

Diseño y Análisis de Algoritmos ICI-522

Sergio Hernández.
PhD computer science

Departamento de Computación e Informática
Universidad Católica del Maule.
`shernandez@ucm.cl`

Sistemas Recomendadores ¹

- Los sistemas recomendadores son herramientas que sugieren items a un usuario.
- Los sitios como YouTube utilizan estos software para encontrar items (videos, aplicaciones, libros, etc..) basado en los intereses del usuario.



Figure: Google Playstore

¹https:

//developers.google.com/machine-learning/recommendation/overview

Filtrado Colaborativo



Figure: Filtrado colaborativo para una matriz de intereses de usuario $A \in \mathbb{R}^{m \times n}$ usando factorización de matrices $U_k \in \mathbb{R}^{m \times k}$ y $V_k \in \mathbb{R}^{n \times k}$, tal que $A \approx U_k \times V_k^T$

Filtrado Colaborativo

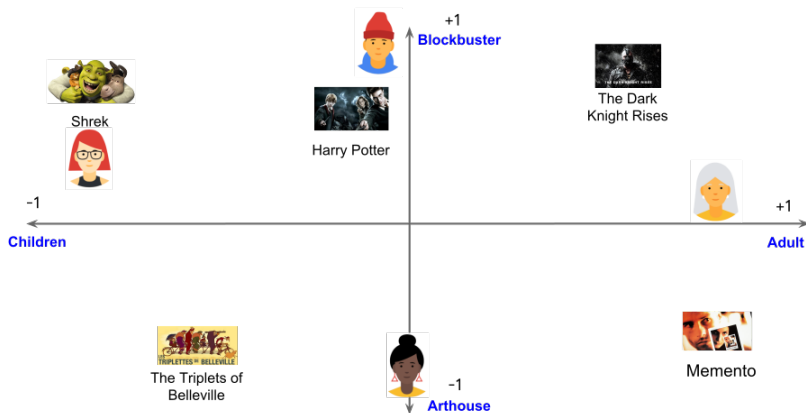


Figure: La idea es encontrar un conjunto de k factores latentes que caractericen a la relación entre los usuarios y los items (embeddings)

Descomposición por Valores Singulares

Sea $A \in \mathbb{R}^{m \times n}$ una matriz de intereses de usuario, $U \in \mathbb{R}^{m \times m}$, $S \in \mathbb{R}^{m \times n}$ y $V \in \mathbb{R}^{n \times n}$

SVD

$$A = U \times S \times V^T$$

Descomposición por Valores Singulares

Sea $A \in \mathbb{R}^{m \times n}$ una matriz de intereses de usuario, $U \in \mathbb{R}^{m \times m}$, $S \in \mathbb{R}^{m \times n}$ y $V \in \mathbb{R}^{n \times n}$

SVD

$$A = U \times S \times V^T$$

Propiedades

Sean $\sigma_1, \dots, \sigma_r$ una secuencia de valores no-negativos en orden decreciente y $r = \text{MIN}(m, n)$

$$U^T \times U = \mathbb{I}_m, \quad S = \begin{pmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \cdots & & & \\ 0 & \cdots & \sigma_r & 0 \end{pmatrix}$$

$$V^T \times V = \mathbb{I}_n$$

SVD Truncada

Ahora consideramos una aproximación de bajo rango utilizando $k < r$ factores latentes:

SVD truncada

$$A \approx U_k \times S_k \times V_k^T$$

SVD Truncada

Ahora consideramos una aproximación de bajo rango utilizando $k < r$ factores latentes:

SVD truncada

$$A \approx U_k \times S_k \times V_k^T$$

Propiedades

$$\tilde{A} = \left[\begin{array}{c|c|c} \mathbf{u}_1 & \cdots & \mathbf{u}_k \end{array} \right] \left[\begin{array}{ccc} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{array} \right] \left[\begin{array}{c} \hline \mathbf{v}_1^T \\ \vdots \\ \hline \mathbf{v}_k^T \end{array} \right]$$

Figure: Aproximación de rango bajo

Análisis de Complejidad

- Para matrices A de tamaño $m \times n$ el costo computacional de obtener la descomposición es $\mathcal{O}(mnr)$
- Para el problema de valores propios, podemos calcular

$$T = \left[\begin{array}{c|c} T_1 & \beta \\ \hline \beta & T_2 \end{array} \right] \quad (1)$$

- La relación de recurrencia está dada por $T(n) = 2T\left(\frac{n}{2}\right) + r^2$
- Por lo tanto, el costo del algoritmo divide y vencerás es $\mathcal{O}(r^2)$

Cuppen, J.J.M. A divide and conquer method for the symmetric tridiagonal eigenproblem. Numer. Math. 36, 177-195 (1980).
<https://doi.org/10.1007/BF01396757>

Valores Propios de Matrices

La relación entre los valores singulares y los valores propios está dada por la matriz de Gram.

```

1 >>> import numpy as np
2 >>> from scipy import linalg
3 >>> A=np.array([[1,4],[5,2]])
4 >>> U,S,V = linalg.svd(A,full_matrices=False,
5     lapack_driver='gesdd')
6 >>> D=np.diag(S)
7 >>> np.matmul(D,D.T)
8 array([[37.31782106,  0.          ],
9        [ 0.          ,  8.68217894]])
10 >>> np.linalg.eig(np.matmul(A,A.T))[0]
11 array([ 8.68217894, 37.31782106])

```

Kevin P. Murphy. 2012. Machine Learning: A Probabilistic Perspective. The MIT Press (pp. 392).

SVD para big data

```
1 >>> import time
2 >>> import numpy as np
3 >>> from scipy import linalg
4
5 >>> A = np.random.random((200000, 100))
6
7 >>> t1=time.time()
8 >>> U,S,V = linalg.svd(A,full_matrices=False,lapack_driver='
    gesdd')
9 >>> print(time.time()-t1)
10 3.5332744121551514
```

SVD para big data

```
1 >>> import dask.array as da
2 >>> import dask
3
4 >>> X = da.random.random((200000, 100), chunks=(50000, 100))
   .persist()
5
6 >>> t2=time.time()
7 >>> U, S, V = da.linalg.svd(X)
8 >>> print(time.time()-t2)
9 0.0014061927795410156
```

```
1 >>> dask.visualize(S)
```

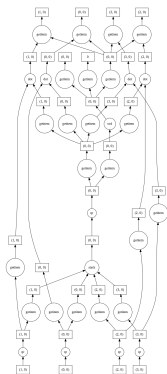


Figure: Grafo de cómputo usando Dask