

# **BÚSQUEDA Y MINERÍA DE INFORMACIÓN**

## **PRÁCTICA 4**

**Sistemas de recomendación y análisis de redes  
sociales**

Tomas Higuera Viso  
Miguel Antonio Núñez Valle  
Pareja 25

## Pregunta 1

- Hemos realizado el ejercicio de **estructuras de datos y recomendación simple**.
- Hemos implementado el **filtrado colaborativo kNN**.
- Hemos implementado la **recomendación basada en contenido**.
- En la ampliación de algoritmos hemos implementado tanto **vecinos próximos orientados a ítem** como el algoritmo basado en **contenido por vecinos próximos**.
- Hemos creado las tablas con las métricas y tiempos de ejecución adjuntadas más abajo, aunque no hemos podido ejecutar nuestra implementación de algunos recomendadores con los conjuntos de datos grandes.
- Hemos creado las dos **redes sociales simuladas**, así como la representación y análisis de todas las topologías.
- Hemos implementado las **métricas de análisis de topologías de red**.

## Pregunta 2

En este ejercicio analizaremos los resultados obtenidos con los distintos recomendadores. Como hemos explicado no hemos podido ejecutar los recomendadores con todos los conjuntos de datos, ya que nuestro ordenador tardaba mucho tiempo en realizar las pruebas con los mismos. Los datos que no hayamos podido obtener los dejaremos representados en la tabla con NaN.

**Tabla conjunto de datos: *Toy dataset***

	<b>Rmse</b>	<b>Precision</b>	<b>Recall</b>	<b>Tiempo</b>
<b>Random</b>	2.486	0.04	0.4	3ms
<b>Majority</b>	2	0.04	0.4	0ms
<b>Average</b>	0.5	0.02	0.2	1ms
<b>User-based kNN</b>	2.226	0.04	0.4	11ms
<b>Normalized user-based kNN</b>	0.508	0.02	0.2	0ms
<b>Item-based NN</b>	1.553	0.04	0.4	2ms
<b>Centroid-based</b>	2.369	0.04	0.4	1ms

<b>Item-based NN (Jaccard)</b>	10.794	0.04	0.4	19ms
--------------------------------	--------	------	-----	------

**Tabla conjunto de datos: *MovieLens latest-small datatest***

	<b>Rmse</b>	<b>Precision</b>	<b>Recall</b>	<b>Tiempo</b>
<b>Random</b>	2.703	0.003	0.001	1.293s
<b>Majority</b>	519.562	0.155	0.780	1.29s
<b>Average</b>	1.092	8.197E-4	3.447E-4	1.394s
<b>User-based kNN</b>	22.237	0.258	0.16	47.14s
<b>Normalized user-based kNN</b>	0.879	0.012	0.008	42.547s
<b>Item-based NN</b>	NaN	NaN	NaN	NaN
<b>Centroid-based</b>	3.844	0.063	0.031	38.639s
<b>Item-based NN (Jaccard)</b>	NaN	NaN	NaN	NaN

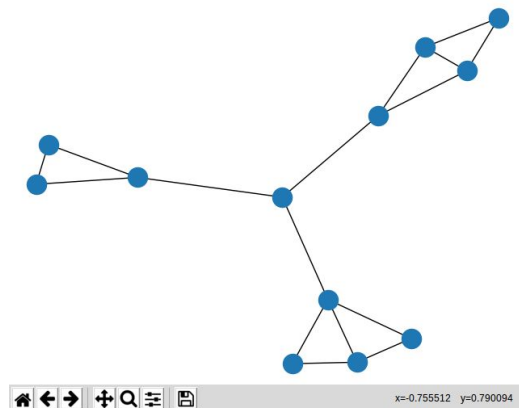
Tras realizar las tablas podemos observar que con conjuntos de datos pequeños (la primera tabla), los resultados que obtenemos con nuestros recomendadores no son muy buenos, sin embargo cuando trabajamos con conjuntos de datos mayores (segunda tabla), los resultados que obtenemos son considerablemente mejores.

Como hemos dicho al principio no hemos podido ejecutar los recomendadores basados en ítem, ya que los tiempos de ejecución eran muy grandes.

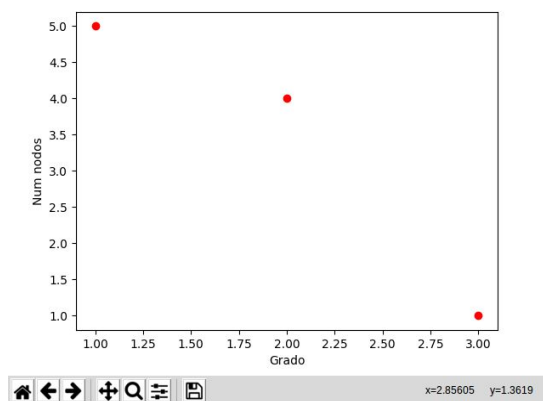
### Pregunta 3

A continuación representaremos las diferentes topologías de red que se nos proporcionan, así como las generadas por nosotros.

#### *small1.csv*

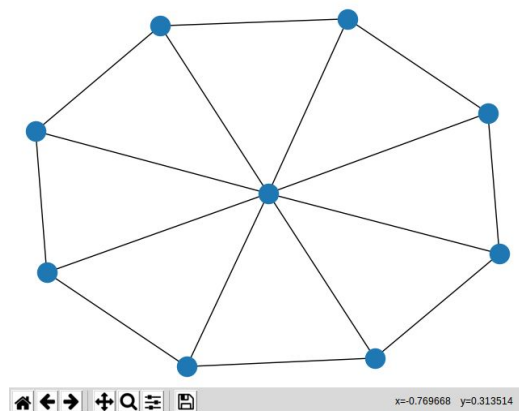


Representación de la topología

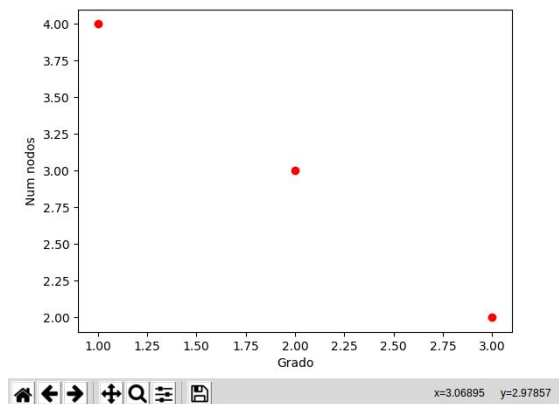


Grafica con x = grado; y = num\_nodos

#### *small2.csv*

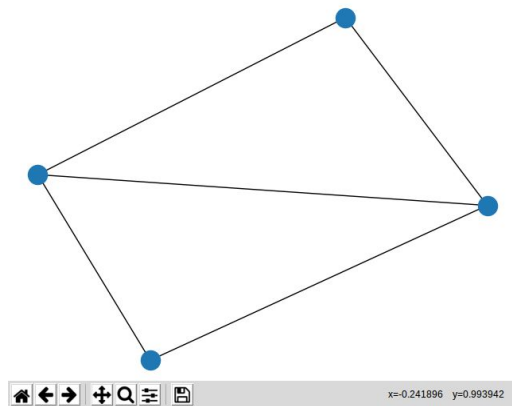


Representación de la topología

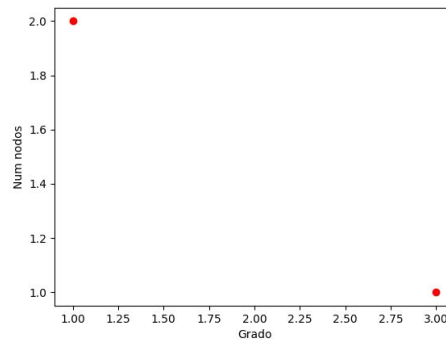


Grafica con x = grado; y = num\_nodos

### ***small3.csv***

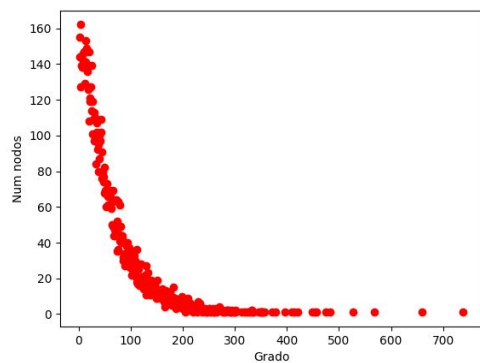


Representación de la topología

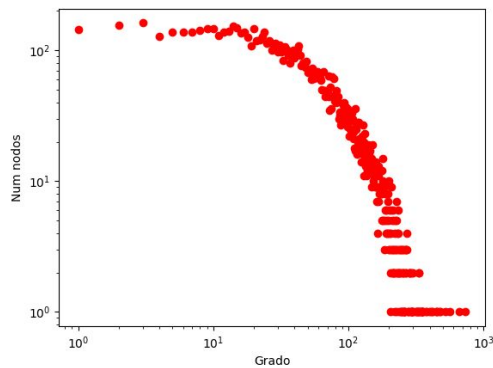


Grafica con x = grado; y = num\_nodos

### ***twitter.csv***

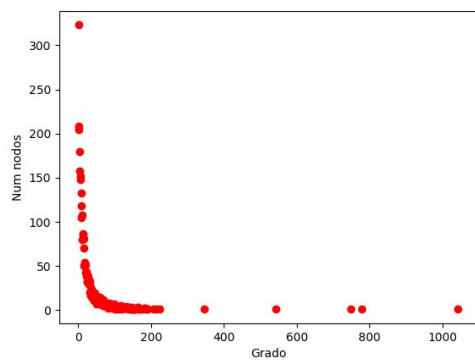


Grafica con x = grado; y = num\_nodos

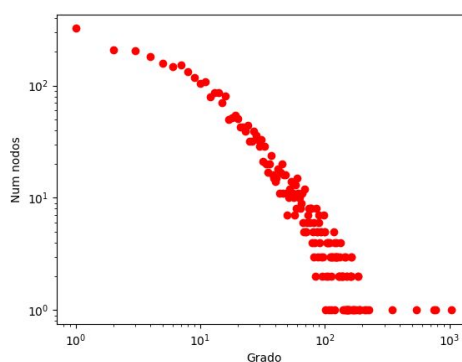


Grafica con x = grado; y = num\_nodos con escala log

### ***facebook.csv***

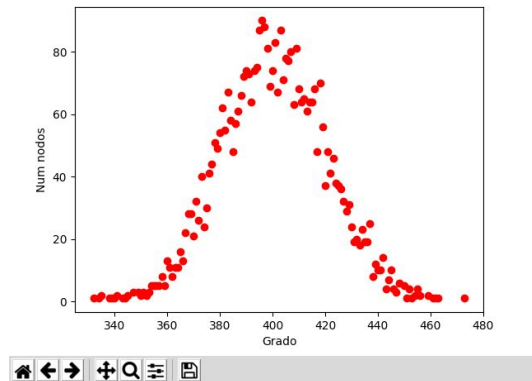


Grafica con x = grado; y = num\_nodos

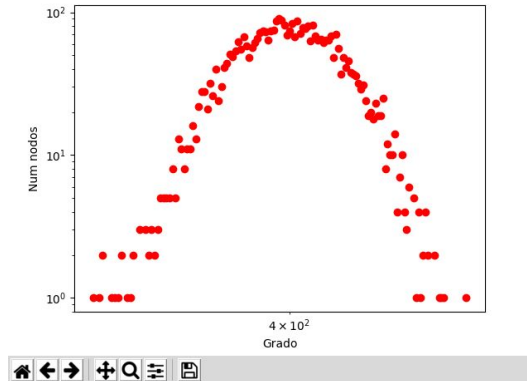


Grafica con x = grado; y = num\_nodos con escala log

### **erdos.csv**

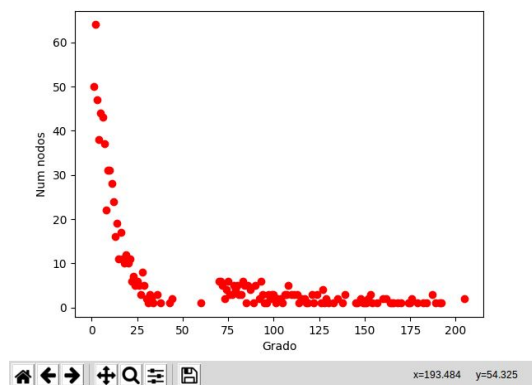


Grafica con x = grado; y = num\_nodos

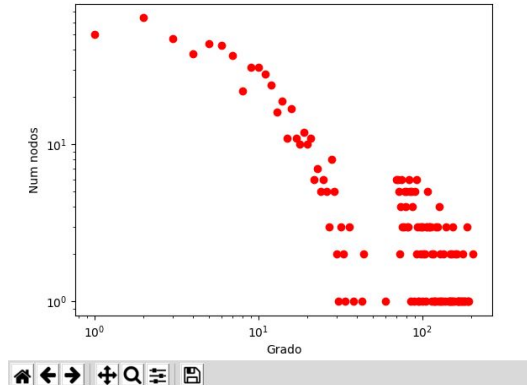


Grafica con x = grado; y = num\_nodos con escala log

### **barabasi.csv**



Grafica con x = grado; y = num\_nodos



Grafica con x = grado; y = num\_nodos con escala log

### **Conclusiones:**

Donde mejor se observa una distribución *power-law* es en las representaciones de twitter, facebook. En la distribución del conjunto de datos de barabasi, también puede apreciarse este tipo de distribución aunque menos clara que en los otros dos conjuntos.

Sobre la paradoja de amistad, la explicaremos en base a los conjuntos de datos de twitter y facebook, ya que son los conjuntos con más datos. Esta paradoja se explica en estos conjuntos diciendo que hay muchas personas con pocos amigos y pocas personas con muchos amigos, lo que al final se resumen con que es fácil ser amigo de personas con más amigos que con menos, ya que estas topologías representan una distribución *power-law*.