

MC970/MO644 - Projeto

Paralelização do Algoritmo Genético para a resolução do problema Min-Sum K-Clustering com PThreads

Felipe Augusto de Castro e Silva, RA 145993

Rafael Grisotto e Souza, RA 192765

1 - O problema Min K-Clustering Sum

O problema Minimum K-Clustering Sum é NP-completo e considera um conjunto finito X de pontos no plano, uma distância $d(x, y) \in \mathbb{R}$ para cada par $x, y \in X$, com todas as distâncias respeitando a regra de desigualdade dos triângulos. A partir disso, pretende-se encontrar uma partição de X em subconjuntos C_1, C_2, \dots, C_k , de forma a minimizar as somas de todas as distâncias entre elementos em um mesmo subconjunto, portanto, sendo k o número de subconjuntos procurados, a meta é minimizar a seguinte expressão:

$$\sum_{i=1}^k \sum_{v_1, v_2 \in C_i} d(v_1, v_2)$$

2 - Algoritmo Genético

A literatura contém poucos exemplos relevantes de abordagem do problema Minimum K-Clustering Sum por algoritmos genéticos, com isso e levando em conta que algoritmos genéticos mostram bons resultados provados empiricamente numa série de problemas e que boas respostas demandam tempos de execução muito grandes foi decidido que tal abordagem paralelizada será explorada neste trabalho.

Alguns outros trabalhos oferecem uma explicação profunda sobre como funcionam os algoritmos genéticos, contudo é importante definir aqui qual será o algoritmo genético base utilizado neste trabalho.

Cada indivíduo da população é um vetor de inteiros em que cada posição corresponde a um ponto de entrada e que pode assumir valores entre 0 e $K - 1$, ou seja, o cluster que tal ponto está localizado. Para o algoritmo base, primeiro é definida uma população inicial aleatória de tamanho passado por parâmetro. Depois são executadas um número suficiente de iterações; para cada iteração, escolhemos aleatoriamente parte da população atual como pais da próxima população, combinamos os pais para a geração de novos indivíduos, mutamos os filhos de forma aleatória. Por último, são selecionados os melhores indivíduos para permanecerem na próxima população.

3 - Desafios

O maior desafio deste trabalho foi conseguir um speedup relevante para o algoritmo genético, devido sua organização: de maneira geral, podemos definir o algoritmo genético como um loop que é executado um grande número de vezes, em que cada vez, algumas etapas são performadas, de forma a gerar uma nova população, que lidere a uma resposta melhor que a já encontrada, mas que dê margem à diversidade, para que não a execução não fique parada em mínimos locais.

A primeira ideia foi a de paralelizar o loop principal do algoritmo, contudo, este é DO-ACROSS, pois a população atual depende da população anterior para ser gerada. Notamos que teríamos que buscar paralelizar cada etapa do algoritmo e que isso poderia não ser eficiente, pois como cada etapa é acessada muitíssimas vezes, o overhead para a criação e gestão de threads poderia superar seus possíveis benefícios.

A partir daqui, rodamos o profiling, que está descrito na seção 4 e observamos que a função que calcula o fitness de cada elemento da população era o grande gargalo da aplicação. Contudo, nossos problemas relacionados ao overhead se agravaram, pois algumas etapas do loop principal chamam um número considerável de vezes a função fitness e ficar criando nova threads poderia ser muito caro. Nossas previsões foram confirmadas quando tentamos paralelizar o código com OpenMP, na qual todos os nossos testes não lideravam speedups e o profiling nos mostrava que estávamos perdendo para o overhead das threads.

Tentamos uma abordagem com pthreads, que começou a se mostrar útil na utilização de instâncias consideradas grandes para o problema, aliado a um processador que suporta 8 threads. Depois disso, tentamos localizar outros trechos do código que poderiam ser paralelizados e foi identificado que a função de mutação continha um loop DO-ALL. Contudo, sua paralelização foi abandonada, pois nada contribuía com o speedup. Depois de alguns resultados razoáveis, decidimos seguir em frente com novas ideias de paralelização.

O que nos veio em mente é que apesar da maioria dos trechos que chamavam a função de fitness eram loops DO-ACROSS, a função que mais chamava o fitness (geração dos pais da nova população), continha um loop DO-ALL. Foram observados speedups para essa abordagem, mas ela se mostrou pior que somente o paralelismo da função fitness, pois ela demanda a utilização de locks.

Por fim, um segundo problema foi abordado: é extremamente difícil encontrar bons parâmetros que liderem seu algoritmo genético à melhor resposta possível. Com isso, decidimos testar a criação de várias populações com vários parâmetros diferentes em um mesmo algoritmo, paralelizar suas execuções, fazendo com que a melhor resposta possa ser utilizada.

4 - Profiling

Utilizamos o pacote gprof e perf para observarmos as funções mais “quentes” do projeto. A partir do gprof, identificamos a função fitness pelo seu maior número de chamadas e a função selection, como são mostradas nos reports do gprof:

```
Flat profile:
Each sample counts as 0.01 seconds.
%   cumulative   self           self         total
time  seconds    seconds   calls   ms/call  ms/call  name
41.48    13.32    13.32   301079    0.04    0.10  fitness(std::vector<int, std::allocator<int> >)
19.62    19.63     6.30  2916918392    0.00    0.00  std::vector<std::pair<int, int>, std::allocator<std::pair<int, int> > >::operator[](unsigned long)
18.19    25.47     5.84  1587034370    0.00    0.00  std::vector<std::pair<int, int>, std::allocator<std::pair<int, int> > >::size() const
14.33    30.67     4.60  3056331458    0.00    0.00  std::vector<int, std::allocator<int> >::operator[](unsigned long)
1.87    30.67     0.00         1   600.20   600.20  std::vector<std::pair<int, int>, std::allocator<std::pair<int, int> > >::vector()
1.81    31.25     0.58  364614799    0.00    0.00  _gnu_cxx::_enable_if<std::_is_integer<int>::__value, double>::__type std::sqrt<int>(int)
```

index	% time	self	children	called	name
					<spontaneous>
[1]	98.1	0.00	31.52		main [1]
		0.01	10.14	1000/1000	selection(std::vector<std::vector<int, std::allocator<int> >, std::allocator<std::vector<int, std::allocator<int> > > > [4]
		0.00	10.14	1000/1000	selectPopulation(std::vector<std::vector<int, std::allocator<int> >, std::allocator<std::vector<int, std::allocator<int> > > >, std::vector<std::vector<int, std::allocator<int> >, std::allocator<int> > > > [5]
8]		0.28	5.41	1000/1000	diversifyPopulation(std::vector<std::vector<int, std::allocator<int> >, std::allocator<std::vector<int, std::allocator<int> > > >, double) [
		0.00	5.06	1000/3000	getBest(std::vector<std::vector<int, std::allocator<int> >, std::allocator<std::vector<int, std::allocator<int> > > > [3]
		0.00	0.24	1000/1000	crossover(std::vector<std::vector<int, std::allocator<int> >, std::allocator<std::vector<int, std::allocator<int> > > > [16]
		0.05	0.66	1079/301079	fitness(std::vector<int, std::allocator<int> > [2]
		0.03	0.02	1000/1000	mutation(std::vector<std::vector<int, std::allocator<int> >, std::allocator<std::vector<int, std::allocator<int> > > >, double) [20]

Na report do perf também é mostrado a função *fitness* como “quente”:

Overhead	Command	Shared Object	Symbol
46.08%	seq_prof	libc-2.23.so	[.] 0x000000000109574
12.71%	seq_prof	libc-2.23.so	[.] _mcount
9.48%	seq_prof	seq_prof	[.] fitness
5.14%	seq_prof	seq_prof	[.] std::vector<std::pair<int, int>, std::allocator<std::pair<int, int> > >::operator[]
3.95%	seq_prof	seq_prof	[.] std::vector<std::pair<int, int>, std::allocator<std::pair<int, int> > >::size
3.42%	seq_prof	seq_prof	[.] std::vector<int, std::allocator<int> >::operator[]
2.04%	seq_prof	libc-2.23.so	[.] 0x00000000010957f
2.02%	seq_prof	libc-2.23.so	[.] 0x0000000001095ff
2.01%	seq_prof	libc-2.23.so	[.] 0x00000000010959c
2.00%	seq_prof	libc-2.23.so	[.] 0x0000000001095e1
2.00%	seq_prof	libc-2.23.so	[.] 0x0000000001095b5
1.97%	seq_prof	libc-2.23.so	[.] 0x000000000109596
1.34%	seq_prof	libc-2.23.so	[.] 0x000000000109560
1.32%	seq_prof	libc-2.23.so	[.] 0x000000000109568
1.24%	seq_prof	seq_prof	[.] mcount@plt
0.67%	seq_prof	libc-2.23.so	[.] 0x000000000109566

Em suma, a partir dos dados do profiling identificamos as funções *fitness* e *selection*, porém pela característica de algoritmo genético, resolvemos paralelizar as populações também e como foi feito está na próxima seção.

5 - Como o código foi paralelizado

Função de *fitness* - a função de *fitness* é composta por dois loops encadeados (complexidade da função $O(n^2)$) e calcula a expressão que deve ser minimizada. Para sua paralelização, foi criada uma nova função chamada *fit*, do tipo ponteiro de void e que tinha como parâmetro também um ponteiro de void. Na função de *fitness*, são atribuídos um número uniforme de intervalos do loop mais externo que é executado por cada thread. A partir daí, uma struct com todos os parâmetros necessários é instanciada e passada para a função *fit*, que é atribuída e executada por uma thread, através da função *pthread_create()*. Por fim, esperamos que todas as threads terminem suas execuções através da função *pthread_join()*, somamos os cálculos de todas as threads em uma variável e retornamos para a função que a chamou.

```

14 vector<int> cromossomeFit;
15 vector<vector<int> > popMut;
16
17 void *fit(void *aux){
18     Aux *la = (Aux *) aux;
19     int ans = 0;
20     int l = la->l;
21     int r = la->r;
22     for(int i = l; i < r && i < points.size(); i++){
23         for(int j = i + 1; j < points.size(); j++){
24             if(cromossomeFit[i] == cromossomeFit[j]){
25                 ans += floor(0.5 + sqrt((points[i].first - points[j].first) * (points[i].first - points[j].first) + (points[i].second - points[j].second) * (points[i].second - points[j].second)));
26             }
27         }
28     }
29     return (void *) ans;
30 }
31
32 pthread_t threads[NUM_THREADS];
33 bool pala = false;
34
35 int fitness(vector<int> &cromossome){
36     cromossomeFit = cromossome;
37     for(int i = 0; i < NUM_THREADS; i++){
38         Aux *aux = (Aux *) malloc(sizeof(aux));
39         aux->l = i * points.size() / NUM_THREADS;
40         aux->r = (i + 1) * points.size() / NUM_THREADS;
41         if(i == NUM_THREADS - 1)
42             aux->r = points.size();
43         pthread_create(&threads[i], NULL, fit, (void *) aux);
44     }
45     int ans = 0;
46     int st;
47     for(int i = 0; i < NUM_THREADS; i++){
48         pthread_join(threads[i], (void **) &st);
49         ans += (int) st;
50     }
51     return ans;
52 }

```

Função de escolha de pais para a próxima população - nesta função, temos um loop que roda por um número de iterações dada como parâmetro que sorteia dois elementos da população atual (note que aqui utilizamos a função `rand_r` para gerar números aleatórios, já que a função `rand()` não é thread safe), escolhe o melhor através do cálculo de seus respectivos fitness e adiciona o melhor numa lista de pais. Para a paralelização desta função, criamos uma nova função chamada `sel`, do tipo ponteiro de void, que recebe um parâmetro também do tipo ponteiro de void. O número de iterações é dividida uniformemente pelo número de threads e a partir daí, uma struct com todos os parâmetros necessários é instanciada e passada para a função `sel`, que é atribuída e executada por uma thread, através da função `pthread_create()`. Por fim, esperamos que todas as threads terminem suas execuções através da função `pthread_join()` e retornamos a lista de pais. É importante ressaltar aqui que a lista de pais é uma variável global e só pode ter um acesso por vez, por isso, utilizamos as funções `pthread_mutex_lock()` e `pthread_mutex_unlock()` para organizar o acesso à lista e evitar qualquer erro.

```

53 typedef struct{
54     int esq, dir;
55 } Aux;
56
57 pthread_mutex_t locked;
58 unsigned int seed;
59
60 void *sel(void *aux){
61     Aux *la = (Aux *) aux;
62     int esq = la->esq;
63     int dir = la->dir;
64     for(int i = esq; i < dir; i++){
65         int l = rand_r(&seed) % popAux.size();
66         int r = rand_r(&seed) % popAux.size();
67         if(fitness(popAux[l]) < fitness(popAux[r])){
68             pthread_mutex_lock(&locked);
69             parents.push_back(popAux[l]);
70             pthread_mutex_unlock(&locked);
71         } else {
72             pthread_mutex_lock(&locked);
73             parents.push_back(popAux[r]);
74             pthread_mutex_unlock(&locked);
75         }
76     }
77 }
78
79 vector<vector<int> > selection(vector<vector<int> > pop){
80     parents.clear();
81     popAux = pop;
82     pthread_t threads[NUM_THREADS];
83     for(int i = 0; i < NUM_THREADS; i++){
84         Aux *aux = (Aux *) malloc(sizeof(aux));
85         aux->esq = i * pop.size() / NUM_THREADS;
86         aux->dir = (i + 1) * pop.size() / NUM_THREADS;
87         pthread_create(&threads[i], NULL, sel, (void *) aux);
88     }
89     for(int i = 0; i < NUM_THREADS; i++){
90         pthread_join(threads[i], NULL);
91     }
92     return parents;
93 }

```

Várias populações - como nos casos anteriores, é criada uma nova função e aqui adicionamos o loop principal do algoritmo genético a ela. Cada set de parâmetros passado como entrada é atribuído a uma thread. A partir daqui, uma struct com todos os parâmetros necessários é instanciada e passada para a função exe, que é atribuída e executada por uma thread, através da função `pthread_create()`. Por fim, esperamos que todas as threads terminem suas execuções através da função `pthread_join()` e printamos a melhor resposta de todas as populações.

```

163 typedef struct{
164     int P, generations, sizePopulation;
165     double mutationRate, divRate;
166 } AuxExe;
167
168 void *exe(void *aux){
169     AuxExe *la = (AuxExe *) aux;
170     int generations = la->generations;
171     int sizePopulation = la->sizePopulation;
172     int P = la->P;
173     double mutationRate = la->mutationRate;
174     int divRate = la->divRate;
175     int bestSol = INT_MAX;
176     vector<vector<int>> pop = initialPopulation(sizePopulation, P);
177
178     for(int i = 0; i < generations; i++){
179         vector<vector<int>> parents = selection(pop);
180         vector<vector<int>> cross = crossover(parents);
181         mutation(cross, mutationRate);
182         pop = selectPopulation(pop, cross);
183         diversifyPopulation(pop, divRate);
184
185         int idx = getBest(pop);
186
187         if(bestSol > fitness(pop[idx])){
188             bestSol = fitness(pop[idx]);
189         }
190     }
191     return (void *) bestSol;
192 }
193
194 pthread_t threads[numPopulations];
195
196 for(int i = 0; i < numPopulations; i++){
197     scanf("%d %lf %lf %d", &generations, &mutationRate, &divRate, &multSize);
198     generations = 10;
199
200     sizePopulation = ceil(log2(points.size())) * multSize;
201
202     if(sizePopulation % 2 == 1) sizePopulation++;
203
204     mutationRate = mutationRate / (double) points.size();
205
206     AuxExe *aux = (AuxExe *) malloc(sizeof(AuxExe));
207     aux->P = P;
208     aux->generations = generations;
209     aux->mutationRate = mutationRate;
210     aux->divRate = divRate;
211     aux->sizePopulation = sizePopulation;
212
213     pthread_create(&threads[i], NULL, exe, (void *) aux);
214 }
215
216 int ans = INT_MAX, st;
217
218 for(int i = 0; i < numPopulations; i++){
219     pthread_join(threads[i], (void **) &st);
220
221     // cout << (int) st << endl;
222 }

```

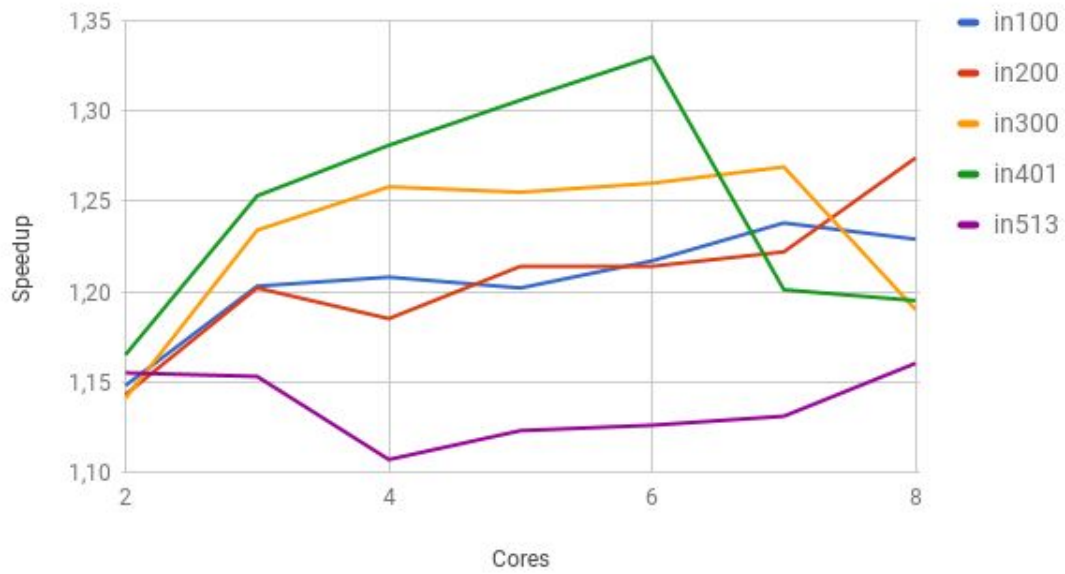
6 - Gráfico de speedup

Fizemos experimentos para os três tipos de paralelização, segue abaixo as figuras:

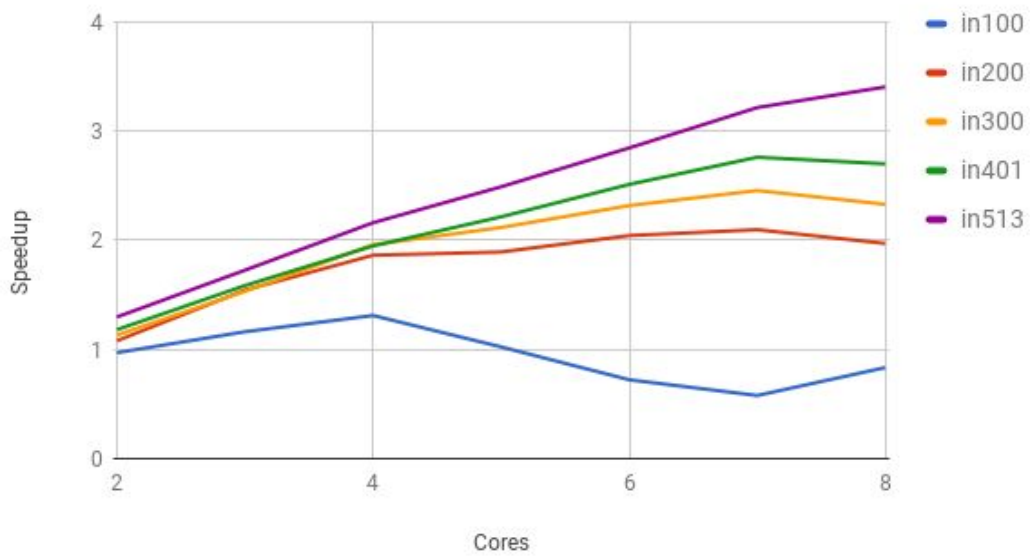
Para todas as formas de paralelização usamos 10 mil a quantidade de gerações para ficar mais real com seu uso prático na literatura, também definimos como 8 cores para o método de *Selection* e *Fitness*. Porém, para o método de múltiplas populações usamos 6 cores para executar 6 populações paralelamente e o serial 6 vezes com configurações dos parâmetros distintos.

Podemos observar nas figuras que todas as instâncias tiveram formas de paralelização tiveram *speedup*. Destacando para o método de múltiplas populações e para a função *fitness* que tiveram *speedup* acima de 3.

Paralelização função Selection



Paralelização função Fitness



Paralelização de Multiplas populações

