

СОФИЙСКИ УНИВЕРСИТЕТ "СВ. КЛИМЕНТ ОХРИДСКИ"  
ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

## Домашна работа 2

ПО

ИЗКУСТВЕН ИНТЕЛЕКТ

СПЕЦ. ИНФОРМАТИКА, 3 КУРС, ЛЕТЕН СЕМЕСТЪР,

УЧЕБНА ГОДИНА 2018/19

ТЕМА: КЛЪСТЕРИЗАЦИЯ С ПОМОЩТА НА МЕТОДА К-MEANS

2 юни 2019 г.

София

*Изготвил:*

Иво Алексеев Стратев

Фак. номер: 45342

Група: 3

# Съдържание

<b>1</b>	<b>Задача</b>	<b>2</b>
<b>2</b>	<b>Описание на използвания метод за решаване на задачата</b>	<b>2</b>
2.1	Класическа версия на алгоритъма k-Means . . . . .	2
2.1.1	Инициализираща фаза . . . . .	2
2.1.2	Стъпка . . . . .	2
2.2	Инициализираща фаза . . . . .	3
2.3	Стъпка . . . . .	3
2.4	Финална фаза . . . . .	4
<b>3</b>	<b>описание на реализацията с псевдокод</b>	<b>4</b>
<b>4</b>	<b>Инструкции за компилиране на програмата</b>	<b>9</b>
<b>5</b>	<b>Примерни резултати</b>	<b>9</b>
5.1	Изходни резултати: . . . . .	10

# 1 Задача

Входните данни са добре известното множество "Iris". Първоначално това множество е публикувано в UCI Machine Learning Repository: Iris Data Set. Всеки ред на таблицата описва цвете ирис, представено с размерите на неговите ботанически части в сантиметри. Таблицата включва описания на 150 примера в термините на 4 атрибута. Представящи ботаническите параметри на цветето. (<http://archive.ics.uci.edu/ml/datasets/Iris>).

Задачата е да се имплементира метода k-Means за групиране на примерите от това множество в три клъстера. Използвайки атрибутите: sepal length, sepal width, petal length и petal width. Множеството съдържа и класови атрибут, които групира примерите в три класа: Iris Setosa, Iris Versicolour и Iris Virginica. След, което да сравнят примерите във всеки от получените клъстери с вече известните класове. Да се използват поне две различни функции за определяне на разстоянието и да определете коя е по-точната за това множество от данни.

## 2 Описание на използвания метод за решаване на задачата

Имплементацията на клъстеризация алгоритъм представлява оптимизирана версия на алгоритъма k-Means. Оптимизациите са две:

- използва се детерминистична инициализираща фаза, вместо типичните рандомизирани инициализации с цел подобряване на сходимостта на алгоритъма.
- стъпката е паралелна, като точките не се разделят, смятат се само новите центроиди. Точките се разделят само в завършителната фаза на реализацията.

### 2.1 Класическа версия на алгоритъма k-Means

В алгоритъма  $k$  е броят на клъстерите, на които искаме да разделим входните данни. Центроид наричаме центъра на всеки клъстер. Той може да бъде намерен чрез използването на различни статистически или псевдостатистически функции. За алгоритъма това е средната точка на клъстера. Тоест математическата функция е *mean*. Понеже алгоритъма е многомерен, то функцията е покомпонентна. Тоест стойността във всяка компонента е средната стойност на съответната компонента на точките от клъстера.

$$Mean(C) = \frac{1}{|C|} \sum_{p \in C} p$$

#### 2.1.1 Инициализираща фаза

Цели да инициализира центроидите. Обикновено се ползват рандомизирани стратегии. Две от най-известните стратегии са:

- На произволен принцип се избират  $k$  точки от множеството.
- На всяка точка по произволен начин се съпоставя клъстер, след което се смятат центроидите на всеки от получените клъстери.

#### 2.1.2 Стъпка

Докато не настъпи определено събитие се повтаря следната логика: Взимат се  $k$  празни множества, представляващи всеки клъстер. Всяко множество съответства на точно един центроид. За всяка точка се пресмята разстоянието до всеки центроид. Точката се премества в съответния клъстер (множество). Пресмятат се медианите на всеки клъстер (множество). Това са новите центроиди. Събитията обикновено са:

- достигнат е определен брой итерации
- Във всеки клъстер по-малко от определен брой точки са били преместени
- Центроидите са се преместили на по-малко от определено разстояние

Следва подобно описание на реализацията.

## 2.2 Инициализираща фаза

Тъй като рандомизираните методи не гарантират по никакъв начин сходимостта на алгоритъма. То е използван следния подход за инициализация, който е със сложност  $\Theta(n)$ , където  $n$  е броя на точките. Тоест сложността е същата като съпоставянето на произволен начин на клъстер. Вземат се първите  $k$  точки от множеството. Това са началните центроиди. След, което за всяка точка се пресмята разстоянието до всеки клъстер. Ако разстоянието между точката и някой клъстер е по-малко от разстоянието от някой центроид до друг, то се обновява някой центроид на най-късо разстояние до точката по формулата:

$$\left( \frac{|C_i|}{|C_i| + 1} \cdot \text{centroid}(C_i) \right) + \left( \frac{1}{|C_i| + 1} \cdot \text{point} \right) \quad (1)$$

където  $C_i$  е избрания клъстер. Ако пък разстоянието от точката до всеки клъстер е по-голямо от разстоянията от всеки клъстер до всеки друг, то се избират два клъстера с минимално разстояние между тях, които да бъдат обединени. Като в резултат центроидите на двата клъстера се заменят с нов центроид пресмет по формулата:

$$\left( \frac{|C_i|}{|C_i| + |C_j|} \cdot \text{centroid}(C_i) \right) + \left( \frac{|C_j|}{|C_i| + |C_j|} \cdot \text{centroid}(C_j) \right) \quad (2)$$

и се добавя нов клъстер с добавянето на нов центроид, който е разглежданата точка.

## 2.3 Стъпка

Лесно се забелязва, че в стъпката на оригиналния вариант на k-Means точките биват разделяни единствено за да могат да бъдат пресметнати новите центроиди. Но това е излишно понеже новите центроиди могат да бъдат пресметнати в движение. При това пресмятането на новите центроиди може да бъде извършено паралелно. За това стъпката в предложената реализация е следната: Множеството от точки се разделя на  $H$  части, за всяка част паралелно се пресмятат временните нови центроиди за съответната част. Тоест за всяка точка във фиксирана част се извършват следните стъпки. Намира се клъстера с най-късо разстояние на центроида си до точката, с най-малък брой точки в себе си. След което ако клъстера е празен, то временния нов центроид е точката. В противен случай се използва формула 1 за обновяване на временния центроид. След като за всяка част бъдат намерени временните нови центроиди се обединяват по формулата:

$$\left( \frac{|temp_i|}{|current_i| + |temp_i|} \cdot \text{centroid}(temp_i) \right) + \left( \frac{|current_i|}{|current_i| + |temp_i|} \cdot \text{centroid}(current_i) \right) \quad (3)$$

Където  $temp_i$  съответства временния клъстер резултат от обединението на центроидите или празното множество.  $current_i$  е клъстера съответстващ на  $i$ -тия центроид от някоя част. Пълната формула, по която се пресмятат новия  $i$ -ти центроид съответстващ на текущия  $i$ -ти центроид е:

$$\sum_{j=1}^H \frac{|temp_{ji}|}{n} \cdot \text{centroid}(temp_{ji}) \quad (4)$$

Където  $temp_{ji}$  съответства на получения центроид съответстващ на  $i$ -тия досегашен центроид от  $j$ -тата част.

## 2.4 Финална фаза

Точките се разпределят в  $k$  клъстера спрямо това до кой центроид се намират най-близо. В реализацията клъстер представлява масив от индекси. Всеки индекс е индексът на съответната точка от входния масив с точки.

## 3 описание на реализацията с псевдокод

Ще използваме следните константи:

- $K$  - брой клъстери (центроиди)
- $N$  - размерност на точките (брой координати)
- $I$  - максимален брой итерации
- $H$  - брой нишки (части, на които да се разделят входните данни)
- $\epsilon$  - достатъчно малко число избрано от потребителя

---

```
1: function K_MEANS(points)
2:   means  $\leftarrow$  INITIALMEANS(points)
3:   for  $i \leftarrow 1$  to  $I$  do
4:     newMeans  $\leftarrow$  CENTROIDS(points, means)
5:     if ARECLOSE(newMeans, means) then break
6:     end if
7:     means  $\leftarrow$  newMeans
8:   end for
9:   return PARTITION(points, means)
10: end function
```

---

---

```
1: function INITIALMEANS(points)
2:   means  $\leftarrow$  FIRSTASMEANS(points)
3:   dists  $\leftarrow$  MEANSDIST(means)
4:   minDist  $\leftarrow$  MINMEANSDIST(dists)
5:   count  $\leftarrow$  points.size
6:   for  $i \leftarrow k + 1$  to count do
7:     point  $\leftarrow$  points[ $i$ ]
8:     distsFromPointToMeans  $\leftarrow$  DISTSTOMEANS(point, means)
9:     minDistPoint  $\leftarrow$  MINDISTTOMEANS(distsFromPointToMeans)
10:    if minDistPoint.dist < distOfMinDist.dist then
11:      UPDITEMEANSWITHPOINT(means, dists, minDist, point, minDistPoint.toCentroid)
12:    else
13:      UPDITEMEANSWITHMERGE(means, dists, minDist, point, distsFromPointToMeans)
14:    end if
15:  end for
16:  return means
17: end function
```

---

---

```

1: function CENTROIDS(points, means)
2:   partMeans  $\leftarrow$  PORTIONMEANS(points, means)
3:   return NEWMEANSFROMPORTIONMEANS(partMeans)
4: end function

```

---



---

```

1: function ARECLOSE(means1, means2)
2:   for all mean : means1 do
3:     if not HASCLOSE(mean, means) then
4:       return false
5:     end if
6:   end for
7:   return true
8: end function

```

---



---

```

1: function PARTITION(points, means)
2:   clusters[K]
3:   for i  $\leftarrow$  1 to K do
4:     clusters[i].centroid  $\leftarrow$  means[i].centroid
5:   end for
6:   count  $\leftarrow$  points.size
7:   for index  $\leftarrow$  1 to count do
8:     point  $\leftarrow$  points[index]
9:     centroidIndex  $\leftarrow$  1
10:    d  $\leftarrow$  DIST(point, means[1].centroid)
11:    for j  $\leftarrow$  2 to K do
12:      dj  $\leftarrow$  DIST(point, means[j].centroid)
13:      if dj < d  $\vee$  dj = d  $\wedge$  clusters[j].size < clusters[centroidIndex].size then
14:        centroidIndex  $\leftarrow$  j
15:        d  $\leftarrow$  dj
16:      end if
17:    end for
18:    PUSH(clusters[centroidIndex].indecies, index)
19:  end for
20:  return clusters
21: end function

```

---



---

```

1: function FIRSTASMEANS(points)
2:   means[K]
3:   for i  $\leftarrow$  1 to K do
4:     means[K]  $\leftarrow$  POINTMEAN(points[i])
5:   end for
6:   return means
7: end function

```

---

---

```

1: function MEANSDISTS(means)
2:    $dists[K][K]$ 
3:   for  $i \leftarrow 1$  to  $K$  do
4:     for  $j \leftarrow 1$  to  $i - 1$  do
5:        $d \leftarrow \text{DIST}(\text{means}[i].\text{centroid}, \text{means}[j].\text{centroid})$ 
6:        $dists[i][j] \leftarrow d$ 
7:        $dists[j][i] \leftarrow d$ 
8:     end for
9:      $dists[i][i] \leftarrow 0$ 
10:  end for
11:  return  $dists$ 
12: end function

```

---



---

```

1: function MINMEANSDIST(dists)
2:    $from \leftarrow 1$ 
3:    $to \leftarrow 2$ 
4:    $d \leftarrow dists[from][to]$ 
5:   for  $i \leftarrow 1$  to  $K$  do
6:     for  $j \leftarrow 1$  to  $i - 1$  do
7:       if  $dists[i][j] < d$  then
8:          $from \leftarrow i$ 
9:          $to \leftarrow j$ 
10:         $d \leftarrow dists[from][to]$ 
11:      end if
12:    end for
13:  end for
14:  return  $\langle from, to, d \rangle$ 
15: end function

```

---



---

```

1: function DISTSTOMEANS(point, means)
2:    $dists[K]$ 
3:   for  $i \leftarrow 1$  to  $K$  do
4:      $dists[i] \leftarrow \text{DIST}(\text{point}, \text{means}[i].\text{centroid})$ 
5:   end for
6:   return  $dists$ 
7: end function

```

---



---

```

1: function MINDISTTOMEANS(dists)
2:    $to \leftarrow 1$ 
3:    $d \leftarrow dists[to]$ 
4:   for  $i \leftarrow 2$  to  $K$  do
5:     if  $dists[i] < d$  then
6:        $to \leftarrow i$ 
7:        $d \leftarrow dists[to]$ 
8:     end if
9:   end for
10:  return  $\langle to, d \rangle$ 
11: end function

```

---

---

```

1: function UPDATEMEANSWITHPOINT(means, dists, minDist, point, centroidIndex)
2:   means[centroidIndex]  $\leftarrow$  UPDATEMEAN(means[centroidIndex], point)
3:   for i  $\leftarrow$  1 to K do
4:     if i  $\neq$  centroidIndex then
5:       d  $\leftarrow$  DIST(means[centroidIndex].centroid, means[i].centroid)
6:       dists[centroidIndex][i]  $\leftarrow$  d
7:       dists[i][centroidIndex]  $\leftarrow$  d
8:       if d < minDist.dist then
9:         minDist  $\leftarrow$  < centroidIndex, i, d >
10:      end if
11:    end if
12:  end for
13: end function

```

---



---

```

1: function UPDATEMEANSWITHMERGE(means, dists, minDist, point, distsToMeans)
2:   means[minDist.from]  $\leftarrow$  MERGEMEANS(means[minDist.from], means[minDist.to])
3:   means[minDist.to]  $\leftarrow$  POINTMEAN(point)
4:   for i  $\leftarrow$  1 to K do
5:     if i  $\neq$  minDist.from then
6:       d  $\leftarrow$  DIST(means[minDist.from].centroid, point)
7:       dists[minDist.from][i]  $\leftarrow$  d
8:       dists[i][minDist.from]  $\leftarrow$  d
9:     end if
10:    if i  $\neq$  minDist.to  $\wedge$  (i  $\neq$  minDist.from) then
11:      dists[minDist.to][i]  $\leftarrow$  distsToMeans[i]
12:      dists[i][minDist.to]  $\leftarrow$  dists[i][minDist.to]
13:    end if
14:  end for
15:  minDist  $\leftarrow$  MINDISTTOMEANS(dists)
16: end function

```

---



---

```

1: function PORTIONMEANS(points, means)
2:   portion  $\leftarrow$   $\frac{\text{points.size}}{H}$ 
3:   futers[H]
4:   for i  $\leftarrow$  1 to H do
5:     start  $\leftarrow$  (i - 1) * portion + 1
6:     end  $\leftarrow$  start + portion
7:     if i = H then
8:       end  $\leftarrow$  points.size
9:     end if
10:    futers[i]  $\leftarrow$  ASYNC(task, points, means, start, end)
11:  end for
12:  return futers
13: end function

```

---



---

```

1: function NEWMEANSFROMPORTIONMEANS(futers)
2:   means  $\leftarrow$  ASYNCRESULT(futers[1])
3:   for h  $\leftarrow$  2 to H do
4:     tempMeans  $\leftarrow$  ASYNCRESULT(futers[h])
5:     for i  $\leftarrow$  1 to K do
6:       if means[i].size = 0  $\wedge$  tempMeans[i].size  $\neq$  0 then
7:         means[i]  $\leftarrow$  tempMeans[i]
8:       else
9:         if means[i].size  $\neq$  0  $\wedge$  tempMeans[i].size  $\neq$  0 then
10:          means[i]  $\leftarrow$  MERGEMEANS(means[i], tempMeans[i])
11:        end if
12:      end if
13:    end for
14:  end for
15:  return means
16: end function

```

---



---

```

1: function TASK(points, means, start, end)
2:   tempMeans[K]
3:   for i  $\leftarrow$  1 to K do
4:     tempMeans[i].size  $\leftarrow$  0
5:   end for
6:   for i  $\leftarrow$  start to end do
7:     point  $\leftarrow$  points[i]
8:     centroidIndex  $\leftarrow$  1
9:     d  $\leftarrow$  DIST(point, means[centroidIndex].centroid)
10:    for j  $\leftarrow$  2 to K do
11:      dj  $\leftarrow$  DIST(point, means[j].centroid)
12:      if dj < d  $\vee$  dj = d  $\wedge$  tempMeans[j].size < tempMeans[centroidIndex].size then
13:        centroidIndex  $\leftarrow$  j
14:        d  $\leftarrow$  dj
15:      end if
16:    end for
17:    if tempMeans[centroidIndex].size = 0 then
18:      tempMeans[centroidIndex]  $\leftarrow$  POINTMEAN(point)
19:    else
20:      tempMeans[centroidIndex]  $\leftarrow$  UPDATEMEAN(tempMeans[centroidIndex], point)
21:    end if
22:  end for
23:  return tempMeans
24: end function

```

---



---

```

1: function HASCLOSE(mean, means)
2:   for i  $\leftarrow$  1 to K do
3:     if DIST(mean.centroid, means[i].centroid) < e then
4:       return true
5:     end if
6:   end for
7:   return false
8: end function

```

---

---

```
1: function POINTMEAN(point)
2:   return < point, 1 >
3: end function
```

---

---

```
1: function UPDATEMEAN(mean, point)
2:   tempMean ← POINTMEAN(point)
3:   return MERGEMEANS(mean, tempMean)
4: end function
```

---

---

```
1: function MERGEMEANS(m1, m2)
2:   count ← m1.size + m2.size
3:   a ←  $\frac{m1.size}{count}$ 
4:   b ←  $\frac{m2.size}{count}$ 
5:   m
6:   for i ← 1 to K do
7:     m[i] ← a.m1.centroid[i] + b.m2.centroid[i]
8:   end for
9:   return < m, count >
10: end function
```

---

## 4 Инструкции за компилиране на програмата

На Linux базирана операционна система се нуждаете от **g++**, поне **5.4.0** версия. За компилиране се преместете в папката на сорс кода и от там изпълнете командата **make**. Тази команда ще стартира компилирането на кода и ще създаде изходен файл **a.out**. За да изпълните изходния файл изпълнете командата **./a.out**.

## 5 Примерни резултати

Използвани са 6 различни функции за определяне на разстоянието.

- L1 дистанция известна още като Манхатанско разстояние.

$$l1(x, y) = \sum_{i=1}^d |x_i - y_i|$$

- L2 дистанция известна още като Евклидово разстояние.

$$l2(x, y) = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

- L3 дистанция.

$$l3(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^3 \right)^{\frac{1}{3}}$$

- L4 дистанция.

$$l_4(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^4 \right)^{\frac{1}{4}}$$

- $L_\infty$  дистанция.

$$l_\infty(x, y) = \max\{|x_i - y_i| \mid i = 1, \dots, d\}$$

- S дистанция.

$$s(x, y) = \sum_{i=1}^d (x_i - y_i)^2$$

За всеки от получените клъстери се определя класа, който доминира в клъстера и се пресмята какво е отношението на точките от този клас към точките в клъстера. За всяка дистанция се пресмята и оценяща функция на чистота на клъстеризацията. Наричана *purity* по следната формула:

$$purity(\{C_1, C_2, C_3\}) = \frac{1}{n} \sum_{i=1}^3 \max\{|C_i \cap Class_l| \mid l \in \{Setosa, Versicolour, Virginica\}\}$$

Тази функция дава оценка за качеството на клъстеризация. За всяка функция на разстояние се извежда името на функцията. След, което за всеки клъстер се намира броя на точките от доминиращия клас към размера на клъстера и се извежда името на доминиращия клас. Накрая се извежда стойността на *purity* функцията за това клъстеризиране.

## 5.1 Изходни резултати:

L1  
 Virginica: 0.0540541  
 Versicolor: 0.234375  
 Setosa: 0  
 0.887417  
 L2  
 Versicolor: 0.222222  
 Setosa: 0  
 Virginica: 0.0526316  
 0.89404  
 L3  
 Versicolor: 0.222222  
 Setosa: 0  
 Virginica: 0.0526316  
 0.89404  
 L4  
 Setosa: 0  
 Virginica: 0.494949  
 Setosa: 0.0454545  
 0.662252  
 L\_infinity  
 Versicolor: 0.5  
 Setosa: 0  
 Setosa: 0.0232558

0.662252

S

Versicolor: 0.222222

Setosa: 0

Virginica: 0.0526316

0.89404