedef typename std::vector<T>::iterator iterator;
7         iterator begin();
8         iterator end();
9         void erase(const T& val);
10        void add(const T&val); //将val添加到向量中
11        Size size() const; //返回_data的长度
12        T& operator[](Size i); //下标索引,建立Schema与data的关系
13    protected:
14        ~Container(){}
15        std::vector<T>_data;
16 };

```

Record和Schema是继承Container类的两个重要的类,他们之间的关系如下:

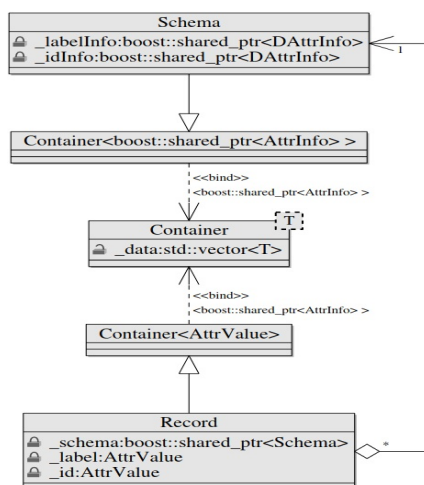


图3 Container关系图

### 3.1.5 Schema类

Schema有两个保护数据成员\_labelInfo,\_idInfo.和一个继承父类的成员\_data,\_data是一个元素为AttrInfo的vector,表示每一个数据的属性(离散/连续).\_labelInfo是一个指向DattrInfo的共享指针,其包含了输入数据的分类情况.

Schema的目的是为一个Record对象设置label和id.set\_id和set\_label函数是为了实现此功能,但是他们又依赖与Record所以

我们在Record类中具体定义.

```
1 //source:clusters/record.hpp
2 class Record;
3 class Schema:public Container<boost::shared_ptr<AttrInfo> >
4 {
5     public:
6         const boost::shared_ptr<DAttrInfo>& labelInfo() const;//标签信息, 整形
7         const boost::shared_ptr<DAttrInfo>& idInfo() const;//id信息, 整形
8         boost::shared_ptr<DAttrInfo>& idInfo();//可以修改成员变量,_labelInfo
9         boost::shared_ptr<DAttrInfo>& labelInfo();//可以修改成员变量,_idInfo
10        void set_label(const boost::shared_ptr<Record>& r,const std::string& val);
11        //设置记录的label
12        void set_id(boost::shared_ptr<Record>& r,const std::string& val);
13        //设置记录的id
14    protected:
15        boost::shared_ptr<DAttrInfo> _labelInfo;
16        boost::shared_ptr<DAttrInfo> _idInfo;
17 };
```

### 3.1.6 Record类

Record继承带参数AttrValue的模板类Container,有四个私有数据成员\_label,\_data,id和\_schema.\_data继承自父类.每一个Record类都有一个指向Schema类的共享指针,可以将类型为AttrValue的数据储存在\_data中,同样每一个record都有一个label和id.Record的构造函数需要传入一个指向Schema的共享指针,并将\_data的长度设置为与\_schema一样,将\_data里的值设置为默认值.我们就可以通过Schema来操控Record,因为Schema的\_data类型为AttrInfo有很多函数如add,set\_c\_val,add\_value等函数可以对离散/类型数据进行操作.所以Record和Schema的关系为通过Schema定义了每一条数据的规范(label,id,每一条属性的类型),然后按照这个规范将数据填充到record中,因为record直接接触的类型是AttrValue.

```
1 //source:clusters/record.hpp
2 class Record:public Container<AttrValue>
3 {
4     public:
5         Record(const boost::shared_ptr<Schema>& schema);//构造函数
6         const boost::shared_ptr<Schema>& schema() const;
7         const AttrValue& labelValue() const;
8         const AttrValue& idValue() const;
9         AttrValue& labelValue();
10        AttrValue& idValue();
11        Size get_id() const;
12        Size get_label() const;
13    private:
14        boost::shared_ptr<Schema> _schema;//通过_schema创建记录
15        AttrValue _label;
16        AttrValue _id;
17 };
```

### 3.1.7 dataset类

上面已经实现了一条数据的储存就是一个Record,我们最终需要n条数据.这里新定义一个类Dataset.很明显按照上面的思路,Record依赖Schema,则Dataset依赖Record.

所以Dataset类继承类型为Record的Container.因为最后我们使用的Dataset类,我们一些我们需要用到的属性可以在这里直接给出.num\_attr(),返回属性的个数,is\_numeric()判断该列属性值是否是连续行(对于Kmeans算法这里需要连续型数据),为了更加方便第获取每一个数据,使用操作符重载.

```
1 //source:datasets/dataset.hpp
2 inline const AttrValue& Dataset::operator()(Size i, Size j) const {
3     return (*_data[i])[j];
4 }
```

```

1 //source:datasets/dataset.hpp
2 class Dataset:public Container<boost::shared_ptr<Record> >
3 {
4     public:
5         Dataset(const boost::shared_ptr<Schema>&); //构造函数，传入含有属性值的schema
6         Size num_attr() const; //返回属性个数
7         const boost::shared_ptr<Schema> &schema() const; //返回_schema
8         const AttrValue& operator()(Size i, Size j) const; //返回第i条第j个属性的值
9         std::vector<Size> get_CM() const;
10        bool is_numeric() const;
11        bool is_categorical() const;
12    protected:
13        boost::shared_ptr<Schema> _schema;
14 };

```

## 3.2 创建一个数据实例

前面关于如何构建dataset相关类已经花了很多时间,下面就让我们实际操作如何创建一个具体的dataset.

假设我们有这样的一组数据:

ID	Attr1	Attr2	Attr3	Label
r1	1.2	A	-0.5	1
r2	-2.1	B	1.5	2
r3	1.5	A	-0.1	1

表3 数据实例

那么我们如何将以上数据用我们的dataset类来表示呢?

```

1 //test/datasettest.cpp
2 #include "../clusters/record.hpp"
3 #include "../datasets/dataset.hpp"
4 #include <iostream>
5 #include <sstream>
6 #include <iomanip>
7 using namespace std;
8 int main()
9 {
10     boost::shared_ptr<Schema> schema(new Schema);
11     boost::shared_ptr<DAttrInfo> labelInfo(new DAttrInfo("Label"));
12     boost::shared_ptr<DAttrInfo> idInfo(new DAttrInfo("id"));
13     schema->labelInfo() = labelInfo;
14     schema->idInfo() = idInfo;
15
16     stringstream ss;
17     boost::shared_ptr<AttrInfo> ai;
18     for(Size j=0;j<3;++j)
19     {
20         ss.str("");
21         ss<<"Attr"<<j+1;
22         if(j==0||j==2)
23         {
24             ai = boost::shared_ptr<CAAttrInfo>(new CAAttrInfo(ss.str()));
25         }
26         else{
27             ai = boost::shared_ptr<DAttrInfo>(new DAttrInfo(ss.str()));
28         }
29         schema->add(ai);
30     }
31     boost::shared_ptr<Dataset> ds(new Dataset(schema));
32     Size val;

```

```

32 boost::shared_ptr<Record> r;
33
34 r = boost::shared_ptr<Record>(new Record(schema));
35 schema->set_id(r, "r1");
36 schema->set_label(r, "1");
37 (*schema)[0]->set_c_val((*r)[0], 1.2);
38 val = (*schema)[1]->cast_to_d().add_value("A");
39 (*schema)[1]->set_d_val((*r)[1], val);
40 (*schema)[2]->set_c_val((*r)[2], -0.5);
41 ds->add(r);
42
43 r = boost::shared_ptr<Record>(new Record(schema));
44 schema->set_id(r, "r2");
45 schema->set_label(r, "2");
46 (*schema)[0]->set_c_val((*r)[0], -2.1);
47 val = (*schema)[1]->cast_to_d().add_value("B");
48 (*schema)[1]->set_d_val((*r)[1], val);
49 (*schema)[2]->set_c_val((*r)[2], 1.5);
50 ds->add(r);
51
52 r = boost::shared_ptr<Record>(new Record(schema));
53 schema->set_id(r, "r3");
54 schema->set_label(r, "1");
55 (*schema)[0]->set_c_val((*r)[0], 1.5);
56 val = (*schema)[1]->cast_to_d().add_value("A");
57 (*schema)[1]->set_d_val((*r)[1], val);
58 (*schema)[2]->set_c_val((*r)[2], -0.1);
59 ds->add(r);
60
61 cout<<"Data: \n";
62 cout<<setw(10)<<left<<"RecordID";
63 for(Size j=0; j<ds->num_attr(); ++j) {
64     stringstream ss;
65     ss<<"Attr(" <<j+1<<")";
66     cout<<setw(10)<<left<<ss.str();
67 }
68 cout<<setw(6)<<left<<"Label"<<endl;
69 for(Size i=0; i<ds->size(); ++i) {
70     cout<<setw(10)<<left<<(*ds)[i]->get_id();
71     for(Size j=0; j<ds->num_attr(); ++j) {
72         if((*schema)[j]->can_cast_to_c()) {
73             cout<<setw(10)<<left
74                 <<(*schema)[j]->get_c_val((*ds)(i,j));
75         } else {
76             cout<<setw(10)<<left
77                 <<(*schema)[j]->get_d_val((*ds)(i,j));
78         }
79     }
80     cout<<setw(6)<<left<<(*ds)[i]->get_label()<<endl;
81 }
82 return 0;
83 }
84

```

输出结果与我们预想的一样:

```

Data:
RecordID  Attr(1)  Attr(2)  Attr(3)  Label
0          1.2      0       -0.5     0
1         -2.1      1        1.5     1
2          1.5      0       -0.1     0

```

### 3.3 构建簇

构建簇的目的就是为了将dataset中的record进行重新组合,所以我们定义一个基类Cluster直接接触Record,有一个数据成员id.

```

1  class Cluster:public Container<boost::shared_ptr<Record> >
2  {
3      public:
4          virtual ~Cluster() {}
5
6          void set_id(Size id);
7          Size get_id() const;
8      protected:
9          Size _id;
10 };
11 inline void Cluster::set_id(Size id) {
12     _id = id;
13 }
14 inline Size Cluster::get_id() const {
15     return _id;
16 }

```

定义一个中心簇,来表示一个簇的中心.中心簇只有一个数据成员\_center即表示中心簇的指向Record的共享指针.

```

1  //clusters/record.hpp
2  class CenterCluster : public Cluster
3  {
4      public:
5          CenterCluster(){}
6          CenterCluster(const boost::shared_ptr<Record>& center); //构造函数传入一个record
7          const boost::shared_ptr<Record>& center() const; //返回中心点的record,不可更改
8      protected:
9          boost::shared_ptr<Record> _center; //成员变量,中心点的record
10 };
11 CenterCluster::CenterCluster(const boost::shared_ptr<Record>& center):_center(center){}
12 const boost::shared_ptr<Record>& CenterCluster::center()
13     const {
14         return _center;
15     }

```

为了实现更多丰富的功能,我们需要再定义一个类Pclustering.

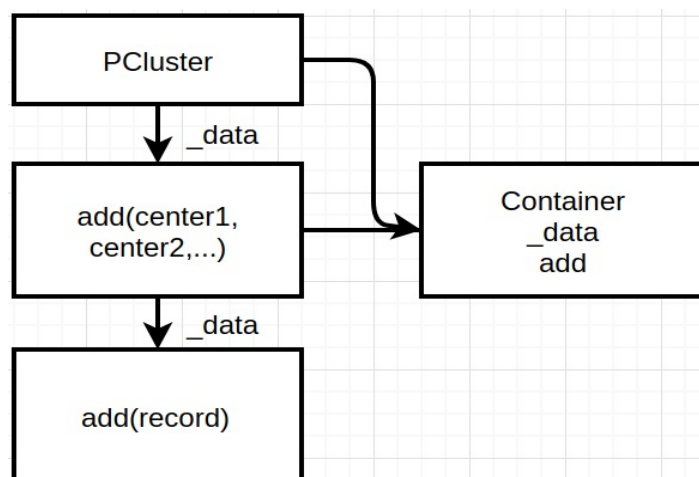


图4 PClustering 关系图

PClustering继承Container,通过add函数添加了中心簇Center.Center也拥有add函数,它添加属于和他同一簇的record,每一个record都有自己的id信息.这样我们就能通过PClustering储存了聚类信息.PClustering的一个数据成员为\_CM,是用来储存每一条record的所属聚类.如:[1,1,1,2,2,2],同一簇拥有相同的数值.calculate函数是用来从\_data中提取相关聚类信息,然后更新\_CM.

```

1 //clusters/record.hpp
2 class PClustering:public Container<boost::shared_ptr<Cluster> >
3 {
4     public:
5         PClustering(); //构造函数
6         friend std::ostream& operator<<(std::ostream& os,
7             PClustering& pc); //操作符重载, 输出聚类结构相关信息
8         void removeEmptyClusters(); //移除空的record
9         void createClusterID(); //创建聚类id
10        void save(const std::string& filename); //保存聚类结果相关信息至文件
11    private:
12        void print(std::ostream &os); //打印聚类结果相关信息
13        void calculate(); //更新_CM和_CMGiven
14        void crosstab(); //将一些聚类结果储存为交叉表
15        bool _bCalculated; //如果数据文件没有标签信息, 则不需要计算_numclustGiven
16        Size _numclust; //聚类数
17        Size _numclustGiven; //文件提供的label数
18        std::vector<Size> _clustsize; //记录每一簇的数据量
19        std::vector<std::string> _clustLabel; //记录原文件中的每个分类的数量
20        std::vector<Size> _CM; //每一条记录数据的所属index
21        std::vector<Size> _CMGiven; //原文件每一条记录所属标签
22        iiiMapB _crosstab; //交叉表储存数据
23 };

```

这里我们介绍一个模板键-值映射类nnmap(utilities/nnmap.hpp),在这里我们用来储存聚类 and 原标签的数量信息. 如有6条数据,计算的\_CM为[1,1,2,2,2,3],所给标签为[0,0,1,1,2,2]. 我们需要通过下面\_crosstab.填充下面的表格

Cluster ID	1	2	3
0	#	#	#
1	#	#	#
2	#	#	#

\_crosstab(1,0)表示聚类为1,标签为0的数量.通过下面的函数,可以为2.同理\_crosstab(2,0)=0,\_crosstab(3,0) = 0.最终可以打印交叉表:

Cluster ID	1	2	3
0	2	0	0
1	0	2	0
2	0	1	1

如果以标签信息为准的化,则(2,2)那个信息有误,每一行只能有一个数据占据,且不能与之前有相同的列.

```

1 //clusters/record.hpp
2 void PClustering::crosstab() {
3     Size c1, c2;
4     for(Size i=0; i<_CM.size();++i) {
5         c1 = _CM[i];
6         c2 = _CMGiven[i];
7         if (_crosstab.contain_key(c1,c2)) {
8             _crosstab(c1,c2) += 1;
9         } else {
10            _crosstab.add_item(c1,c2,1);
11        }
12    }
13 }

```

## 3.4 K-Means算法

### 3.4.1 算法思路

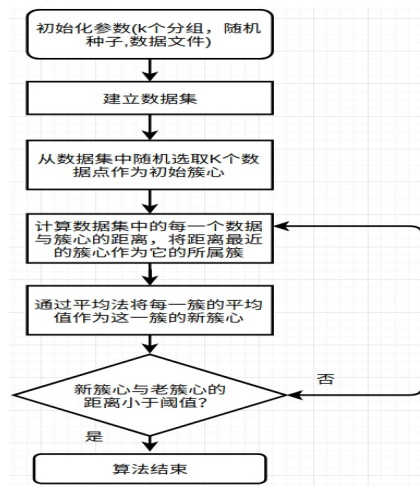


图5 K-Means算法流程图

### 3.4.2 并行化思路

我们使用一种序列 - 均值算法的思路.即计算所有记录 $n$ 和所有中心之间的距离. $p$ 个进程,让每一个参与计算的进程处理  $n/p$ 条数据.主要步骤如下:

- (a)主进程:读取数据文件,并将数据块发送至每一个进程.
- (b)主进程:初始化簇中心,并将这些簇中心发送至每一个程.
- (c)所有进程:计算所给数据块与簇中心的距离,并将这些数据块归属到与它距离最近的中心.
- (d)所有进程:更新新的簇中心.
- (e)所有进程:重复(c)和(4)直至满足停止条件.
- (f)主进程:收集聚类结果.

reduce 是将其他进程汇聚到一个进程.

all\_reduce是将一个进程广播到所有进程.

### 3.4.3 MPIKmean类

将所有的中心簇的数据编码成一个向量\_clusters,这样可以很方便第从一个进程发送至其他进程.同样\_data表示所有的数据的值.

```

1 //source:mainalgorithm/mpikmean.hpp
2 class MPIKmean
3 {
4     public:
5         Arguments& getArguments();//获取初始参数
6         const Results& getResults() const;//获取结果_CM
7         void reset() const;//清除结果
8         void clusterize();//执行计算(初始化,更新,迭代,...)
9     protected:
10         void setupArguments();//设置初始参数
11         void fetchResults() const;//获取结果
12         virtual void initialization() const;//随机初始中心簇
13         virtual void iteration() const;//迭代更新
14         virtual Real dist(Size i, Size j) const;//返回与中心簇的距离
15         mutable vector<Real> _centers;//中心簇的属性值
16         mutable vector<Real> _data;//数据值
17         mutable Size _numObj;//分发给每一个进程的数据量
18         mutable Size _numAttr;//数据属性量
19         mutable vector<Size> _CM;//数据的所属簇index
20
21         mutable vector<boost::shared_ptr<CenterCluster> >
22             _clusters;//中心簇
23         mutable Real _error;//簇之间的总距离
24         mutable Size _numiter;//迭代次数
25         mutable Results _results;//结果
26         boost::shared_ptr<Dataset> _ds;//dataset
27         Arguments _arguments;
28         Size _numclust;//聚类数目
29         Size _maxiter;//最大迭代数目
30         Size _seed;//种子
31         boost::mpi::communicator _world;//mpi通信
32 };

```

主进程负责初始化中心簇(4-33行),一旦中心簇被初始化,就会将中心簇(\_centers)和每个进程的数据的数目(\_numRecords)和属性数(\_numAttr)发送给所有进程(34-36行).

一旦这些数据被进程接收到,每个进程就会划分自己的数据块数量和剩余量(37-38行).首先主进程会将第一个数据块分配给自己(40-49行),剩余的数据通过 `send` 函数发送给其他进程(51-63行).其他进程通过 `recv` 进行接收数据(67行).

```

1 void MPIKmean::initialization() const {
2     Size numRecords;
3     Size rank = _world.rank();
4     if (rank == 0) {
5         numRecords = _ds->size();
6         _numAttr = _ds->num_attr();
7         _centers.resize(_numclust * _numAttr);
8         vector<Integer> index(numRecords,0);
9         for(Size i=0;i<index.size();++i){
10             index[i] = i;
11         }
12         boost::shared_ptr<Schema> schema = _ds->schema();
13         boost::mtrand generator(_seed);
14         for(Size i=0;i<_numclust;++i){
15             boost::uniform_int<> uni_dist(0, numRecords-i-1);
16             boost::variate_generator<boost::mtrand&,
17                 boost::uniform_int<> >
18                 uni(generator, uni_dist);
19             Integer r = uni();
20             boost::shared_ptr<Record> cr = boost::shared_ptr
21                 <Record>(new Record>(*(_ds)[r]));
22             boost::shared_ptr<CenterCluster> c =
23                 boost::shared_ptr<CenterCluster>(
24                     new CenterCluster(cr));
25             c->set_id(i);
26             _clusters.push_back(c);
27             for(Size j=0; j<_numAttr; ++j) {
28                 _centers[i*_numAttr + j] =

```



```

29         (*schema)[j]->get_c_val((*_ds)(r,j));
30     }
31     index.erase(index.begin()+r);
32 }
33 }
34 boost::mpi::broadcast(_world, _centers, 0);
35 boost::mpi::broadcast(_world, numRecords, 0);
36 boost::mpi::broadcast(_world, _numAttr, 0);
37 Size nDiv = numRecords / _world.size();
38 Size nRem = numRecords % _world.size();
39 if(rank == 0) {
40     boost::shared_ptr<Schema> schema = _ds->schema();
41     _numObj = (nRem > 0) ? nDiv+1 : nDiv;
42     _data.resize(_numObj * _numAttr);
43     _CM.resize(_numObj);
44     for(Size i=0; i<_numObj; ++i) {
45         for(Size j=0; j<_numAttr; ++j) {
46             _data[i*_numAttr + j] =
47                 (*schema)[j]->get_c_val((*_ds)(i, j));
48         }
49     }
50     Size nCount = _numObj;
51     for(Size p=1; p<_world.size(); ++p) {
52         Size s = (p< nRem) ? nDiv + 1 : nDiv;
53         vector<Real> dv(s*_numAttr);
54         for(Size i=0; i<s; ++i) {
55             for(Size j=0; j<_numAttr; ++j) {
56                 dv[i*_numAttr+j] =
57                     (*schema)[j]->get_c_val(
58                         (*_ds)(i+nCount,j));
59             }
60         }
61         nCount += s;
62         _world.send(p, 0, dv);
63     }
64 } else {
65     _numObj = (rank < nRem) ? nDiv+1 : nDiv;
66     _CM.resize(_numObj);
67     _world.recv(0,0,_data);
68 }
69 }

```

进行初始化之后,就开始迭代中心簇.

首先定义一个单元素的 vector 来控制循环(2行).在 while 循环内,定义三个局部变量 nChangedLocal,newCenters,newSize . 每一个进程将会处理自己的数据块与每一个中心簇的距离(11-30行).变量 newCenters 包含了一个聚类中所有数据的和. newSize 包含了一个聚类中的数据的数量.一旦所有的数据通过并行处理完毕. all\_reduce 方法将会对所有的进程的数据进行收集,如

```
all_reduce(_world, nChangedLocal, nChanged,vplus<Size>());
```

对所有进程中的nChangedLocal进行相加(通过操作符vplus,具体见源文件定义),只有有一个进程的nChangedLocal>0(中心簇未收敛)则nChange都会>0,整个迭代都会继续进行.(31-36行),在对这些数据收集之后会更新\_center(37-41行).

收敛之后,所有进程会将聚类的index\_CM发送给主进程.主进程会将自己的\_CM添加进去就形成了整个数据集的\_CM(47-58行).

```

1 void MPIKmean::iteration() const {
2     vector<Size> nChanged(1,1); //初始化nChanged, 表示中心簇是否有变化.
3     _numiter = 1; //初始迭代次数
4     while(nChanged[0] > 0) {
5         nChanged[0] = 0;
6         Size s;
7         Real dMin,dDist;
8         vector<Size> nChangedLocal(1,0);
9         vector<Real> newCenters(_numclust*_numAttr,0.0);
10        vector<Size> newSize(_numclust,0);
11        for(Size i=0;i<_numObj;++i) {
12            dMin = MAX_REAL;
13            for(Size k=0;k<_numclust;++k) {
14                dDist = dist(i, k);
15                if (dMin > dDist) {
16                    dMin = dDist;
17                    s = k;
18                }
19            }
20            for(Size j=0; j<_numAttr; ++j) {
21                newCenters[s*_numAttr+j] +=
22                    _data[i*_numAttr+j];
23            }
24            newSize[s] +=1;
25
26            if (_CM[i] != s){
27                _CM[i] = s;
28                nChangedLocal[0]++;
29            }
30        }
31        all_reduce(_world, nChangedLocal, nChanged,
32                    vplus<Size>());
33        all_reduce(_world, newCenters, _centers,
34                    vplus<Real>());
35        vector<Size> totalSize(_numclust,0);
36        all_reduce(_world, newSize, totalSize, vplus<Size>());
37        for(Size k=0; k<_numclust; ++k) {
38            for(Size j=0; j<_numAttr; ++j) {
39                _centers[k*_numAttr+j] /= totalSize[k];
40            }
41        }
42        ++_numiter;
43        if (_numiter > _maxiter){
44            break;
45        }
46    }
47    if(_world.rank() > 0) {
48        _world.send(0,0,_CM);
49    } else {
50        for(Size p=1; p<_world.size(); ++p) {
51            vector<Size> msg;
52            _world.recv(p,0,msg);
53            for(Size j=0; j<msg.size(); ++j) {
54                _CM.push_back(msg[j]);
55            }
56        }
57    }
58 }

```

其他几个函数的定义就不再赘述,相信通过看源文件一定可以看懂.

### 3.4.5 主函数

从前面建立数据集,到构建簇类,编写一些辅助类到算法的应用.最后我们需要用一个实际的文件进行聚类.代码如下:

```

1 //source:mainalgorithm/mpikmeanmain.cpp

```

```

2  #include<boost/timer.hpp>
3  #include<boost/mpi.hpp>
4  #include<boost/program_options.hpp>
5  #include<iostream>
6  #include<sstream>
7  #include<iomanip>
8  #include<functional>
9  #include "mpikmean.hpp"
10 #include "../utilities/datasetreader.hpp"
11
12 using namespace std;
13 using namespace boost::program_options;
14 namespace mpi=boost::mpi;
15 int main(int ac, char* av[]){
16     try{
17         mpi::environment env(ac, av);
18         mpi::communicator world;
19         options_description desc("Allowed options");
20         desc.add_options()
21             ("help", "produce help message")
22             ("datafile", value<string>(), "the data file")
23             ("k", value<Size>()->default_value(3),
24              "number of clusters")
25             ("seed", value<Size>()->default_value(1),
26              "seed used to choose random initial centers")
27             ("maxiter", value<Size>()->default_value(100),
28              "maximum number of iterations")
29             ("numrun", value<Size>()->default_value(1),
30              "number of runs");
31         variables_map vm;
32         store(parse_command_line(ac, av, desc), vm);
33         notify(vm);
34         if (vm.count("help") || ac==1) {
35             cout << desc << "\n";
36             return 1;
37         }
38         Size numclust = vm["k"].as<Size>();
39         Size maxiter = vm["maxiter"].as<Size>();
40         Size numrun = vm["numrun"].as<Size>();
41         Size seed = vm["seed"].as<Size>();
42         string datafile;
43         if (vm.count("datafile")) {
44             datafile = vm["datafile"].as<string>();
45         } else {
46             cout << "Please provide a data file\n";
47             return 1;
48         }
49         boost::shared_ptr<Dataset> ds;
50         if (world.rank() ==0) {
51             DatasetReader reader(datafile);
52             reader.fill(ds);
53         }
54         boost::timer t;
55         t.restart();
56         Results Res;
57         Real avgiter = 0.0;
58         Real avgerror = 0.0;
59         Real dMin = MAX_REAL;
60         Real error;
61         for(Size i=1; i<=numrun; ++i) {
62             MPIKmean ca;
63             Arguments &Arg = ca.getArguments();
64             Arg.ds = ds;
65             Arg.insert("numclust", numclust);
66             Arg.insert("maxiter", maxiter);
67             Arg.insert("seed", seed);
68             if (numrun == 1) {
69                 Arg.additional["seed"] = seed;

```

```

70         } else {
71             Arg.additional["seed"] = i;
72         }
73         ca.clusterize();
74         if(world.rank() == 0) {
75             const Results &tmp = ca.getResults();
76             avgiter +=
77                 boost::any_cast<Size>(tmp.get("numiter"));
78             error = boost::any_cast<Real>(tmp.get("error"));
79             avgererror += error;
80             if (error < dMin) {
81                 dMin = error;
82                 Res = tmp;
83             }
84         }
85     }
86     double seconds = t.elapsed();
87     if(world.rank() == 0) {
88         avgiter /= numrun;
89         avgererror /= numrun;
90         std::cout<<"completed in "<<seconds
91             <<" seconds"<<std::endl;
92         std::cout<<"number of processes: "
93             <<world.size()<<std::endl;
94         PClustering pc =
95             boost::any_cast<PClustering>(Res.get("pc"));
96         std::cout<<pc<<std::endl;
97         std::cout<<"Number of runs: "<<numrun<<std::endl;
98         std::cout<<"Average number of iterations: "
99             <<avgiter<<std::endl;
100        std::cout<<"Average error: "<<avgererror<<std::endl;
101        std::cout<<"Best error: "<<dMin<<std::endl;
102        std::string prefix;
103        size_t ind = datafile.find_last_of('.');
104        if(ind != std::string::npos) {
105            prefix = datafile.substr(0,ind);
106        } else {
107            prefix = datafile;
108        }
109        std::stringstream ss;
110        ss<<prefix<<"-kmean-k"<<numclust<<"-s"<<seed<<".txt";
111        pc.save(ss.str());
112    }
113    return 0;
114 } catch (std::exception& e) {
115     std::cout<<e.what()<<std::endl;
116     return 1;
117 } catch (...) {
118     std::cout<<"unknown error"<<std::endl;
119     return 2;
120 }
121 }

```

编译:

```
mpic++ -o mpikmean mpikmeanmain.cpp -L/usr/local/lib -lboost_program_options -lboost_mpi -lboost_seriali
```

运行:

```
mpirun -n 8 ./mpikmean --datafile=../testdata/15000points.csv --k=10 --numrun=50
```

运行结果如第一章所示.

可以比较使用不同的进程数目的不同运行时间,使用多进程确实可以提高运行速度,但是因为I/O操作会占用一些时间,运行效率

并没有出现倍数的提升.

对于小数据集,I/O操作的开销与数据计算开销相差无几,多进程没有明显优势.对于大数据集,I/O操作开销会小于数据计算的时间,这时候多进程会带来效率上的提升.

## 四，实验总结

到此，我们的K-Means算法的实验就到此结束了.由于考虑到整个内容的繁杂度,有很多小的细节可能没有拿出来细讲,如果小伙伴对有些地方没有弄懂,希望自己能够继续从源码中寻找答案.虽然我们最后只实现了一个简单的聚类算法,但前面介绍的关于构建聚类数据集却具有一定的通用性,对于其他聚类算法也很适用,如果小伙伴愿意尝试其他聚类算法,也可以按照此思路进行改写.并行处理是一种技巧,如果使用恰当,能够给计算效率带来很大的提升,本例的并行处理思路同样可以推广到其他算法当中.

感谢你能够看到最后,希望你有所收获!