## trab2-knnMPI: Determinar conjuntos de K vizinhos mais próximos (KNN) com processos MPI

-----

Histórico:

- v1.0 a versão inicial, - v2.1 adicionado pequenos detalhes, retirado um erro no texto

Data do enunciado: 06/nov/2023 Data da entrega: **19/nov/2023** 

PS: note que esse trabalho reaproveita codigo dos exercicios feitos nas duas ultimas aulas bem como parte do trab 1

Esse trab pode ser em grupo de no MAXIMO 2 alunos o nome dos autores devem estar no relatório e no progr. Fonte A entrega será via UFPR Virtual

-----

Objetivo do trabalho 2:

O objetivo deste exercício é fazer um trabalho que use paralelismo em MPI no cluster com rede ethernet e verificar a aceleração da nossa solução.

## O trabalho:

A entrada para nossa função paralela são dois conjuntos de pontos de D dimensoes cada ponto. São pontos de muitas dimensões.

Dados os conjuntos de n pontos P e Q, denominados respectivamente

P = Pontos da base (dataset)

(matriz de floats)

Q = Pontos de consulta, send |Q| o tamanho de Q.

(matriz de floats)

Dado um número inteiro k

Para cada ponto em Q, queremos determinar quais são seus k vizinhos mais próximos no conjunto P.

Fazer uma funçao knn que recebe como parâmetros: Q, nq, P, n, D, k

sendo: nq = numero de pontos em Q

n = numero de pontos em P

D = número de dimensoes dos pontos

k tamanho dos conjuntos de pontos vizinhos buscados

Note que o resultado (saída) da funçao knn deve ter nq conjuntos de tamanho k Cada conjunto deve ter o indice de seus k vizinhos mais próximos

Ou seja, a saida de knn pode ser uma matriz de nq linhas por k colunas.

A figura na pagina seguinte ilustra o problema.

Nosso programa MPI deve obter as matrizes Q e P.

O conteúdo dessas matrizes poderia ser obtido de arquivos MAS

como o conjunto P pode ser grande, a leitura pode demorar,

então vamos gerar esses conjunto com dados de maneira aleatoria no processo 0 MPI (raiz).

Usando uma função geraConjuntoDeDados (C, nc, D)

que gera um conjuto qualquer C de nc pontos aleatorios de D dimensoes.

Todos as nossas coordenadas dos porntos em Q e P devem ser float.

Somente o processo 0 deve gerar os conjuntos Q e P.

Para nossas experiencias queremos um conjunto Q de 128 pontos de 300 dimensoes e um conjunto P de 400mil pontos de 300 dimensoes

- SOBRE A IMPLEMENTACAO:
- Sua solução DEVE usar a adequadamente implementação do HEAP (chave, valor) e operação deccreaseMax feitos no trab 1
- Voce DEVE usar a operação scatter de MPI na matriz **Q** da raiz para os nodos e a operação broadcast para enviar a matriz **P** para todos os nodos, e gather para "juntar" o resultado R na RAIZ
- todas as distancias calculadas, comparadas e mantidas no heap são distancias quadraticas (distancia ao quadrado d^2, ou seja, NAO vamos usar operacoes de raiz quadrada e nosso programa rodará mais rapido
- Todos os resultados reportados deve ser medidos no nosso cluster XEON e você deve fornecer os scripts para compilacao e para rodar suas experiencias, bem como as planilhas openoffice com seus resultados e graficos ALEM do relatório
- Ordem dos parametros para rodar o programa

```
\begin{split} & \text{mpirun -np 4 knn-mpi } & \text{nq npp d k} \\ & \text{onde: } & \text{nq} = \text{numero de pontos em } Q \quad \text{(tamanho do conjunto de pesquisa)} \\ & \text{npp} = \text{numero de pontos em } P \quad \text{(dataset)} \\ & \text{d} \quad = \text{numero de dimensoes dos pontos (dimensionalidade)} \end{split}
```

k = tamanho de cada conjunto de vizinhos (k vizinhos por ponto)

ex:

mpirun -np 4 knn-mpi 128 400000 300 1024

Calcula os 1024 vizinhos mais próximos de 128 pontos de 300 dimensoes em base de dados de 400mil pontos

## SOBRE AS EXPERIÊNCIAS:

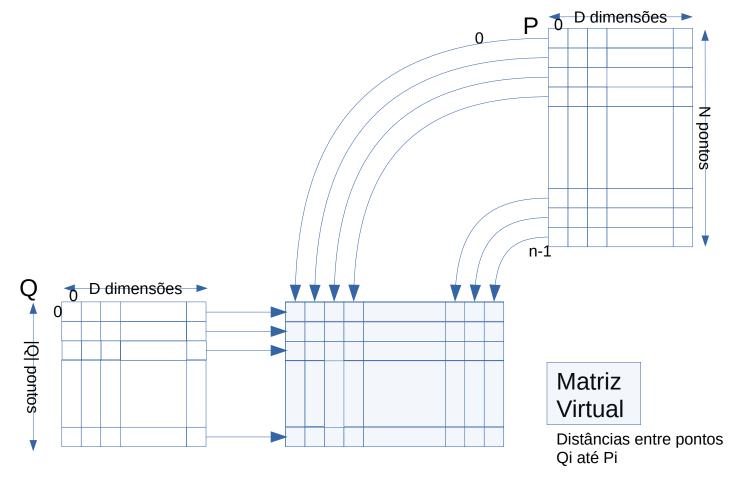
- TODAS os resultados reportados DEVEM ser
- Faremos a computação total para 3 casos apenas:
- Experiencia 1:
- Rodar o programa para APENAS 1 processo MPI e medir o tempo da computação de knn
- Experiencia 2:
- Rodar o programa para 4 processos MPI no mesmo host e medir o tempo da computação de knn
- Experiencia 3:
- Rodar o programa para 4 processos MPI em hosts diferentes e medir o tempo da computação de knn
- Reportar a aceleração (speedup) da sua implementação paralela em relação a 1 processo apenas, nos casos da experiencia 2 e 3.
- SOBRE O RELATÓRIO do trabalho:
  - Fazer um relatório (PDF) com suas conclusões. Descreva as características importantes da sua CPU usada. Inclua em apendice ao relatório o texto da saida do comando lscpu bem como a figura do comando lstopo para sua CPU usada.
- SOBRE A funçao de vefificação a ser chamada (**próxima página**)

- SOBRE A funçao de vefificação a ser chamada
  - Incluir ao final do seu programa, ANTES de terminar o MPI e ANTES de desalocar suas matrizes, para rodar APENAS no processo 0 MPI, a funcao verificaKNN(...) que DEVE ter o mesmo prototipo abaixo. Voce deve colocar a funcao inicial de verificacao abaixo em um arquivo chamado verificaKNN.c e fazer o include desse arquivo no INICIO do seu programa.
  - A função inicial está abaixo, esta apenas vai imprimir o vetor de resultados. Para a verificação adequada o prof irá substituir essa versão por outra que fara a verificação.

```
void verificaKNN( float *Q, int nq, float *P, int n, int D, int k, int *R ) {
    // note que R tem nq linhas por k colunas, para qualquer tamanho de k (colunas)
    // entao é linearizado para acesso como um VETOR
    printf( " ------- VERIFICA KNN ------ " );
    for( int linha=0; linha<nq; linha++ ) {
        printf( "knn[%d]: ", linha );
        for( int coluna=0; coluna<k; coluna++ )
            printf( "%d ", R[ linha*k+coluna ] );
        printf( "\n" );
    }
}</pre>
```

OBS: será na página seguinte temos a ilustraçao do prolema

## ilustração do problema



Nao precisa ser armazenada em memoria, basta calcular as distancias indicadas entre os pontos

Para comparar distancias, podemos evitar a raiz quadrada comparando o quadrado das distancias (sem extrair a raiz)