

DATOS MASIVOS II

DESCOMPOSICIÓN CUR

Blanca Vázquez y Gibran Fuentes-Pineda

Octubre 2020

- Surge en los años 1997-1998 por Stewart Stewart y Goreinov-Tyrtysnikov-Zamarashkin.
- Este método se formaliza en el 2004 por Drineas-Kannan-Mahoney
- Surge como una alternativa a las descomposiciones deterministas para grandes volúmenes de datos.

¿POR QUÉ SURGE LA DESCOMPOSICIÓN CUR?

Bolsa de palabras

	w_1	w_2	w_3	w_4	w_5	...	w_n
doc_1	0	0	1	0	0	0	1
doc_2	0	1	0	0	0	0	0
doc_3	0	0	0	1	0	0	1
doc_4	1	0	0	0	0	0	0
....
doc_n	0	0	0	0	1	0	0

¡El reto de las matrices dispersas!

- Es un método de reducción, que se expresa en términos de un sub-conjunto de las variables originales.
- Es una alternativa aleatoria para SVD
- A diferencia de SVD que genera una descomposición exacta, CUR realiza una aproximación.

DESCOMPOSICIÓN DE CUR

Dada una matriz M de $m \times n$, se define la descomposición de CUR de M como:

$$M \approx CUR$$

Donde:

C es una selección aleatoria de r columnas de M , y forma la matriz de $m \times r$

R es una selección aleatoria de r filas de M , y forma la matriz de $r \times n$

U es una matriz que se construye a partir de C y R

Para construir U, deben seguirse estos pasos:

- Generar la matriz W de $r \times r$ que es la intersección de las columnas escogidas para C y de las filas escogidas para R.
- Calcular la SVD para W, es decir: $W = X\Sigma Y^T$
- Calcular la pseudo-inversa de Moore-Penrose (Σ^+) de la diagonal de la matriz Σ
- Calcular $U = Y(\Sigma^T)^2 X^T$

La descomposición CUR de M , necesita:

- Una matriz C que está compuesta de algunas columnas de M
- Una matriz R que está compuesta de algunas filas de M

La descomposición CUR de M , necesita:

- La selección se basa en la importancia de las columnas y de las filas
- Para medir la importancia se usa el cuadrado de la norma de Frobenius

$$f = \sum_{i,j} m_{i,j}^2$$

Cada vez que seleccionamos una fila, la probabilidad p_i de que la fila i sea seleccionada es:

$$\Sigma_j m_{i,j}^2 / f$$

Cada vez que seleccionamos una columna, la probabilidad q_j de que la columna j sea seleccionada es:

$$\Sigma_i m_{i,j}^2 / f$$

Tabla 1: Matriz M^1

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

¹Ejemplo tomado de Jure Leskovec, 2011.

Tabla 2: Matriz M

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

- La suma de los cuadrados de los elementos de M = 243

Tabla 3: Matriz M

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

- Cada una de las columnas de M, A y S tiene una norma cuadrada de Frobenius de $1^2 + 3^2 + 4^2 + 5^2 = 51$
- La probabilidad de cada una de las columnas para ser seleccionada es: $51/243 = 0.210$

Tabla 4: Matriz M

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

- Cada una de las columnas de C y T tiene una norma cuadrada de Frobenius de $4^2 + 5^2 + 2^2 = 45$
- La probabilidad de cada una de las columnas para ser seleccionada es: $45/243 = 0.185$

Tabla 5: Matriz M

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

- La probabilidad de las columnas M, A y S es 0.210 y de las columnas C y T es 0.185

Tabla 6: Matriz M

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

- Para la fila 1 (Joe) $1^2 + 1^2 + 1^2 = 3$
- La probabilidad de que la fila sea seleccionada es $3/243 = 0.012$

Tabla 7: Matriz M

	Matrix	Alien	Star Wars	Casablanca	Titanic	Frobenius	Prob
Joe	1	1	1	0	0	3	0.12
Jim	3	3	3	0	0	27	0.111
John	4	4	4	0	0	48	0.198
Jack	5	5	5	0	0	75	0.309
Jill	0	0	0	4	4	32	0.132
Jenny	0	0	0	5	5	50	0.206
Jane	0	0	0	2	2	8	0.033
Frobenius	51	51	51	45	45		
Prob	0.210	0.210	0.210	0.185	0.185		

- Supongamos que $r = 2$ (no. columnas y filas a seleccionar)
- Columnas más importantes: Alien y Casablanca
- Cada columna seleccionada debe escalarse: dividiendo sus elementos entre $\sqrt{rq_j}$

$$\textit{Alien} = \begin{bmatrix} 1 \\ 3 \\ 4 \\ 5 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \textit{Casablanca} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 4 \\ 5 \\ 2 \end{bmatrix}$$

$$Alien = \begin{bmatrix} 1 \\ 3 \\ 4 \\ 5 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Calculamos $\sqrt{rq_2}$

Donde $q_2 = 0.210(51/243)$, $r = 2$

Por lo tanto: $\sqrt{rq_2} = \sqrt{2 * 0.210} = 0.648$

$$Alien = \begin{bmatrix} 1/0.648 \\ 3/0.648 \\ 4/0.648 \\ 5/0.648 \\ 0/0.648 \\ 0/0.648 \\ 0/0.648 \end{bmatrix}, Alien_{esc} = \begin{bmatrix} 1.54 \\ 4.63 \\ 6.17 \\ 7.72 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$Alien_{esc}$ es la primera columna de la matriz C .

$$\text{Casablanca} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 4 \\ 5 \\ 2 \end{bmatrix}$$

Calculamos $\sqrt{rq_4}$

Donde $q_4 = 0.185(45/243)$, $r = 2$

Por lo tanto: $\sqrt{rq_4} = \sqrt{2 * 0.185} = 0.608$

$$Casablanca = \begin{bmatrix} 0/0.608 \\ 0/0.608 \\ 0/0.608 \\ 0/0.608 \\ 4/0.608 \\ 5/0.608 \\ 2/0.608 \end{bmatrix}, Casablanca_{esc} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 6.58 \\ 8.22 \\ 3.29 \end{bmatrix}$$

$Casablanca_{esc}$ es la segunda columna de la matriz C .

1ER RESULTADO: MATRIZ C

$$C = \begin{bmatrix} 1.54 & 0 \\ 4.63 & 0 \\ 6.17 & 0 \\ 7.72 & 0 \\ 0 & 6.58 \\ 0 & 8.22 \\ 0 & 3.29 \end{bmatrix}$$

- Supongamos que $r = 2$ (no. columnas y filas a seleccionar)
- Filas más importantes: Jack y Jenny
- Cada fila seleccionada debe escalarse: dividiendo sus elementos entre $\sqrt{rp_i}$

$$Jenny = \begin{bmatrix} 0 & 0 & 0 & 5 & 5 \end{bmatrix}$$

$$Jack = \begin{bmatrix} 5 & 5 & 5 & 0 & 0 \end{bmatrix}$$

$$Jenny = \begin{bmatrix} 0 & 0 & 0 & 5 & 5 \end{bmatrix}$$

Calculamos $\sqrt{rp_6}$

Donde $p_6 = 0.206(50/243)$, $r = 2$

Por lo tanto: $\sqrt{rp_6} = \sqrt{2 * 0.206} = 0.642$

$$Jenny = \begin{bmatrix} 0/0.642 & 0/0.642 & 0/0.642 & 5/0.642 & 5/0.642 \end{bmatrix}$$

$$Jenny_{esc} = \begin{bmatrix} 0 & 0 & 0 & 7.79 & 7.79 \end{bmatrix}$$

$Jenny_{esc}$ es la primera fila de la matriz R.

$$Jack = \begin{bmatrix} 5 & 5 & 5 & 0 & 0 \end{bmatrix}$$

Calculamos $\sqrt{rp_4}$

Donde $p_4 = 0.309(75/243)$, $r = 2$

Por lo tanto: $\sqrt{rp_4} = \sqrt{2 * 0.309} = 0.786$

$$Jack = \begin{bmatrix} 5/0.786 & 5/0.786 & 5/0.786 & 0/0.786 & 0/0.786 \end{bmatrix}$$

$$Jack_{esc} = \begin{bmatrix} 6.36 & 6.36 & 6.36 & 0 & 0 \end{bmatrix}$$

$Jack_{esc}$ es la segunda fila de la matriz R.

$$R = \begin{bmatrix} 0 & 0 & 0 & 7.79 & 7.79 \\ 6.36 & 6.36 & 6.36 & 0 & 0 \end{bmatrix}$$

CONSTRUYENDO LA MATRIZ U

- U es una matriz de $r \times r$
- Para construir U, es necesario construir la matriz W.
- W es la intersección resultante entre filas y columnas que se usaron para construir C y R.
- Una vez construida W, se calcula su SVD, es decir:
$$W = X\Sigma Y^T$$
- Cada elemento de la matriz Σ debe reemplazarse usando la pseudoinversa de Moore-Penrose
- Para obtener U, se calcula

$$U = Y(\Sigma^+)^2 X^T$$

CONSTRUYENDO LA MATRIZ U

		Columnas para construir C				
		Matrix	Alien	Star Wars	Casa	Titanic
Filas para construir R	Jenny	0	0	0	5	5
	Jack	5	5	5	0	0

CONSTRUYENDO LA MATRIZ U

		Columnas para construir C				
		Matrix	Alien	Star Wars	Casa	Titanic
Filas para construir R	Jenny	0	0	0	5	5
	Jack	5	5	5	0	0

$$W = \begin{bmatrix} 0 & 5 \\ 5 & 0 \end{bmatrix}$$

El siguiente paso es calcular la SVD de W

$$W = X\Sigma Y^T$$

La SVD de $W = X\Sigma Y^T$:

$$W = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Usando la matriz Σ , calculamos la pseudoinversa de Moore-Penrose Σ^+

$$\Sigma = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$$

- Cada elemento de la diagonal de Σ se reemplaza por $1/\sigma_i$
- Si el elemento de la diagonal es 0, se deja igual.

Usando la matriz Σ , calculamos la pseudoinversa de Moore-Penrose Σ^+

$$\Sigma = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}, \Sigma^+ = \begin{bmatrix} 1/5 & 0 \\ 0 & 1/5 \end{bmatrix}$$

Para calcular U, debemos calcular:

$$U = Y(\Sigma^+)^2 X^T$$

$$W = (X) \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} (\Sigma^+) \begin{bmatrix} 1/5 & 0 \\ 0 & 1/5 \end{bmatrix} (Y^T) \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Para calcular U, debemos calcular:

$$U = Y(\Sigma^+)^2 X^T$$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1/5 & 0 \\ 0 & 1/5 \end{bmatrix}^2 \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Recordemos, X e Y son simétricas, por lo que son sus propias transpuestas.

CONSTRUYENDO LA MATRIZ U

Para calcular U, debemos calcular:

$$U = Y(\Sigma^+)^2 X^T$$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1/5 & 0 \\ 0 & 1/5 \end{bmatrix}^2 \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1/25 & 0 \\ 0 & 1/25 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 1/25 & 0 \\ 0 & 1/25 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 0 & 1/25 \\ 1/25 & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 0 & 1/25 \\ 1/25 & 0 \end{bmatrix}$$

DESCOMPOSICIÓN CUR DE M

Dada una matriz M de $m \times n$, se define la descomposición de CUR de M como:

$$M \approx CUR$$

$$M = \begin{bmatrix} 1.54 & 0 \\ 4.63 & 0 \\ 6.17 & 0 \\ 7.72 & 0 \\ 0 & 6.58 \\ 0 & 8.22 \\ 0 & 3.29 \end{bmatrix} \begin{bmatrix} 0 & 1/25 \\ 1/25 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 7.79 & 7.79 \\ 6.36 & 6.36 & 6.36 & 0 & 0 \end{bmatrix}$$

- La descomposición CUR es un método de aproximación de la matriz original, construida a partir de un subconjunto de valores originales.
- Es un método de reducción de variables
- Es un método alternativo para SVD para matrices dispersas.