

Open-MP: Generación de Matriz de Covarianzas

Pedro Vladimir Hernández Serrano, José Alfredo Méndez Barrera y Elisa Hernández Rodríguez

I. INTRODUCCIÓN

El siguiente documento aborda la descripción de la codificación, compilación y ejecución de una matriz de covarianzas en paralelo y de manera secuencial. La ejecución secuencial se realiza en C. Por otra parte, la ejecución en paralelo es posible a través de la herramienta Open-MP.

El uso de esta matriz, particularmente en *big data*, requiere poder de cómputo y eficiencia en el código. En la siguiente sección se describe la forma en que una matriz de covarianzas puede ejecutarse.

II. INVESTIGACIÓN DEL PROBLEMA

La matriz de Covarianzas es muy utilizada en distintos problemas, ya que ayuda a conocer la relación que existe entre las variables que componen la matriz origen.

La covarianza es una medida estadística de la relación lineal entre dos variables. Si la covarianza entre dos variables es alta, un cambio en una variable significa un cambio lineal similar en la variable relacionada. Cuando la covarianza es muy negativa, las dos variables están inversamente relacionadas entre sí.

En particular, esta matriz juega un papel importante en la minería de datos y para clusterización por medio de **componentes principales basados en las covarianzas**.

- El análisis de componentes principales lo utiliza para ordenar las variables de manera decreciente con respecto a su variabilidad y así el número de variables puede ser reducido a aquellas de mayor alcance.
- La clusterización requiere la recuperación de información de grandes bases de datos por lo que es conveniente reducir la dimensionalidad del conjunto de datos, y la idea principal es la de implementar una técnica para formar una matriz de covarianza de los datos y encontrar los vectores propios. Los componentes principales se utilizan para formar un nuevo espacio, más pequeño dimensionalmente con puntos de datos proyectados de los datos originales.
- La matriz de covarianza también se utiliza en la minería de datos. La clasificación es una forma de crear normas que rigen la categoría de un objeto de datos. Estas reglas se construyen mediante el examen de un pequeño conjunto de datos ya categorizados. Los algoritmos utilizan estas reglas para separar los nuevos valores de los datos en clases existente.

De manera formal podemos definir la covarianza matemática dadas las variables X, Y como se muestra en la ecuación 1 de igual manera podemos ver en la ecuación 4 una forma alternativa.

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])], \quad (1)$$

En donde $E[X]$ es el valor esperado de X .

$$\begin{aligned} \text{cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY - XE[Y] - E[X]Y + E[X]E[Y]] \\ &= E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y] \\ &= E[XY] - E[X]E[Y]. \end{aligned}$$

Que también puede ser expresado como sigue:

$$\text{cov}(X, Y) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} (x_i - x_j) \cdot (y_i - y_j) \quad (2)$$

$$= \frac{1}{n^2} \sum_i \sum_{j>i} (x_i - x_j) \cdot (y_i - y_j) \quad (3)$$

$$= \frac{1}{n} \sum_i (x_i - \bar{x}) \cdot (y_i - \bar{y}) \quad (4)$$

Cada elemento de la matriz es una suma de productos de los elementos en la matriz de entrada tomados de una ventana de datos dinámica.

II-A. PASOS PARA GENERAR LA MATRIZ

El programa desarrollado permite que el usuario incluya una matriz dentro de un documento en formato csv. Las características que debe cumplir el archivo son:

1. Nombrarlo file.csv
2. Separar los elementos con comas
3. Cada fila de la matriz debe escribirse en una nueva fila en el archivo

A continuación en I se presenta un ejemplo del file.csv que debe contruirse.

TABLE I
FILE.CSV: EJEMPLO

20,10,1,3,4,7,8,4,23,56,844,24
3,12,13,56,1,6,7,98,46,34,21,36
7,4,5,9,30,3,700,6,45,2,93,16

II-A.1. SECUENCIAL: El proceso para calcular la matriz de covarianzas de forma secuencial está en un *script* en C. Este proceso imprime a pantalla la matriz origen, la matriz de los productos de las desviaciones y la matriz de covarianzas.

Algorithm 1: Matriz de Varianzas y Covarianzas: Método Secuencial

Result: Matriz de covarianzas de las variables de una matriz (M) de p filas y q columnas implementando un algoritmo secuencial

```
1 initialization;
2 M una matriz de datos de p filas y q columnas
3 Calcular la suma de los valores de cada una de las q
  variables de M.
4 while Mientras no se terminen las columnas do
5   while Mientras no se terminen las filas do
6     suma_j += M(i,j);
7   end
8 end
9 Procedemos a calcular los valores esperados por
  variable. while Mientras no se terminen las columnas
  do
10  E(j)=suma_j/p;
11 end
12 Procedemos a calcular la matriz de los productos de las
  desviaciones.
13 while Mientras no se terminen las columnas do
14   while Mientras no se terminen las columnas do
15     while Mientras no se terminen las filas do
16       restasX_i = M(i,j) - E(j); restasY_i = M(i,z) -
       E(z); M_multiXY(j,z) += restasX_i*restasY_i
17     end
18   end
19 end
20 Procedemos a calcular la matriz de las covarinzas.
21 while Mientras no se terminen las columnas do
22   while Mientras no se terminen las columnas do
23     M_covarianzas(j,z) = M_multiXY(i,j)/q
24   end
25 end
```

II-A.2. PARALELO:

OPEN-MP: El cómputo en paralelo tiene dos arquitecturas principales: *shared memory* y *clusters*. En este proyecto, la matriz es realizada con una arquitectura de *shared memory*: Open-MP.

Es la API estándar para la escritura de la memoria compartida de las aplicaciones en paralelo en C, C++ y Fortran.

En qué consiste?

- Directivas del compilador
- Subrutinas/funciones de tiempo de ejecución
- Variables de entorno

El algoritmo construido tiene algunas características que el usuario debe distinguir:

1. *pragma omp parallel* se utiliza para bifurcar *threads* adicionales para llevar a cabo el trabajo. El *thread* original se denota como *thread* principal con ID 0.

a) *private*: los datos de la región paralela son

privados para cada *thread*, lo que significa que cada *thread* tendrá una copia local que la usa como variable temporal. Una variable privada no es inicializada y tampoco se mantiene fuera de la región paralela. Por definición, los contadores de iteraciones en Open-MP es privado.

b) *shared*: los datos de la región paralela son compartidos, lo que significa que son visibles y accesibles por todos los *threads*. Por definición, todas las variables que trabajan en la región paralela son compartidas excepto los contadores de iteraciones.

2. *schedule*: Las iteraciones se asignan a los *thread* de acuerdo con el método de programación definido en el *schedule*, en particular en este proyecto se asigno un *schedule static*. En este último todos los *thread* son asignados antes de que las iteraciones sean ejecutadas.

3. *chunk*: Siguiendo la elección de *schedule static*, el loop es dividido en cada *thread* en partes iguales.

En este proyecto es necesario que el usuario considere que el número de *threads* debe ser divisor del número de columnas.

III. DISCUSIÓN

El cálculo de la matriz de covarianzas a través de Open-MP se realizó con éxito. No obstante, el uso del algoritmo construido requiere que el usuario elimine las instrucciones de impresión para hacer más eficiente el código. Por otra parte, se sugiere incorporar la comparación entre paralelismo con memoria compartida y métodos como MPI que posibilitan ejecutar paralelismo a través de la formación de clusters.

REFERENCIAS

- [1] Chilson J. Parallel Computation of High Dimensional Robust Correlation and Covariance Matrices, University of British Columbia, Thesis, March, 2004.
- [2] Open-MP Loop Scheduling, Ronald W Green (Intel). September 4, 2012

Algorithm 2: Matriz de Varianzas y Covarianzas: Mtodo Paralelo

Result: Matriz de covarianzas de las variables de una matriz (M) de p filas y q columnas utilizando paralelismo de cores por medio de la implementacin de multithreading processing a travs de la extensi3n OpenMP

```
1 initialization;
2 M una matriz de datos de p filas y q columnas
3 nthreads : N3mero de threads elegidos
4 chunk: N3mero de elementos consecutivos que se le
  asignan a cada thread dentro de un for
5 Obtenemos el n3mero de threads que eligi3 el usuario.
6 Calculamos el chunk en funci3n al nthreads y a la
  cantidad de campos (q) que tiene la matriz fuente (M).
7 Inicializamos la directiva parallel de pragma omp
  estableciendo las variables que ser3n compartidas para
  todos los threads y las que mantendr3n valores en
  privado.
8 Iniciamos un for schedule con la cantidad de calculada
  de chunk y procedemos a calcular la suma de los valores
  de cada una de las q variables de M.
9 while Mientras no se terminen las columnas do
10 |   while Mientras no se terminen las filas do
11 | |   suma_j +=M(i,j);
12 |   end
13 end
14 Iniciamos un for schedule con la cantidad de calculada
  de chunk y procedemos a calcular los valores esperados
  por variable. while Mientras no se terminen las
columnas do
15 |   E(j)=suma_j/p;
16 end
17 Iniciamos un for schedule con la cantidad de calculada
  de chunk y procedemos a calcular la matriz de los
  productos de las desviaciones.
18 while Mientras no se terminen las columnas do
19 |   while Mientras no se terminen las columnas do
20 | |   while Mientras no se terminen las filas do
21 | | |   restasX_i = M(i,j) - E(j); restasY_i = M(i,z) -
  | | |   E(z); M_multiXY(j,z) += restasX_i*restasY_i
22 | |   end
23 |   end
24 end
25 Iniciamos un for schedule con la cantidad de calculada
  de chunk y procedemos a calcular la matriz de las
  covarianzas.
26 while Mientras no se terminen las columnas do
27 |   while Mientras no se terminen las columnas do
28 | |   M_covarianzas(j,z) = M_multiXY(i,j)/q
29 |   end
30 end
```
