

Final Project : Build Recommendations system using movies data

Personalized Product Recommendations
with Neo4j

Recommendations

Personalized product recommendations can increase conversions, improve sales rates and provide a better experience for users. In this Neo4j Browser guide, we'll take a look at how you can generate graph-based real-time personalized product recommendations using a dataset of movies and movie ratings, but these techniques can be applied to many different types of products or content.

Graph-Based Recommendations

Generating personalized recommendations is one of the most common use cases for a graph database. Some of the main benefits of using graphs to generate recommendations include:

1. Performance. Index-free adjacency allows for calculating recommendations in real time, ensuring the recommendation is always relevant and reflecting up-to-date information.
2. Data model. The labeled property graph model allows for easily combining datasets from multiple sources, allowing enterprises to unlock value from previously separated data silos.

Data sources:

¥ We are going to import dataset from this file `data/all-plain.cypher`

¥ Pre-requis

1. Create empty database named recommendations
2. use recommendations
3. Import data from file `data/all-plain.cypher`

! hints :

! enable apoc import data

! Read `apoc.cypher.runFile` documentation

! Mandatory " Use `apoc.cypher.runFile` to import data into your new database

The Open Movie Graph Data Model

The Property Graph Model

The data model of graph databases is called the labeled property graph model.

Nodes: The entities in the data.

Labels: Each node can have one or more label that specifies the type of the node.

Relationships: Connect two nodes. They have a single direction and type.

Properties: Key-value pair properties can be stored on both nodes and relationships.

Eliminate Data Silos

In this use case, we are using graphs to combine data from multiple sources.

Product Catalog: Data describing movies comes from the product catalog silo.

User Purchases / Reviews: Data on user purchases and reviews comes from the user or transaction source.

By combining these two in the graph, we are able to query across datasets to generate personalized product recommendations.

Nodes

Movie, Actor, Director, User, Genre are the labels used in this example.

Relationships

ACTED_IN, IN_GENRE, DIRECTED, RATED are the relationships used in this example.

Properties

title, name, year, rating are some of the properties used in this example.

Memo on Cypher

In order to work with our labeled property graph, we need a query language for graphs.

Graph Patterns

Cypher is the query language for graphs and is centered around graph patterns. Graph patterns are expressed in Cypher using ASCII-art like syntax.

Nodes

Nodes are defined within parentheses (). Optionally, we can specify node label(s): (:Movie)

Relationships

Relationships are defined within square brackets []. Optionally we can specify type and direction:

```
<code>(:Movie)<strong><-[:RATED]-</strong>(:User)</code>
```

Variables

Graph elements can be bound to variables that can be referred to later in the query:

```
<code>(<strong>m</strong>:Movie)<-[:<strong>r</strong>:RATED]-(<strong>u</strong>:User)</code>
```

Predicates

Filters can be applied to these graph patterns to limit the matching paths. Boolean logic operators, regular expressions and string comparison operators can be used here within the WHERE clause, e.g. WHERE m.title CONTAINS 'Matrix'

Aggregations

There is an implicit group of all non-aggregated fields when using aggregation functions such as count.

Use the [Cypher Refcard](#) as a syntax reference.

WORK TO DO

Dissecting a Cypher Statement

Let's implement a Cypher query that answers the question "How many reviews does each Matrix movie have?". Don't worry if this seems complex, we'll build up our understanding of Cypher as we move along.

Int: Replace ??? by the correct values

```
MATCH (m: ???)-[:???]-(u: ???)
WHERE m.??? CONTAINS '???'
WITH m, count(*) AS reviews
RETURN m.title AS movie, reviews
ORDER BY reviews DESC LIMIT 5;
```

2/ After you completed previous request and tested it, create your own User defined procedure to do the same work.

Personalized Recommendations

Now let's start generating some recommendations. There are two basic approaches to recommendation algorithms.

Content-Based Filtering

Recommend items that are similar to those that a user is viewing, rated highly or purchased previously.

1/ "Find Items similar to the item you're looking at now"

TODO: Create Cypher query

2/ After you completed previous request and tested it, create your own User defined procedure to do the same work.

Collaborative Filtering

Use the preferences, ratings and actions of other users in the network to find items to recommend.

1/ "Get Users who got this item, also got that other item."

TODO: Create Cypher query

2/ After you completed previous request and tested it, create your own User defined procedure to do the same work.

Content-Based Filtering

The goal of content-based filtering is to find similar items, using attributes (or traits) of the item. Using our movie data, one way we could define similarity is movies that have common genres.

Similarity Based on Common Genres

1/ Find movies most similar to Inception based on shared genres

```
// Find similar movies by common genres  
TODO: Create Cypher query
```

2/ After you completed previous request and tested it, create your own User defined procedure to do the same work.

Personalized Recommendations Based on Genres

If we know what movies a user has watched, we can use this information to recommend similar movies:

1/ Recommend movies similar to those the user has already watched

```
TODO: Create Cypher query  
  
// Content recommendation by overlapping genres
```

2/ After you completed previous request and tested it, create your own User defined procedure to do the same work.

Weighted Content Algorithm

Of course there are many more traits in addition to just genre that we can consider to compute similarity, such as actors and directors. Let's use a weighted sum to score the recommendations based on the number of actors (3x), genres (5x) and directors (4x) they have in common to boost the score:

Compute a weighted sum based on the number and types of overlapping traits

```
Hint: Find similar movies by common genres  
TODO: Create Cypher query
```

Content-Based Similarity Metrics

So far we've used the number of common traits as a way to score the relevance of our recommendations. Let's now consider a more robust way to quantify similarity, using a similarity metric. Similarity metrics are an important component used in generating personalized recommendations that allow us to quantify how similar two items (or as we'll see later, how similar two users preferences) are.

Jaccard Index

The Jaccard index is a number between 0 and 1 that indicates how similar two sets are. The Jaccard index of two identical sets is 1. If two sets do not have a common element, then the Jaccard index is 0. The Jaccard is calculated by dividing the size of the intersection of two sets by the union of the two sets.

We can calculate the Jaccard index for sets of movie genres to determine how similar two movies are.

What movies are most similar to Inception based on Jaccard similarity of genres?

TODO: Create Cypher query

Apply this same approach to all "traits" of the movie (genre, actors, directors, etc.):

TODO: Create Cypher query

Collaborative Filtering & Leveraging Movie Ratings

Notice that we have user-movie ratings in our graph. The collaborative filtering approach is going to make use of this information to find relevant recommendations.

Steps:

1. Find similar users in the network (our peer group).
2. Assuming that similar users have similar preferences, what are the movies those similar users like?

Show all ratings by Misty Williams

```
// Show all ratings by Misty Williams  
TODO: Create Cypher query
```

Find Misty's average rating

```
// Show average ratings by Misty Williams  
TODO: Create Cypher query
```

What are the movies that Misty liked more than average?

```
// What are the movies that Misty liked more than average?  
TODO: Create Cypher query
```

Collaborative Filtering & The Wisdom of Crowds

Simple Collaborative Filtering

Here we just use the fact that someone has rated a movie, not their actual rating to demonstrate the structure of finding the peers. Then we look at what else the peers rated, that the user has not rated themselves yet.

```
MATCH (u:User {name: 'Cynthia Freeman'})-[:RATED]->  
      (:Movie)-[:RATED]-(peer:User)  
MATCH (peer)-[:RATED]->(rec:Movie)  
WHERE NOT EXISTS { (u)-[:RATED]->(rec) }  
RETURN rec.title, rec.year, rec.plot  
LIMIT 25
```

Of course this is just a simple approach, there are many problems with this query, such as not normalizing based on popularity or not taking ratings into consideration. We'll do that next, looking at movies being rated similarly, and then picking highly rated movies and using their rating and frequency to sort the results.

```
MATCH (u:User {name: 'Cynthia Freeman'})-[:RATED]->  
      (:Movie)-[:RATED]-(peer:User)
```

```

WHERE abs(r1.rating-r2.rating) < 2 // similarly rated
WITH distinct u, peer
MATCH (peer)-[r3:RATED]->(rec:Movie)
WHERE r3.rating > 3
& AND NOT EXISTS { (u)-[:RATED]->(rec) }
WITH rec, count(*) as freq, avg(r3.rating) as rating
RETURN rec.title, rec.year, rating, freq, rec.plot
ORDER BY rating DESC, freq DESC
LIMIT 25

```

In the next section, we will see how we can improve this approach using the kNN method.

Only Consider Genres Liked by the User

Many recommender systems are a blend of collaborative filtering and content-based approaches:

For a particular user, what genres have a higher-than-average rating? Use this to score similar movies

```

TODO :
// 1 compute mean rating

// 2 find genres with higher than average rating and their number of rated movies

// 3 find movies in those genres, that have not been watched yet

```

Collaborative Filtering & Similarity Metrics

We use similarity metrics to quantify how similar two users or two items are. We've already seen Jaccard similarity used in the context of content-based filtering. Now, we'll see how similarity metrics are used with collaborative filtering.

Cosine Distance

Jaccard similarity was useful for comparing movies and is essentially comparing two sets (groups of genres, actors, directors, etc.). However, with movie ratings each relationship has a weight that we can consider as well.

Cosine Similarity

The cosine similarity of two users will tell us how similar two users' preferences for movies are. Users with a high cosine similarity will have similar preferences.

Find the users with the most similar preferences to Cynthia Freeman, according to cosine similarity

```
TODO
// Most similar users using Cosine similarity
```

We can also compute this measure using the [Cosine Similarity algorithm](#) in the Neo4j Graph Data Science Library.

Find the users with the most similar preferences to Cynthia Freeman, according to cosine similarity function

```
TODO:
hint : gds.similarity.cosine
```

Collaborative Filtering & Similarity Metrics

Pearson Similarity

Pearson similarity, or Pearson correlation, is another similarity metric we can use. This is particularly well-suited for product recommendations because it takes into account the fact that different users will have different mean ratings: on average some users will tend to give higher ratings than others. Since Pearson similarity considers differences about the mean, this metric will account for these discrepancies.

Find users most similar to Cynthia Freeman, according to Pearson similarity

TODO

We can also compute this measure using the [Pearson Similarity algorithm](#) in the Neo4j Graph Data Science Library.

Find users most similar to Cynthia Freeman, according to the Pearson similarity function

TODO

Collaborative Filtering & Neighborhood-Based Recommendations

kNN & K-Nearest Neighbors

Now that we have a method for finding similar users based on preferences, the next step is to allow each of the k most similar users to vote for what items should be recommended.

Essentially:

"Who are the 10 users with tastes in movies most similar to mine? What movies have they rated highly that I haven't seen yet?"

kNN movie recommendation using Pearson similarity

TODO

Further Work

Optional Exercises

Extend these queries:

Temporal component

Preferences change over time, use the rating timestamp to consider how more recent ratings might be used to find more relevant recommendations.

Keyword extraction

Enhance the traits available using the plot description.
How would you model extracted keywords for movies?

Image recognition using posters

There are several libraries and APIs that offer image recognition and tagging.

Since we have movie poster images for each movie, how could we use these to enhance our recommendations?