



# Generación de recomendaciones de ítems musicales basado en las valoraciones implícitas y las similitudes de los usuarios utilizando Hadoop para procesamientos masivos y escalables

Presentada por

**Mervyn Xavier Macías Martrus**

**Freddy Fernando de La Rosa**

**Octubre 2009**

# AGENDA

- Sistema de recomendaciones
- Algoritmos de recomendaciones
- Dataset
- Diseño
- Esquema map-reduce
- Implementación
- Alternativas de Ejecución
- Resultados
- Recomendaciones
- Conclusiones

# Sistemas de Recomendaciones

- Un sistema de recomendaciones es un tipo específico de filtro de información que sugiere a los usuarios ítems o productos concretos basándose en sus preferencias

# Ejemplos de Sistemas de Recomendaciones

- iTunes



- Pandora



- Last.fm



- Facebook



- Snooth



# Componentes

- Usuarios
- Perfiles
- Recomendadores
- Valoración

# Mecanismos de Retroalimentación

- **Explícita:** Cuando se aplica la recomendación explícita, el sistema otorga al usuario la oportunidad de calificar, dentro de un rango predefinido, los ítems que ha utilizado
- **Implícita:** El sistema obtiene retroalimentación implícita capturando la interacción del usuario sin que él lo note.

# Métodos de Recomendación

- Recomendación colaborativa
- Recomendación basada en contenido
- Recomendación basada en demografía
- Recomendación basada en utilidad
- Recomendación basada en conocimiento

# Recomendación Basada en Contenido

- En este caso las recomendaciones son hechas exclusivamente en base a los ítems que el usuario ha elegido en el pasado.

# Recomendación Colaborativa

- Cuando se aplica recomendación colaborativa, las recomendaciones son hechas exclusivamente en base a los usuarios con gustos similares.
- Hay sistemas que aplican más de un método de recomendación:
  - FAB
  - SELECT

# Desventajas de la Recomendación Basada en Contenido

- Se hace un análisis automático de los ítems, considerando atributos predefinidos pero se deja de lado otros atributos relevantes.
- Las sugerencias son hechas en función del historial de ítems elegidos por el usuario, por lo tanto el usuario sólo puede recibir recomendaciones que concuerden con su perfil

# Desventajas de la Recomendación Colaborativa

- Cuando un nuevo ítem se añade al sistema, como no ha sido evaluado por ningún usuario, no hay forma de recomendarlo.
- Si un usuario particular no se identifica con los gustos de ningún otro usuario del sistema, no es factible hallar vecinos cercanos, y por lo tanto hacer recomendaciones.
- Se requiere un mínimo de usuarios para elaborar las predicciones

## Ventajas de Recomendación Colaborativa sobre Basada en Contenido

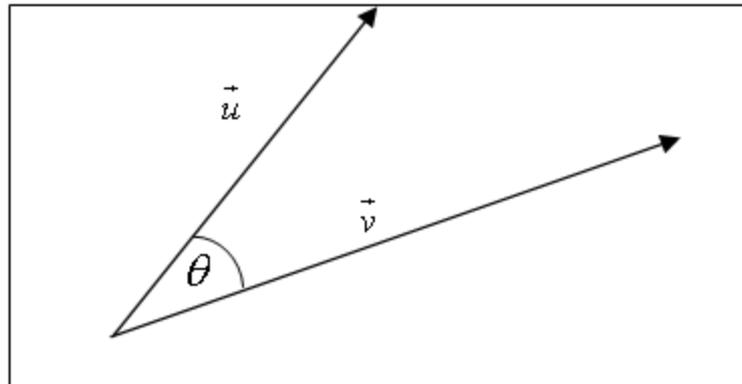
- En la recomendación colaborativa, los usuarios utilizan sus criterios para evaluar los ítems, cubriendo así más características de los ítems que evalúan.
- En recomendación colaborativa es posible que el usuario reciba recomendaciones que no se alineen a su perfil (pero sí a sus gustos).

# Algoritmos de Recomendación Colaborativa

- Los algoritmos de recomendación colaborativa se basan en el producto punto entre dos vectores y en las fórmulas de correlación .

# Producto Punto

$$\vec{u} \cdot \vec{v} = |\vec{u}| |\vec{v}| \cos \theta$$



# Producto Punto

- Cuando  $\theta$  es cero, los vectores apuntan en la misma dirección. Así, para valores de  $\theta$  cercanos a cero, los vectores tienden a apuntar en la misma dirección.

# Producto Punto

- Ejemplo:

	Ítems				
Usuario	A	B	C	D	F
1	4	5	0	0	0
2	5	6	0	5	0
3	5	0	4	0	8
4	0	0	5	0	0

# Correlación

- El coeficiente de correlación de **Pearson** mide la relación lineal entre dos variables cuantitativas.

$$\text{sim}(u, v) = \frac{\sum_{i=1}^m (r_{u,i} - r_u)(r_{v,i} - r_v)}{\sigma_u \sigma_v}$$

# Algoritmos de Recomendación Basado en Usuario

- Similitud

$$Sim(\vec{u}_i, \vec{u}_j) = \frac{u_i \cdot u_j}{\|u_i\| \|u_j\|}$$

- Recomendación

$$p_{a,i} = \frac{\sum_{k=1}^n r_{u_k,i} \cdot Sim(\vec{a}, \vec{u}_k)}{\sum_{k=1}^n Sim(\vec{a}, \vec{u}_k)}$$

# Algoritmos de Recomendación Basado en Usuario

- Con la aplicación de la fórmula anterior para cada par de usuarios del sistema se genera una matriz de similitud. Por ejemplo, si la primera fila de la matriz de similitud contiene la siguiente información:
- $u1: 0.5 \ u2 \ | \ 0.3 \ u3 \ | \ 0.2 \ u4 \ | \ 0.6 \ u7 \ | \ 0.9 \ u8$

# Algoritmos de Recomendación Basado en Ítems

- El principio es el mismo que en el caso de las recomendaciones basadas en usuario, la diferencia es que en este caso buscamos similitudes entre ítems en lugar de buscar similitudes entre usuarios.

# Análisis de Algoritmos

- En general, para un dataset con  $n$  usuarios y  $m$  ítems, para cada usuario se deben realizar  $n-1$  comparaciones, en total  $n(n-1)$ . En el peor de los casos cada comparación implica  $m$  operaciones. Así, el tiempo de ejecución es del orden de  $mn^2$ .

# Datasets

- Un dataset es una colección de datos presentados, por lo general, en forma tabular, de tal manera que cada columna representa una variable particular y cada fila corresponde a un miembro del dataset. La fila contiene los valores de las variables correspondientes a cada columna. Cada valor recibe el nombre de dato.

# Ejemplos de Dataset

- MusicBrainz es un dataset de música mantenido por la comunidad de usuarios. Contiene datos como el nombre del artista, título del álbum lanzado y la lista de pistas que aparecen en un álbum lanzado.
- El dataset de Audioscrobbler contiene una columna de usuarios, una de artistas y el número de veces que cada usuario escuchó a cada artista.

# Dataset de Audioscrobbler

	<b>id_user text</b>	<b>id_artist text</b>	<b>value double precis</b>
<b>825</b>	1000002	82	9
<b>826</b>	1000002	831	36
<b>827</b>	1000002	833	5
<b>828</b>	1000002	860	32
<b>829</b>	1000002	893	1
<b>830</b>	1000002	930	3
<b>831</b>	1000002	949	86
<b>832</b>	1000002	958	5
<b>833</b>	1000002	969	2
<b>834</b>	1000002	976	5
<b>835</b>	1000002	979	86
<b>836</b>	1000002	988	23
<b>837</b>	1000002	999	6
<b>838</b>	1000019	1000010	11
<b>839</b>	1000019	1000028	5
<b>840</b>	1000019	1000033	2
<b>841</b>	1000019	1000036	5
<b>842</b>	1000019	1000054	1

# Herramientas

- Map-Reduce

- Hadoop

- Mahout

- Amazon Web Services - EC2

- Amazon Web Services - S3



# Diseño

- Requerimiento de hardware
- Requerimiento de software
  1. Putty
  2. Eclipse
  3. Firefox
    - Amazon S3 Organizer
    - Elasticfox
  4. Vmware
  5. Filezilla



# Plug-in Firefox

The screenshot shows the AWS Management Console interface for the 'us-east-1' region. The 'Your Instances' tab is active, displaying a table of EC2 instances. The table columns include Reservation ID, Owner, Instance ID, Availability Zone, Architecture, Instance Profile, VPC, Subnet, State, and Public DNS. All instances shown are in a 'running' state.

Reservation ID	Owner	Instanc...	A...	AKI	ARI	VPC	Su...	State	Public DNS
r-193f0b70	8729809627...	i-9611f9fe	a...	aki...	ari...			running	ec2-75-101-23
r-893007e0	8729809627...	i-d05db...	a...	aki...	ari...			running	ec2-67-202-46
r-373f085e	8729809627...	i-d85cb...	a...	aki...	ari...			running	ec2-75-101-23
r-373f085e	8729809627...	i-da5cb...	a...	aki...	ari...			running	ec2-75-101-20
r-373f085e	8729809627...	i-dc5cb...	a...	aki...	ari...			running	ec2-75-101-23
r-373f085e	8729809627...	i-de5cb...	a...	aki...	ari...			running	ec2-75-101-20
r-373f085e	8729809627...	i-d05cb...	a...	aki...	ari...			running	ec2-75-101-18
r-373f085e	8729809627...	i-d25cb...	a...	aki...	ari...			running	ec2-67-202-8...
r-373f085e	8729809627...	i-d45cb...	a...	aki...	ari...			running	ec2-67-202-22

The screenshot shows the 'Your Groups' page in the AWS Management Console. A 'Grant New Permission' dialog box is open, allowing the user to add a new permission for the 'srmTesis-master' security group. The dialog is configured for an 'External' group with 'HTTP' protocol and 'TCP/IP' type, on port 5570. The 'Host/Network Details' section shows the 'Host' radio button selected, with the address '200.25.197.94/32'. Buttons for 'Get My Host Address' and 'Get My Network Range' are visible.

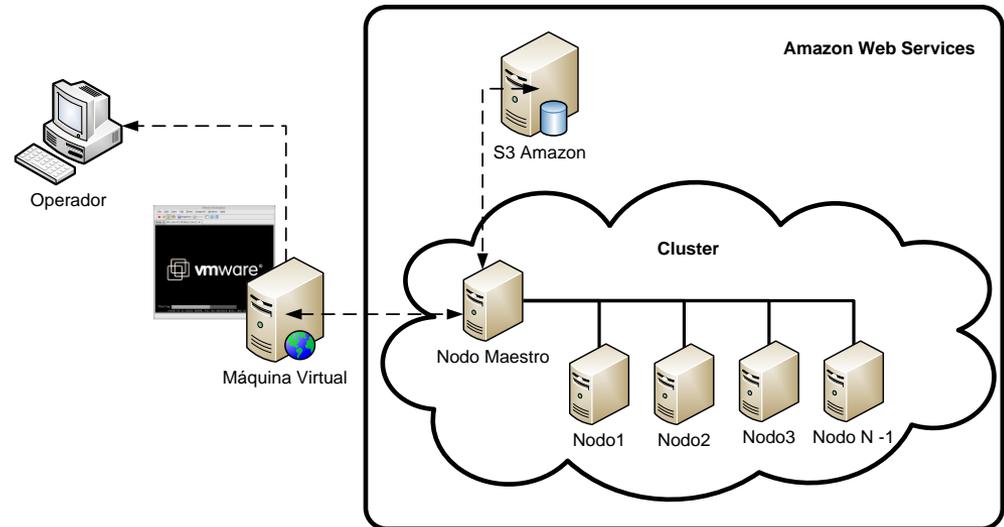
The screenshot shows two Firefox plugins. The 'S3 Firefox Organizer' plugin is on the left, displaying a 'Remote View' of a folder structure with items like 'blogbonsai', 'cdrs', 'distribuidos001', 'enxml', 'espol', and 'esxml'. The 'S3 Account Manager' plugin is on the right, showing 'S3 Account Preferences' for the account '8729-8096-2776'. It displays the 'Access Key' as '0G0RWXRKETPMX6VZ0AR2' and the 'Secret Key' as a series of dots. There are 'Save', 'Remove', and 'Clear' buttons, along with a hint: 'hint: Press 'Clear' button to add a new account'.

# Vmware Workstation

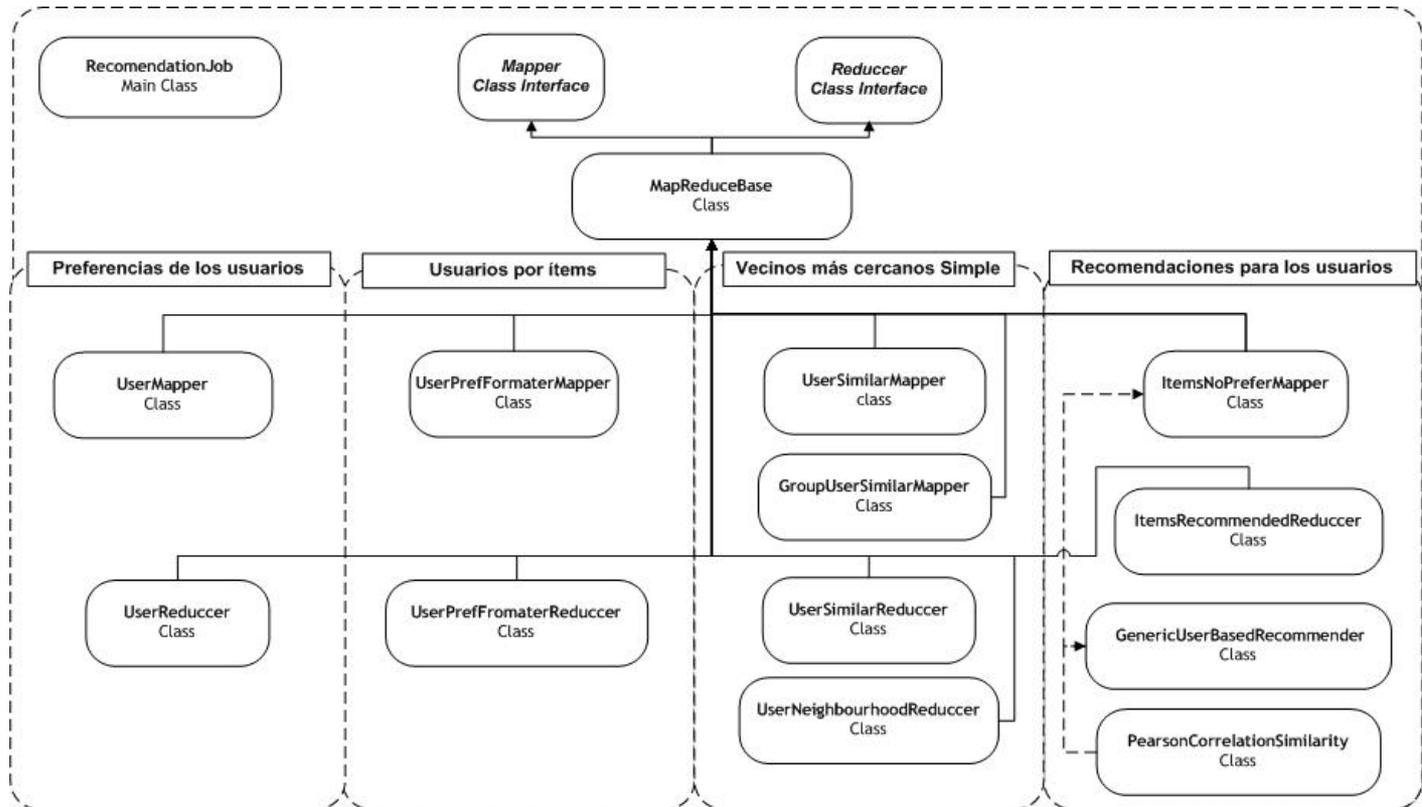
<a href="#">Name</a>	<a href="#">Creator</a>	<a href="#">Host OS(s)</a>	<a href="#">License</a>	<a href="#">Can boot an OS on another disk partition as guest</a>	<a href="#">USB</a>	<a href="#">GUI</a>	<a href="#">Live memory allocation</a>	<a href="#">3D acceleration</a>	<a href="#">Snapshot of running system</a>	<a href="#">Live migration</a>
<a href="#">Virtual Server 2005 R2</a>	<a href="#">Microsoft</a>	Windows 2003, XP	Proprietary					No		
<a href="#">Windows Virtual PC</a>	<a href="#">Microsoft</a>	<a href="#">Windows 7</a>	Proprietary	No	partially	Yes				
<a href="#">Virtual PC 2007</a>	<a href="#">Microsoft</a>	<a href="#">Windows Vista, XP</a>	Proprietary	No	No	Yes	No	No		
<a href="#">VMware Server</a>	<a href="#">VMware</a>	Windows, Linux	Proprietary			Yes		No	Yes	Yes
<a href="#">VMware Workstation 6.0</a>	<a href="#">VMware</a>	Windows, Linux	Proprietary	Yes	Yes	Yes	Yes	<a href="#">Experimental support for DirectX 8</a>	Yes	
<a href="#">VMware Player 2.0</a>	<a href="#">VMware</a>	Windows, Linux	Proprietary	No	Yes	Yes	Yes	<a href="#">Supported with VMGL</a>		
<a href="#">Sun xVM VirtualBox</a>	<a href="#">Sun Microsystems</a>	Windows, Linux, Mac OS X (Intel), Solaris, eComStation	GPL version 2	<a href="#">Partial (since version 1.4, but unsupported)</a>	<a href="#">Partial (for any host OS except Solaris)[8]</a>	Yes	Yes	<a href="#">OpenGL 2.0</a>	Yes (only on the closed-source edition)	Yes (only on the closed-source edition)

# Infraestructura Tecnológica

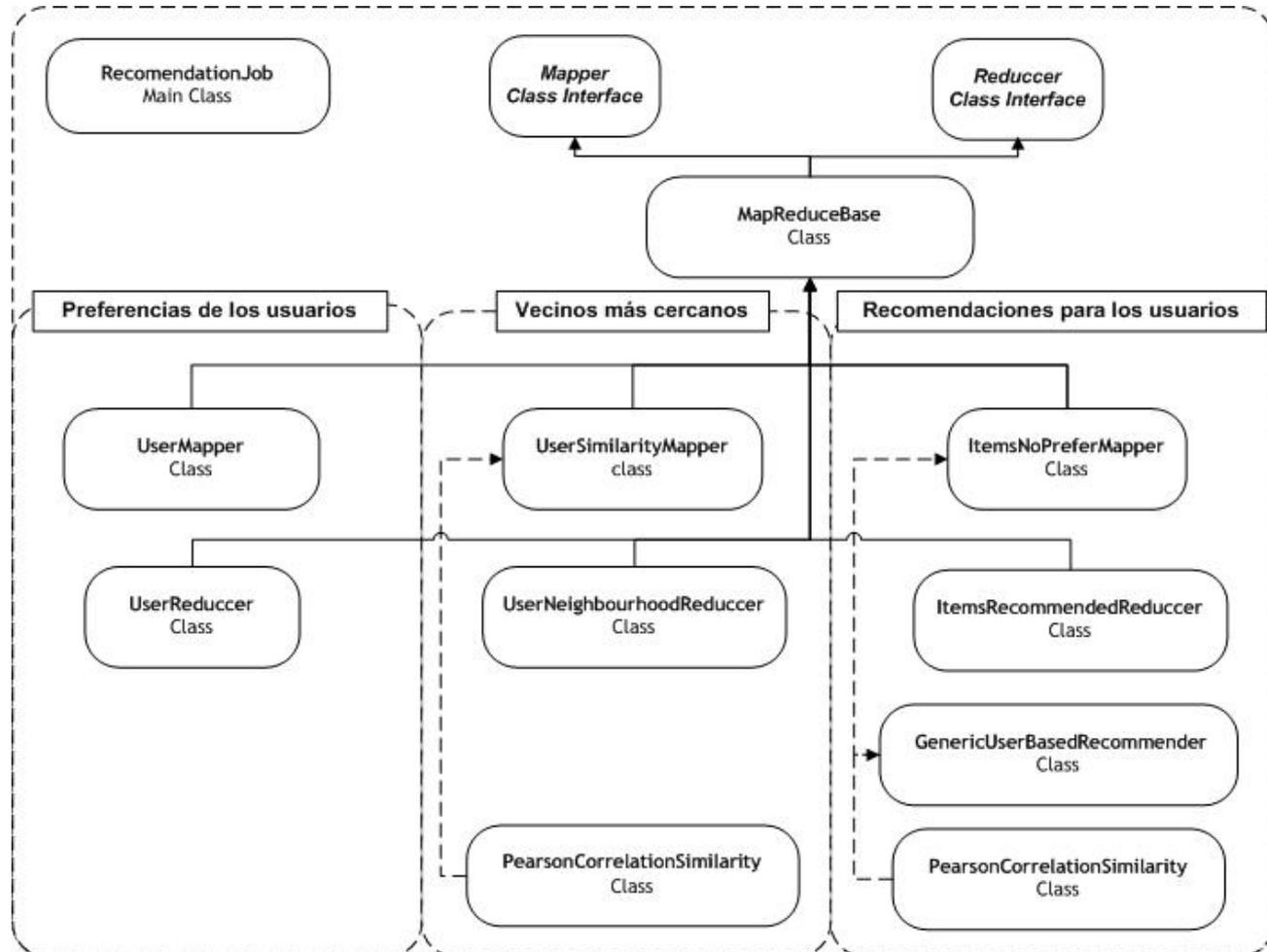
- Operador
- Máquina Virtual
- Nube EC2
- Nodo Maestro
- Nodos Esclavos
- Amazon S3



# Producto Punto



# Pearson

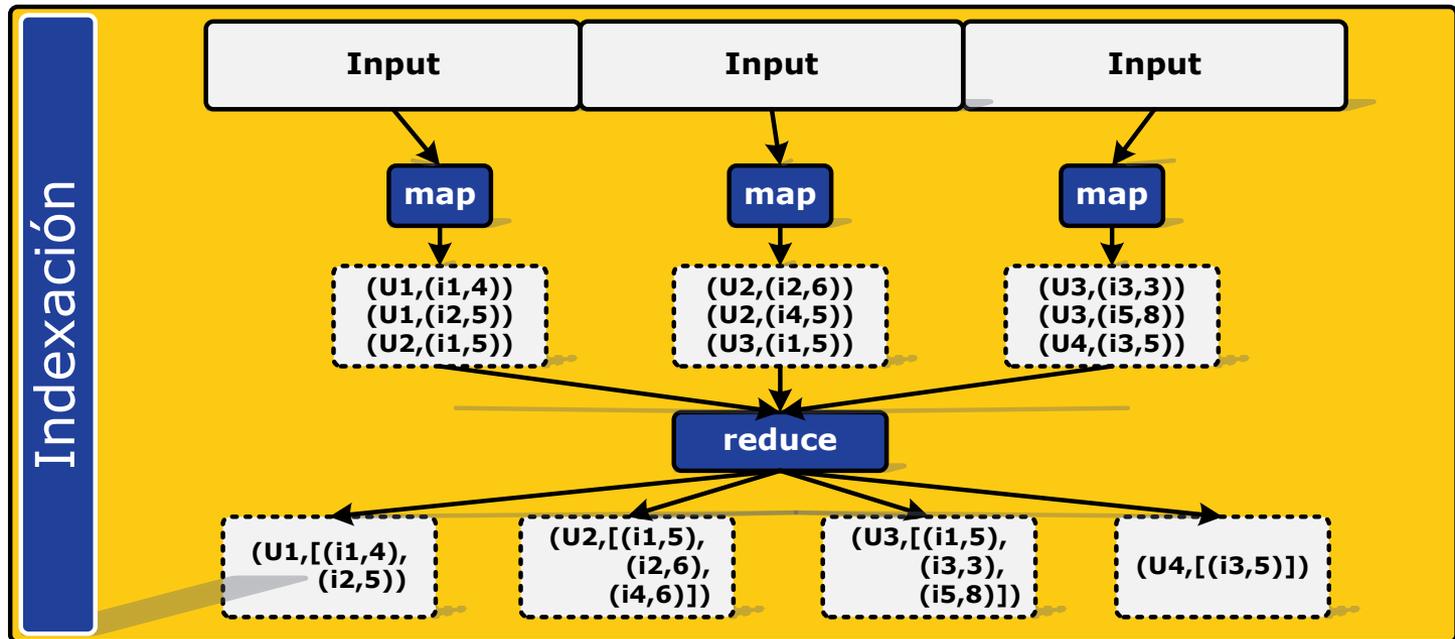


# Esquema map-reduce

- MapReduce es un modelo de programación para el procesamiento de grandes cantidades de información
- Función “**map**” procesa pares clave/valor para generar un conjunto de pares intermedios clave/valor, y una función “**reduce**” que agrupa todos los valores intermedios asociados con la misma clave

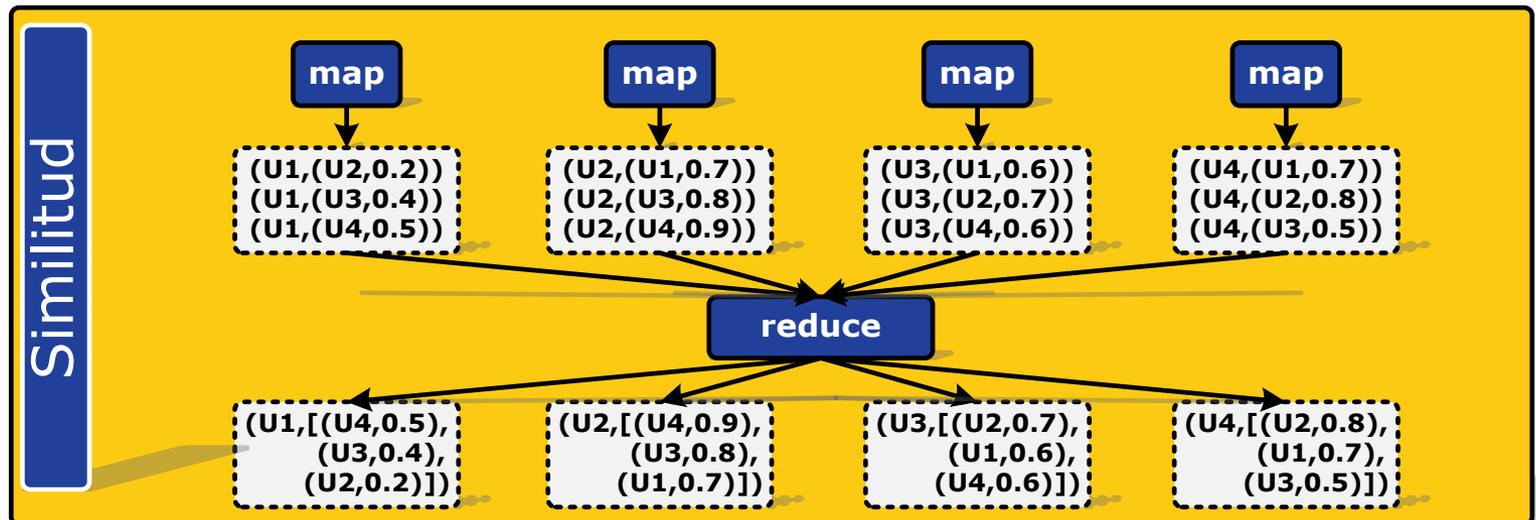
# Esquema map-reduce: Indexación

- Pre-procesamiento de datos



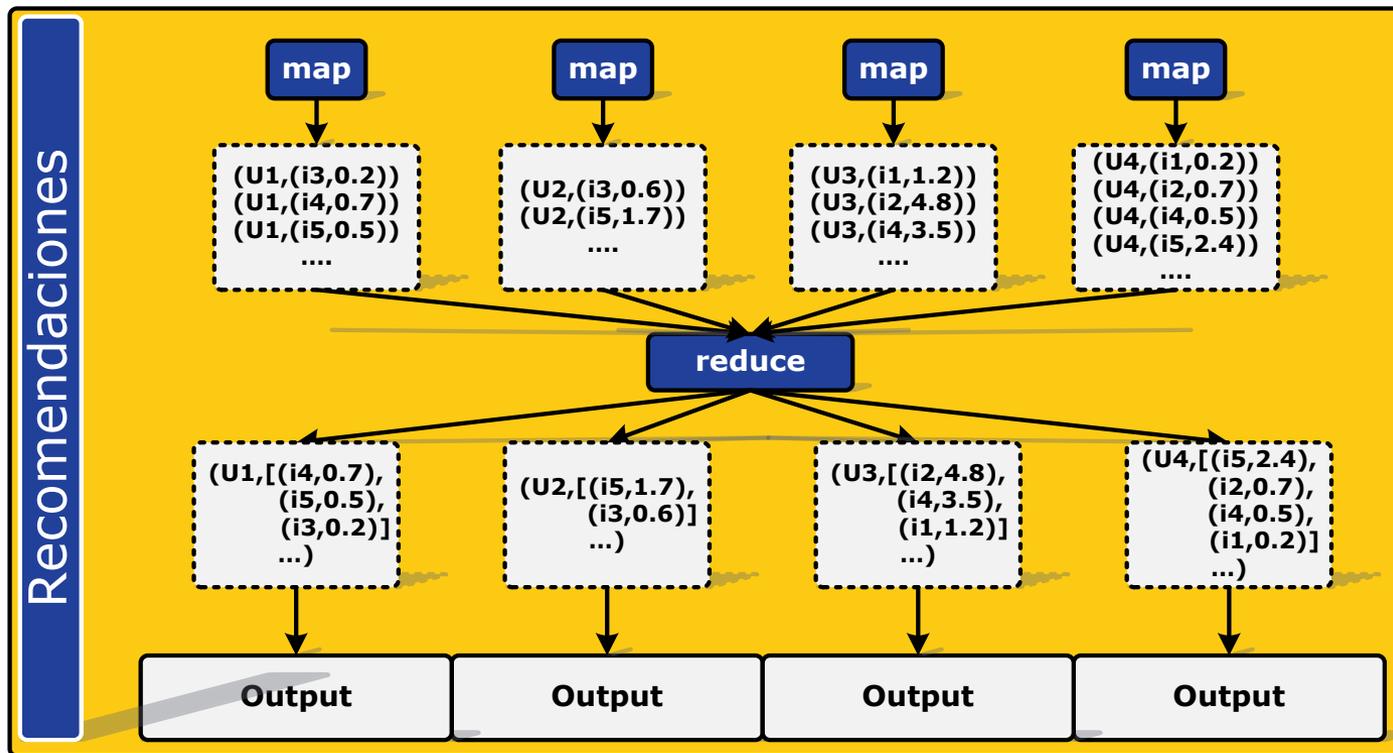
# Esquema map-reduce: Similitud

- Elección de los vecinos más cercanos

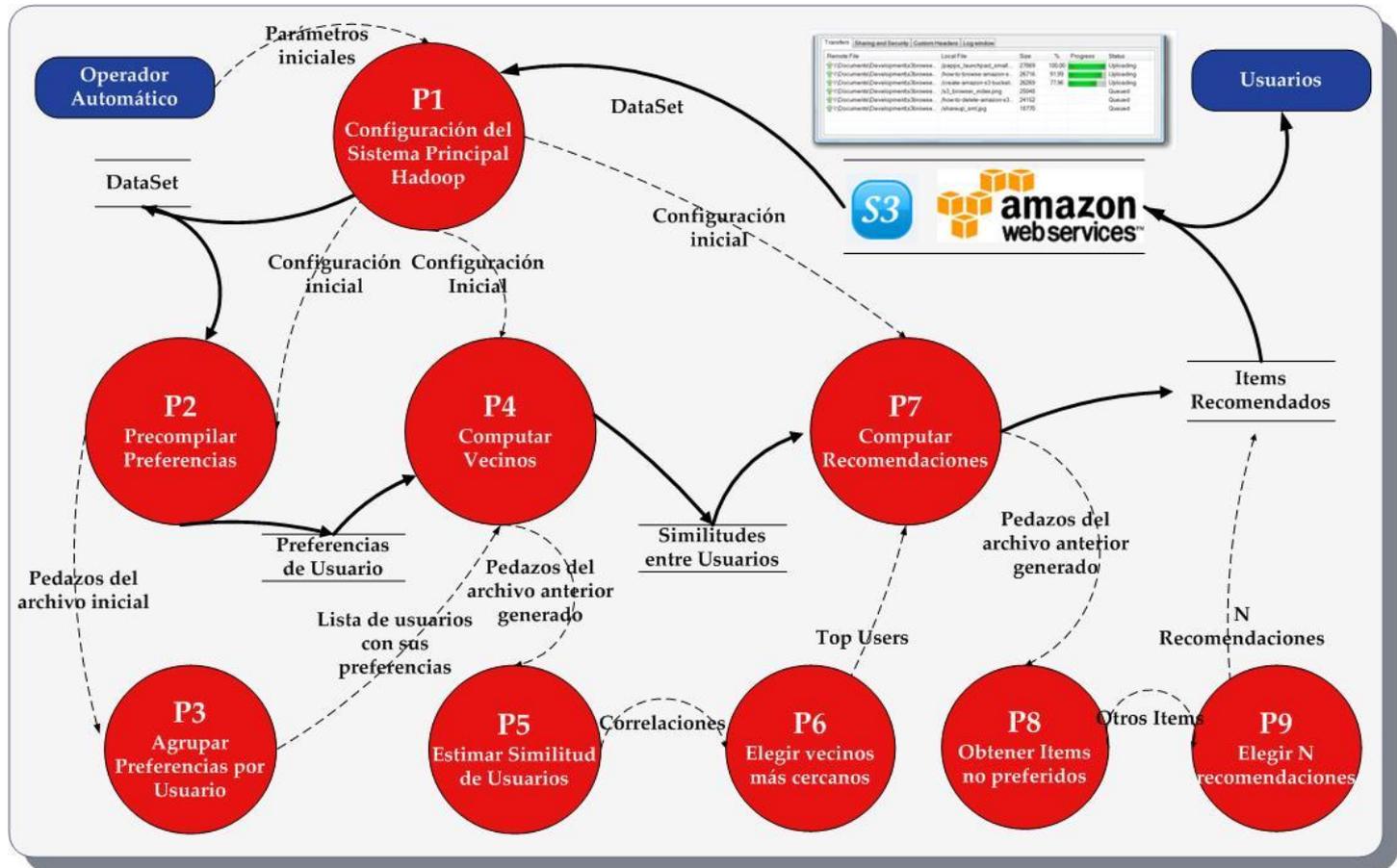


# Esquema map-reduce: Recomendaciones

- Selección de ítems nunca antes visto.



# Flujo de datos

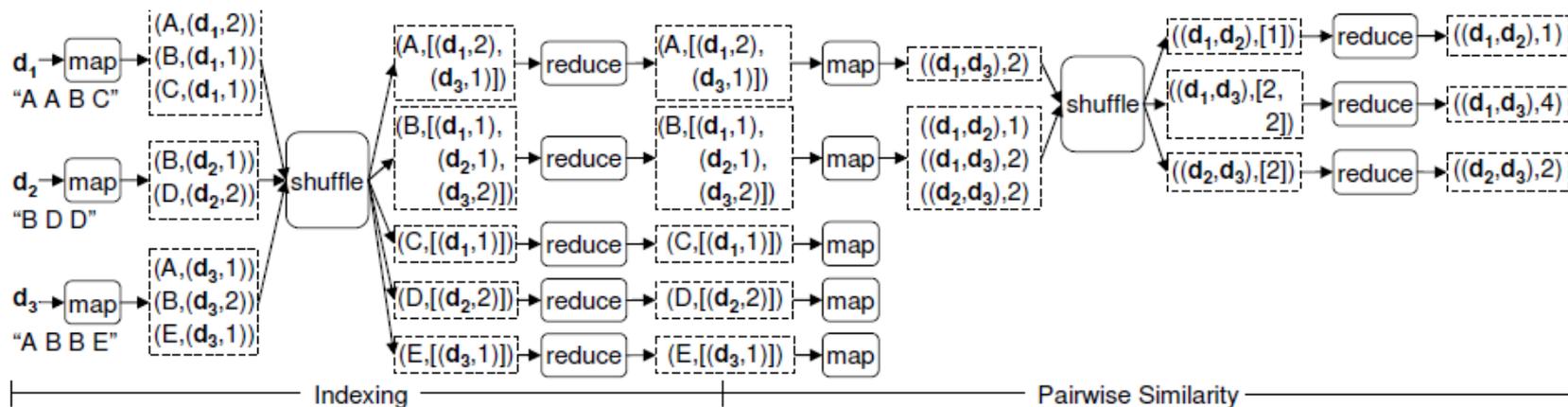


# Recomendaciones usando Pearson

- **Proceso 1:** Configuración del Sistema Principal Hadoop
- **Proceso 2:** Pre-compilar Preferencias
- **Proceso 3:** Agrupar Preferencias por Usuario
- **Proceso 4:** Computar Vecinos
- **Proceso 5:** Estimar Similitud de Usuarios
- **Proceso 6:** Elegir vecinos más cercanos
- **Proceso 7:** Computar Recomendaciones
- **Proceso 8:** Obtener ítems no preferidos
- **Proceso 9:** Elegir N recomendaciones

# Recomendaciones usando producto punto

- **Proceso 4:** Agrupar usuarios por ítem
- **Proceso 5:** Agrupar por tuplas de usuarios



Computo de los valores de similitud simple presentado por Tamer Elsayed, Jimmy Lin,<sup>†</sup> and Douglas W. Oard

# Configuraciones

- Parámetros de inicio:
  - <input>** Es la ruta de origen HFDS donde se almacena el dataset.
  - <output>** Es la ruta de destino HDFS donde se almacenan los resultados.
  - <nvecindad>** Es el número de vecinos más cercanos.
  - <nrecomendaciones>** Es el número de recomendaciones por usuario.

# Configuraciones

- Variables de inicio clave/valor para el **map o reduce** (opcional). Ejemplo:

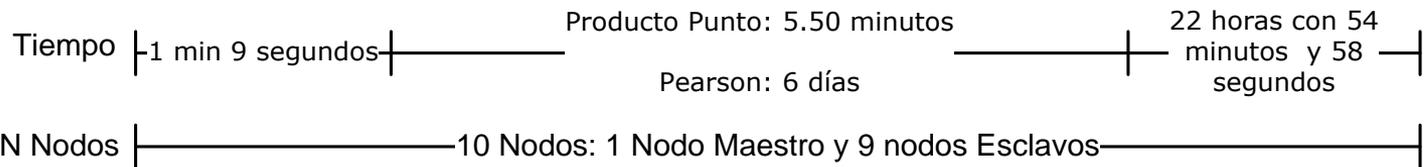
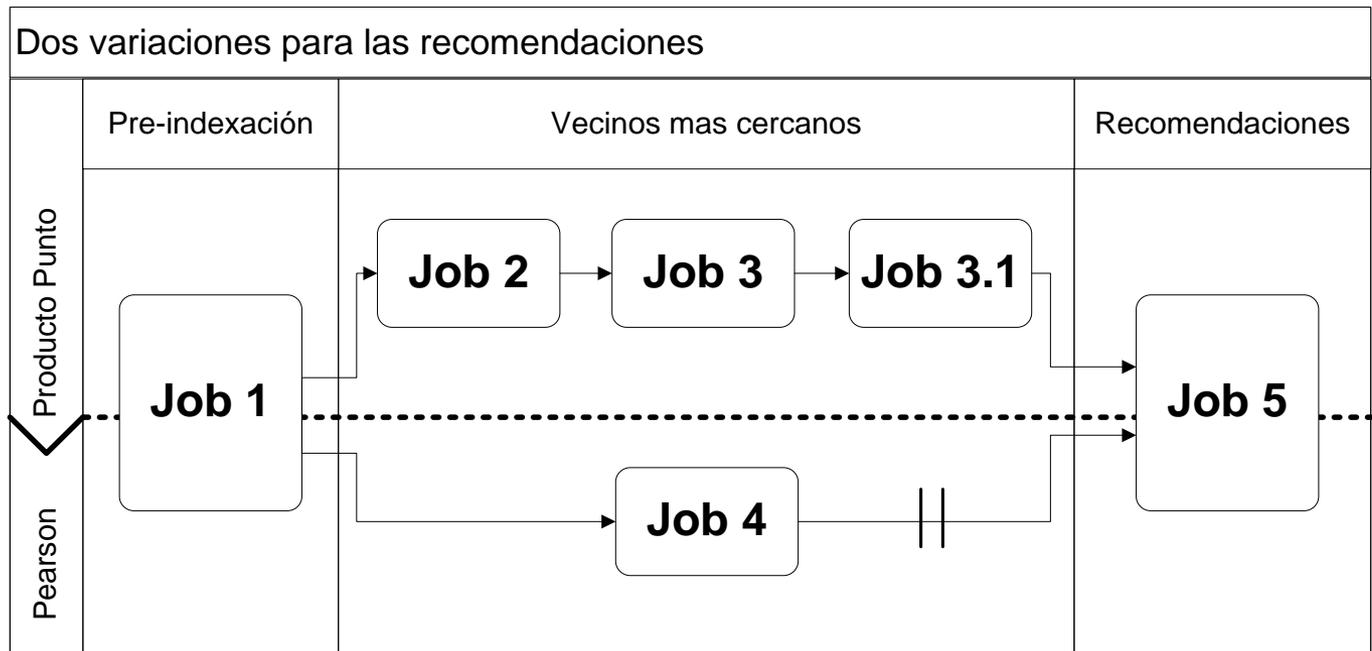
- `jobConf.set (UserNeighbourhoodReducer.NEAREST_NUSER_NEIGHBORHOOD, nvecindad);`

- `jobConf.set (ItemsRecommendedReducer.RECOMMENDATIONS_PER_USER, nrecomendaciones);`

- Carga de archivos en la cache distribuida (opcional). Ejemplo:

- `DistributedCache.addCacheFile (new URI (inputPathDis), jobConf);`

# Ejecución de Jobs con 2 alternativas

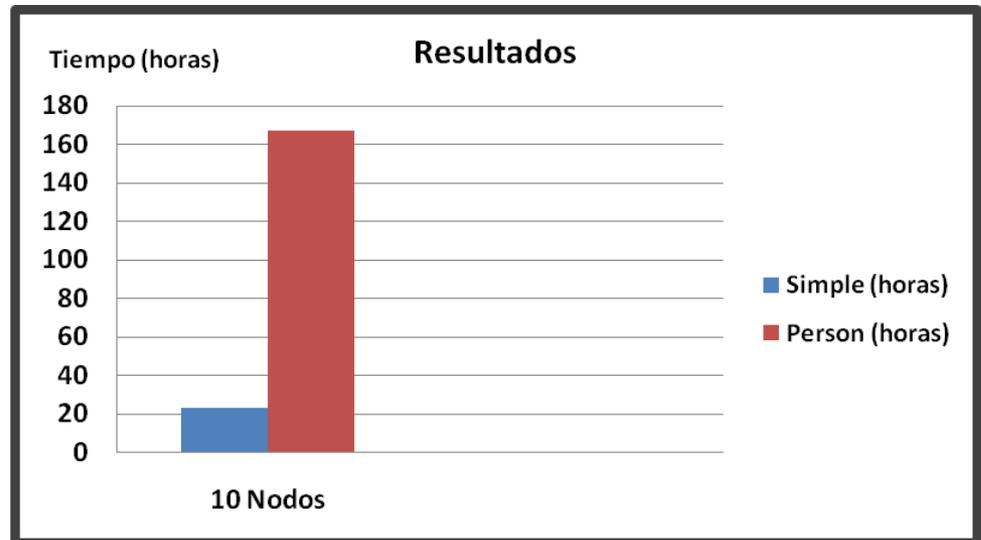


# Tiempos de respuesta

Job	map tasks	reduce tasks	Tiempo map-reduce (min)
2	4	10	1.25
3	10	10	3.4
3.1	10	10	0.85
		Total	5.50

Proceso de obtención de vecinos más cercanos basado en una correlación simple.

Reducción de un 99.93% con respecto a la alternativa de la correlación de Pearson.



# Costos

Descripción	Simple	Con Pearson
Constante (nodos)	\$0.20	\$0.20
Número de nodos en cluster	10	10
Duración cluster levantado (Horas)	24	169
Subtotal Nodos	\$48.00	\$338.00
Constante (Upload por GB)	\$0.10	\$0.10
Tamaño Upload (GB)	0.00000	0.00000
Subtotal Upload	\$0.00	\$0.00
Constante (Download por GB)	\$0.17	\$0.17
Tamaño Download (GB)	1	1
Subtotal Download	\$0.17	\$0.17
Costo Total Sesion	\$48.17	\$338.17

La correlación de Pearson de la librería de Mahout tiene un costo mucho mayor que utilizando la correlación simple .

Costo total del proyecto

Categoría	Valor
Saldo Inicial	\$200.00
Costo Total EC2 (US)	\$125.73
Costo Total Storage (US)	\$0.81
Costo Total Proyecto	\$126.54
Saldo Disponible	\$73.46

# Conclusiones

1. Los tiempos de ejecución de los algoritmos implementados aumentan con el cuadrado de la cantidad de usuarios y por ello presentan problemas de escalabilidad.
2. El costo del algoritmo de Pearson es tan alto en términos monetarios, que supera 8 a 1 al algoritmo basado en producto punto.
3. Aplicando producto punto no se obtienen recomendaciones para usuarios que han calificado ítems que otros usuarios no han calificado, Pearson resuelve este problema aplicando inferencias.

# Recomendaciones

1. La poca cantidad de nodos compartidos es una limitante. El producto punto se constituye en una alternativa viable para el procesamiento de las similitudes.
2. Para el seguimiento futuro de estos casos se puede buscar una alternativa al método de producto punto en el que se apliquen inferencias y utilice el esquema map-reduce para escalabilidad.