

Recommendations with Neo4j and Python

Grant Beasley Twitter - @grantbeasley
Blog - grant592.github.io

Senior Data Scientist at Sage

About Me



SaleCycle Sage



HISIsUS



Content

- Context
- Recommendations
- Graph model
- Image embeddings
- Collaborative filtering
- Bringing it all together

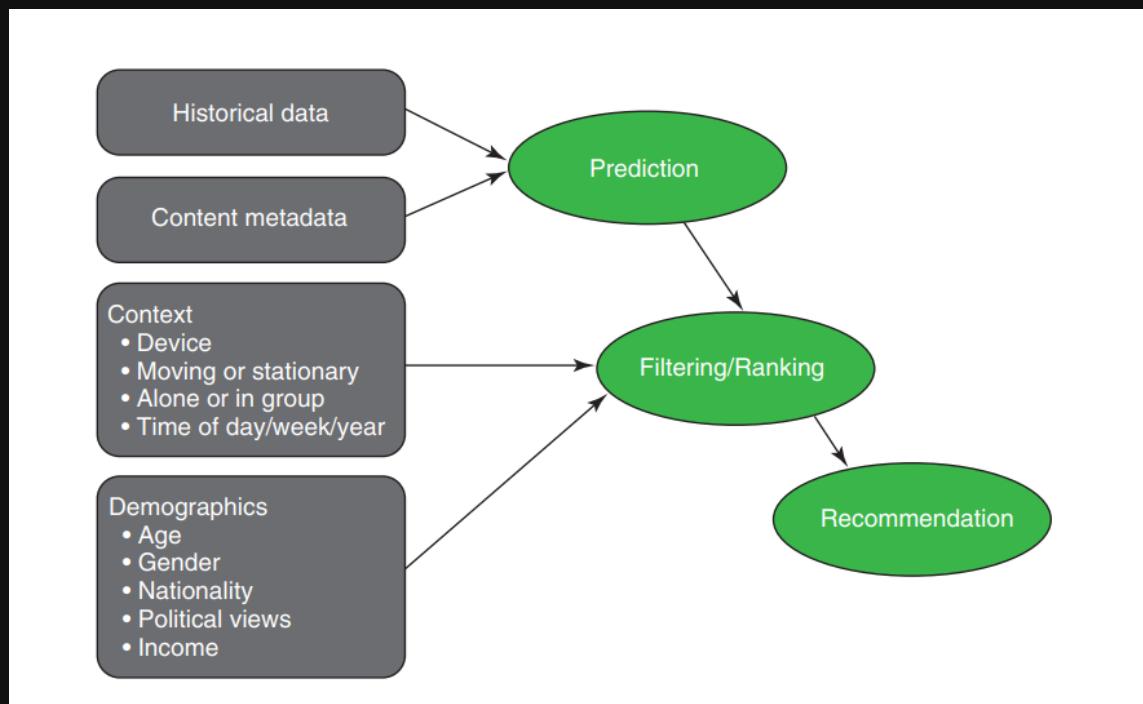
- Kim Falk - Practical Recommender Systems
- Ethan Rosenthal - blog

Recommendations

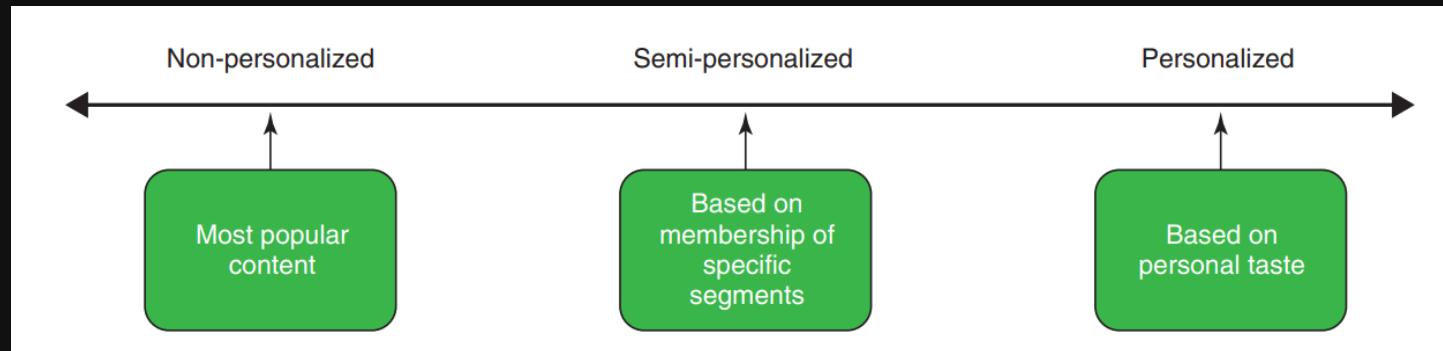
“A recommender system calculates and provides relevant content to the user based on knowledge of the user, content, and interactions between the user and the item.”

- The long tail
- Netflix

Taxonomy of Recommender System



Levels of Personalisation



Other Considerations

- Tackling the cold start problem
- Grey sheep



Background and Why

- Kaggle Competition
- Combine different approaches
- Neo4j deep dive
- Safety Notice!



I LIKE THE PLASMIC NATURE
OF YOUR DATA MODELLING. NICE.

The Data Model

Loading with Cypher

Create indexes, load data, create relationships.

```
1 CREATE CONSTRAINT uniqueProduct IF NOT EXISTS
2 FOR (n:Product) REQUIRE (n.code) IS UNIQUE;
3
4 CREATE CONSTRAINT uniqueColourGroup IF NOT EXISTS
5 FOR (n:ColourGroup) REQUIRE (n.code) IS UNIQUE;
6
7 CREATE CONSTRAINT uniqueColourValue IF NOT EXISTS
8 FOR (n:PerceivedColourValue) REQUIRE (n.id) IS UNIQUE;
9
10 CREATE CONSTRAINT uniqueColourMaster IF NOT EXISTS
11 FOR (n:PerceivedColourMaster) REQUIRE (n.id) IS UNIQUE;
12
13 CREATE CONSTRAINT uniqueProductType IF NOT EXISTS
14 FOR (n:ProductType) REQUIRE (n.number) IS UNIQUE;
15
16 CREATE CONSTRAINT uniqueProductGroup IF NOT EXISTS
```

Loading with Cypher

Create indexes, load data, create relationships.

```
1 CREATE CONSTRAINT uniqueProduct IF NOT EXISTS
2 FOR (n:Product) REQUIRE (n.code) IS UNIQUE;
3
4 CREATE CONSTRAINT uniqueColourGroup IF NOT EXISTS
5 FOR (n:ColourGroup) REQUIRE (n.code) IS UNIQUE;
6
7 CREATE CONSTRAINT uniqueColourValue IF NOT EXISTS
8 FOR (n:PerceivedColourValue) REQUIRE (n.id) IS UNIQUE;
9
10 CREATE CONSTRAINT uniqueColourMaster IF NOT EXISTS
11 FOR (n:PerceivedColourMaster) REQUIRE (n.id) IS UNIQUE;
12
13 CREATE CONSTRAINT uniqueProductType IF NOT EXISTS
14 FOR (n:ProductType) REQUIRE (n.number) IS UNIQUE;
15
16 CREATE CONSTRAINT uniqueProductGroup IF NOT EXISTS
```

Loading with Cypher

```
1 LOAD CSV WITH HEADERS FROM 'file:///articles.csv' AS row
2 MERGE (ga:GraphicalAppearance {number: row.graphical_appearance_no})
3 ON CREATE SET
4   ga.name = row.graphical_appearance_name
5 MERGE (dept:Department {number: row.department_no})
6 ON CREATE SET
7   dept.name = row.department_name
8 MERGE (sect:Section {number: row.section_no})
9 ON CREATE SET
10  sect.name = row.section_name
11 MERGE (indx:Index {code: row.index_code})
12 ON CREATE SET
13  indx.name = row.index_name
14 MERGE (indxGrp:IndexGroup {number: row.index_group_no})
15 ON CREATE SET
16  indxGrp.name = row.index_group_name
```

Loading with Cypher

```
34 ON CREATE SET
35   prod.name = row.prod_name,
36   prod.description = row.detail_desc
37
38 WITH ga, dept, sect, idx, idxGrp, garmGrp, pcv, pcm, cg, pro
39 MERGE (prod)-[:HAS_COLOUR_VALUE]->(pcv)
40 MERGE (prod)-[:HAS_MASTER_COLOUR]->(pcm)
41 MERGE (prod)-[:HAS_COLOUR]->(cg)
42 MERGE (prod)-[:SUBSET_OF]->(prodType)
43 MERGE (prodType)-[:SUBSET_OF]->(prodGrp)
44 MERGE (prod)-[:HAS_APPEARANCE_TYPE]->(ga)
45 MERGE (prod)-[:FROM_DEPARTMENT]->(dept)
46 MERGE (dept)-[:FROM_SECTION]->(sect)
47
48 MERGE (prod)-[:HAS_INDEX]->(idx)
49 MERGE (idx)-[:HAS_GROUP]->(idxGrp)
50 MERGE (idx)-[:HAS_GARMENT_GROUP]->(garmGrp)
```

Loading with Cypher

Customer data

```
1 CREATE CONSTRAINT uniqueCustomer IF NOT EXISTS FOR (n:Customer)
2
3 LOAD CSV WITH HEADERS FROM 'file:///customers.csv' AS row
4 CALL {
5   WITH row
6   MERGE (c:Customer {id: row.customer_id})
7   ON CREATE SET
8     c.club_member_status =
9     CASE
10       WHEN row.club_member_status IS NULL THEN 'NOT_ACTIVE'
11       ELSE row.club_member_status
12     END,
13     c.active =
14     CASE
15       WHEN row.Active IS NULL THEN 0
16       ELSE row.Active
```

Loading with Cypher

Customer data

```
1 CREATE CONSTRAINT uniqueCustomer IF NOT EXISTS FOR (n:Customer)
2
3 LOAD CSV WITH HEADERS FROM 'file:///customers.csv' AS row
4 CALL {
5   WITH row
6   MERGE (c:Customer {id: row.customer_id})
7   ON CREATE SET
8     c.club_member_status =
9     CASE
10       WHEN row.club_member_status IS NULL THEN 'NOT_ACTIVE'
11       ELSE row.club_member_status
12     END,
13     c.active =
14     CASE
15       WHEN row.Active IS NULL THEN 0
16       ELSE row.Active
```

Loading with Cypher

Customer data

```
24     CASE
25         WHEN row.fashion_news_frequency IS NULL THEN 'NONE'
26         ELSE row.fashion_news_frequency
27     END,
28     c.age =
29     CASE
30         WHEN row.age IS NULL THEN 100
31         ELSE row.age
32     END,
33     c.postal_code =
34     CASE
35         WHEN row.postal_code IS NULL THEN 'NONE'
36         ELSE row.postal_code
37     END
38
39 } IN TRANSACTIONS OF 1000 ROWS;
```

Loading with Cypher

Second pass of the data.

```
1 LOAD CSV WITH HEADERS FROM 'file:///transactions_train.csv' AS
2 CALL {
3   WITH row
4     MATCH (c:Customer {id: row.customer_id})
5     MATCH (p:Product {code: row.article_id})
6     MERGE (c)-[:PURCHASED]{
7       t_date: date(row.t_date),
8       price: row.price,
9       sales_channel: row.sales_channel_id}]->(p)
10 } IN TRANSACTIONS OF 1000 ROWS;
```

Bucketizing the data

Age categories

```
1 //  
2 WITH [[16, 25], [26, 35], [36, 45], [46, 55], [56, 65], [65, 99], [100, 100]]  
3 UNWIND ageGroups as ag  
4 CREATE (:AgeGroup {lower: ag[0], upper: ag[1]});  
5  
6 // Create relationship between age group and customer  
7 MATCH (c:Customer)  
8 MATCH (a:AgeGroup) WHERE toInteger(c.age) >= a.lower AND toInteger(c.age) <= a.upper  
9 CREATE (c)-[:IS AGE GROUP]->(a);
```

EDA

EDA

Label

Total Count

EDA

Label	Total Count
-------	-------------

Customers	1,371,980
-----------	-----------

EDA

Label	Total Count
Customers	1,371,980
Product	105,542

EDA

Label	Total Count
Customers	1,371,980
Product	105,542
Purchases	28,813,419

EDA

EDA

Label

Average **Std Dev**

EDA

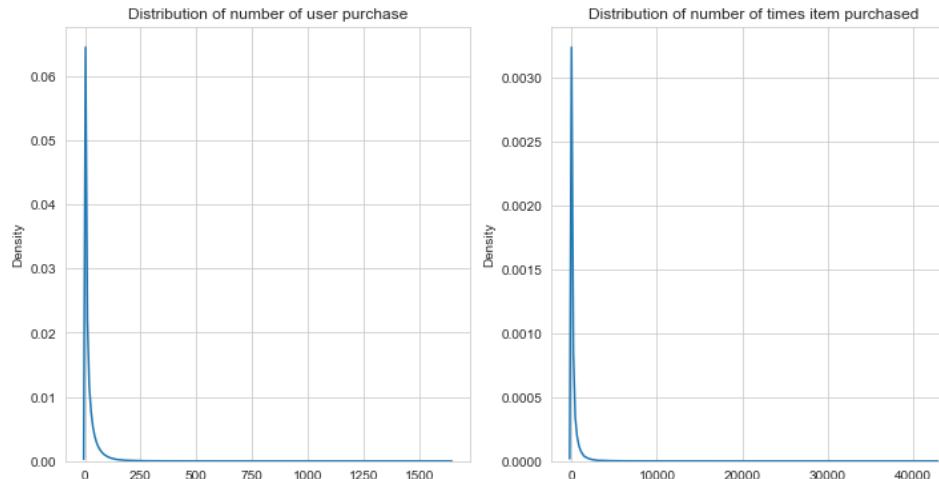
Label	Average	Std Dev
Customer Purchases	21.2	34.6

EDA

Label	Average	Std Dev
Customer Purchases	21.2	34.6
Product Purchased	275.6	696.9

EDA

Label	Average	Std Dev
Customer Purchases	21.2	34.6
Product Purchased	275.6	696.9



EDA

Most Popular Product Types

```
1 with driver.session() as session:  
2     query = """  
3         MATCH (p:Product)-[:SUBSET_OF]->(n:ProductType)  
4         RETURN n.name as product_type, COUNT(p) as count  
5     """  
6     data = list(session.run(query))  
7  
8     df = pd.DataFrame([dict(record) for record in data])  
9  
10    top_prods = df.sort_values('count', ascending=False).head()
```

EDA

Most Popular Product Types

```
1 with driver.session() as session:  
2     query = """  
3         MATCH (p:Product)-[:SUBSET_OF]->(n:ProductType)  
4         RETURN n.name as product_type, COUNT(p) as count  
5     """  
6     data = list(session.run(query))  
7  
8     df = pd.DataFrame([dict(record) for record in data])  
9  
10    top_prods = df.sort_values('count', ascending=False).head()
```

EDA

Most Popular Product Types

```
1 with driver.session() as session:  
2     query = """  
3         MATCH (p:Product)-[:SUBSET_OF]->(n:ProductType)  
4         RETURN n.name as product_type, COUNT(p) as count  
5     """  
6     data = list(session.run(query))  
7  
8     df = pd.DataFrame([dict(record) for record in data])  
9  
10    top_prods = df.sort_values('count', ascending=False).head()
```

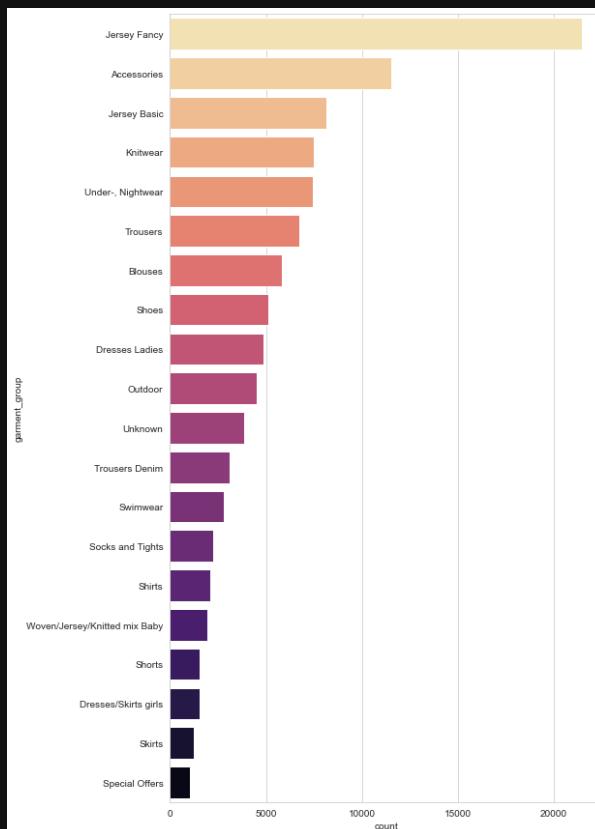

EDA

Most Popular Garment Types

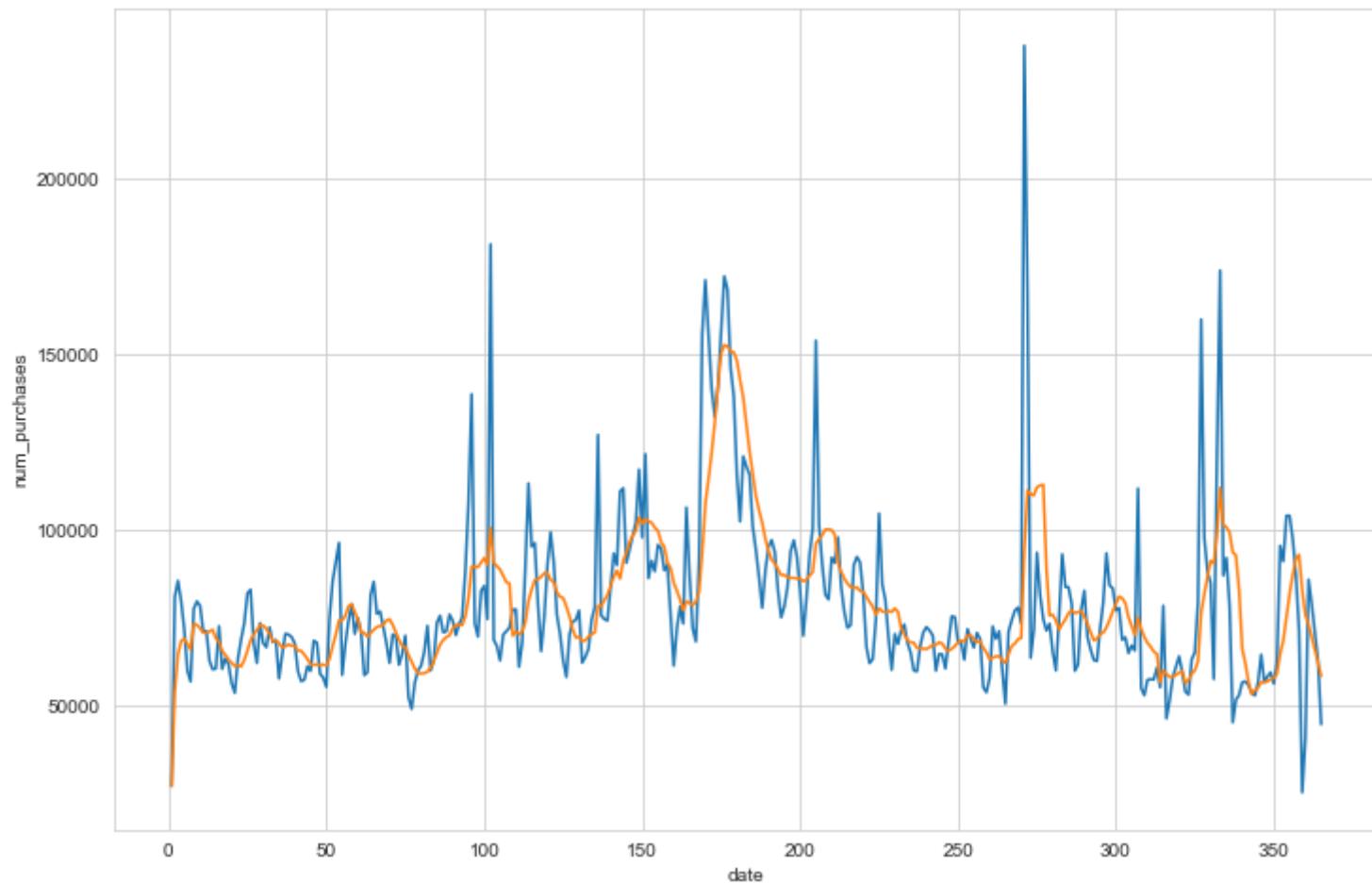
```
1 with driver.session() as session:  
2     query = """  
3         MATCH (p:Product)-[:HAS_GARMENT_GROUP]->(g:GarmentGroup)  
4         RETURN g.name as garment_group, COUNT(g) as count  
5     """  
6     data = list(session.run(query))
```

EDA

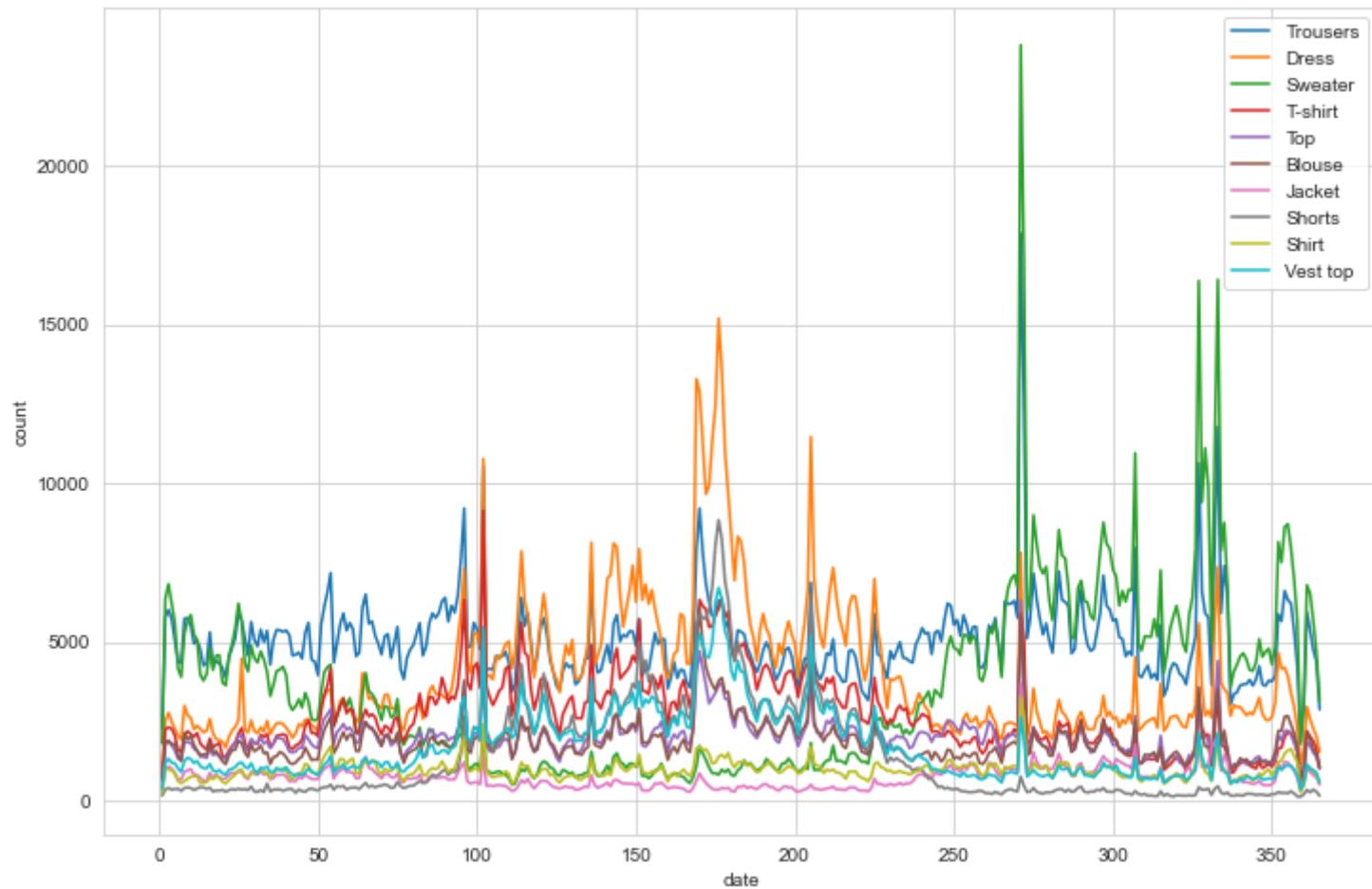
Most Popular Garment Types



Seasonal Effects



Seasonal Effects





Age Effects

Is they just don't make 'em like they used to
No! Nobody ever made them like this!

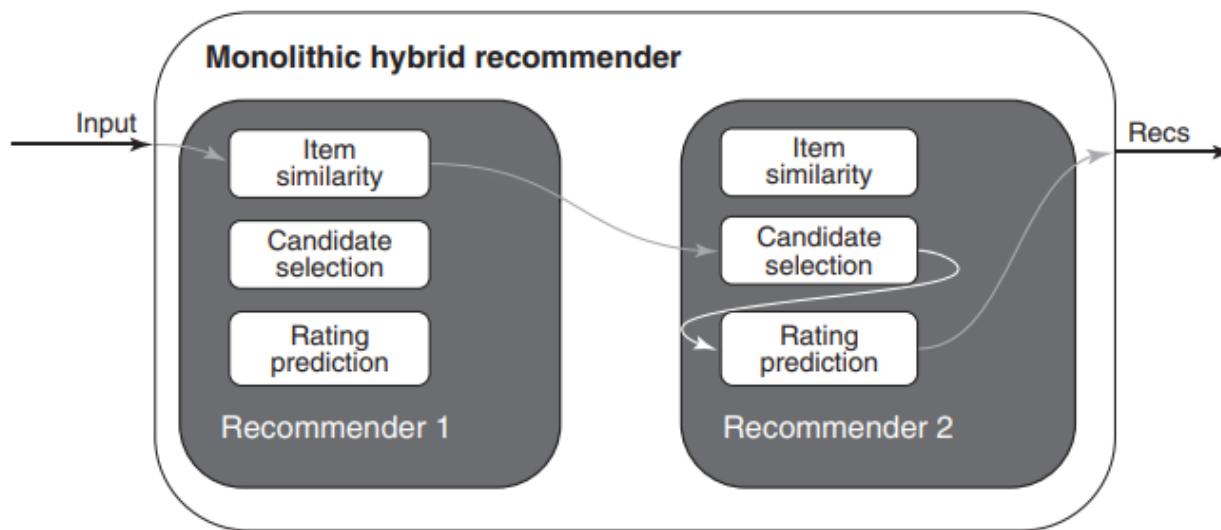
A photograph of Steve Harvey, a well-known American comedian, actor, and television host. He is shown from the chest up, wearing a dark suit jacket over a white shirt. His head is tilted back, eyes wide open with a look of surprise or excitement, and his mouth is slightly open. His hands are raised in front of him, palms facing up, as if he has just discovered something unexpected. He is positioned behind a light-colored wooden desk. The background is a warm, orange-toned wall.

Exploration over...let's get to work

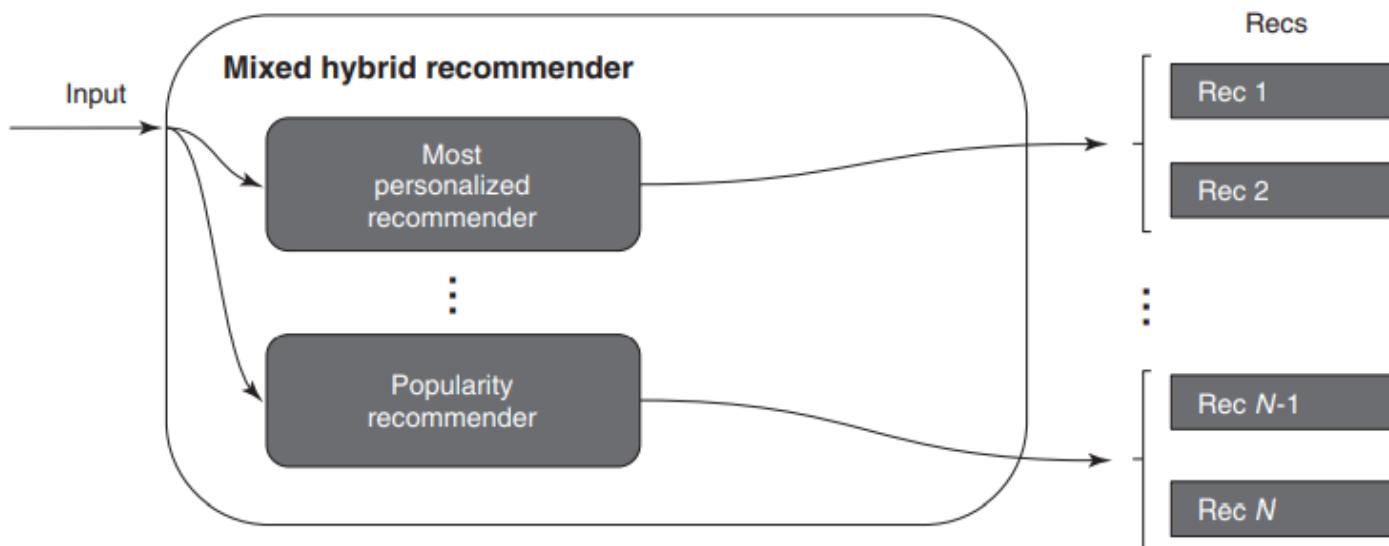
Considerations

- How to make best use of graph database?
- How to approach recommendations

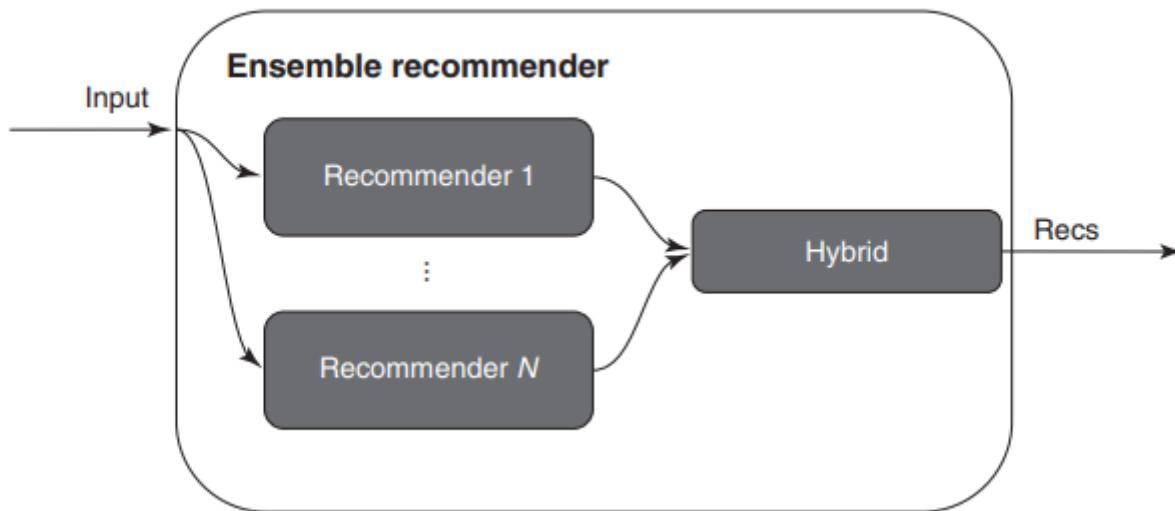
Recommender System



Recommender System

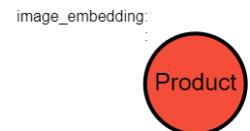
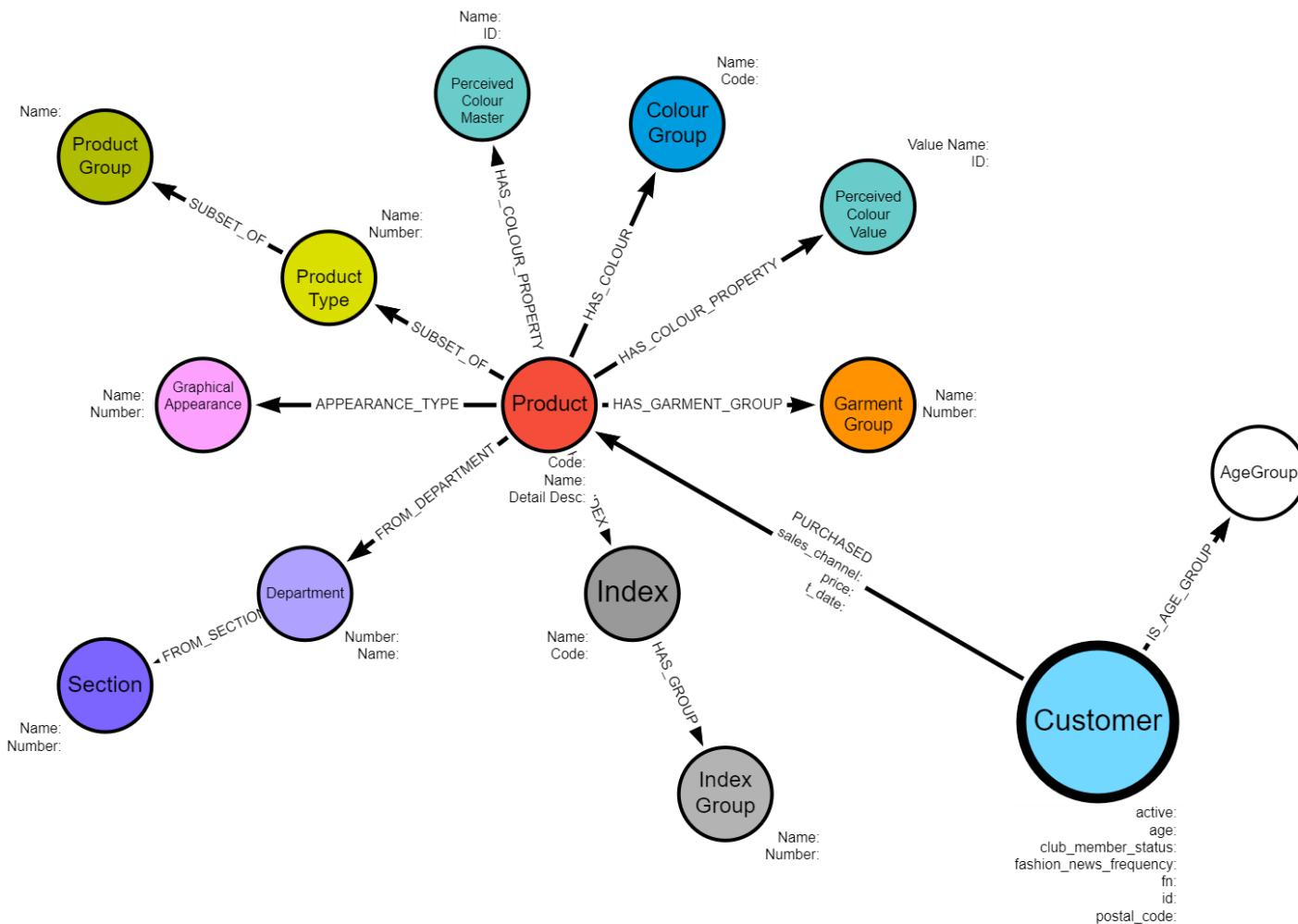


Recommender System



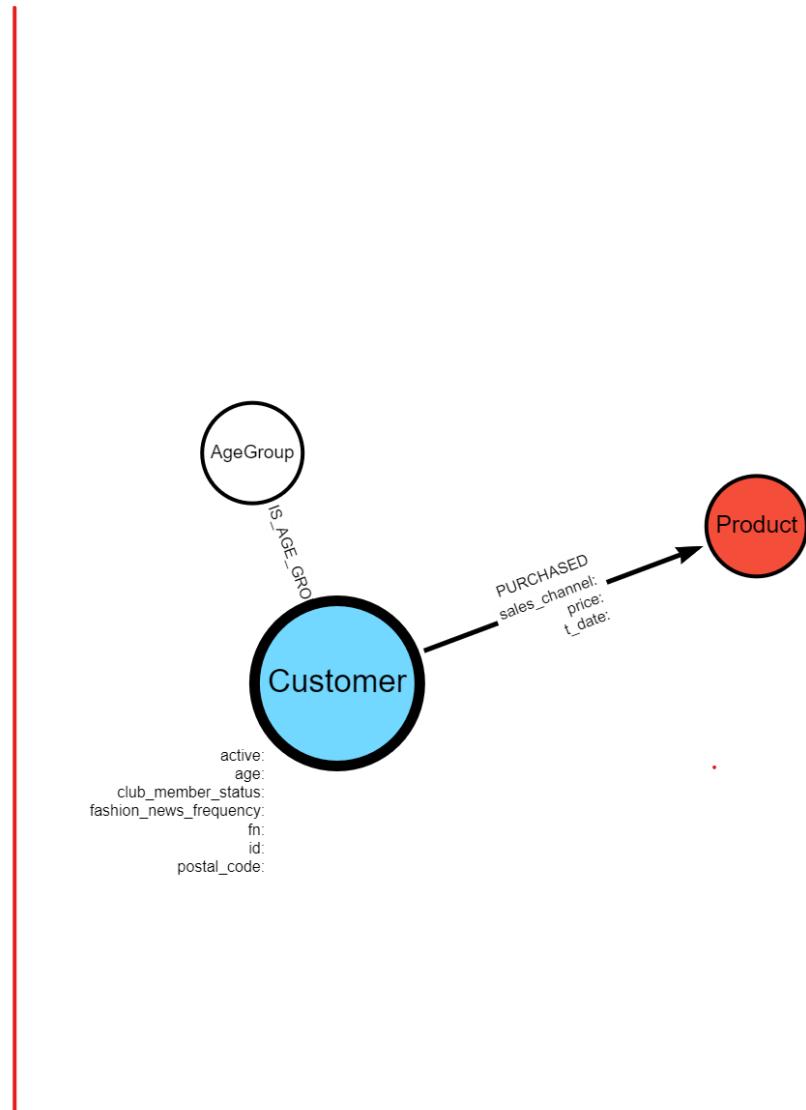
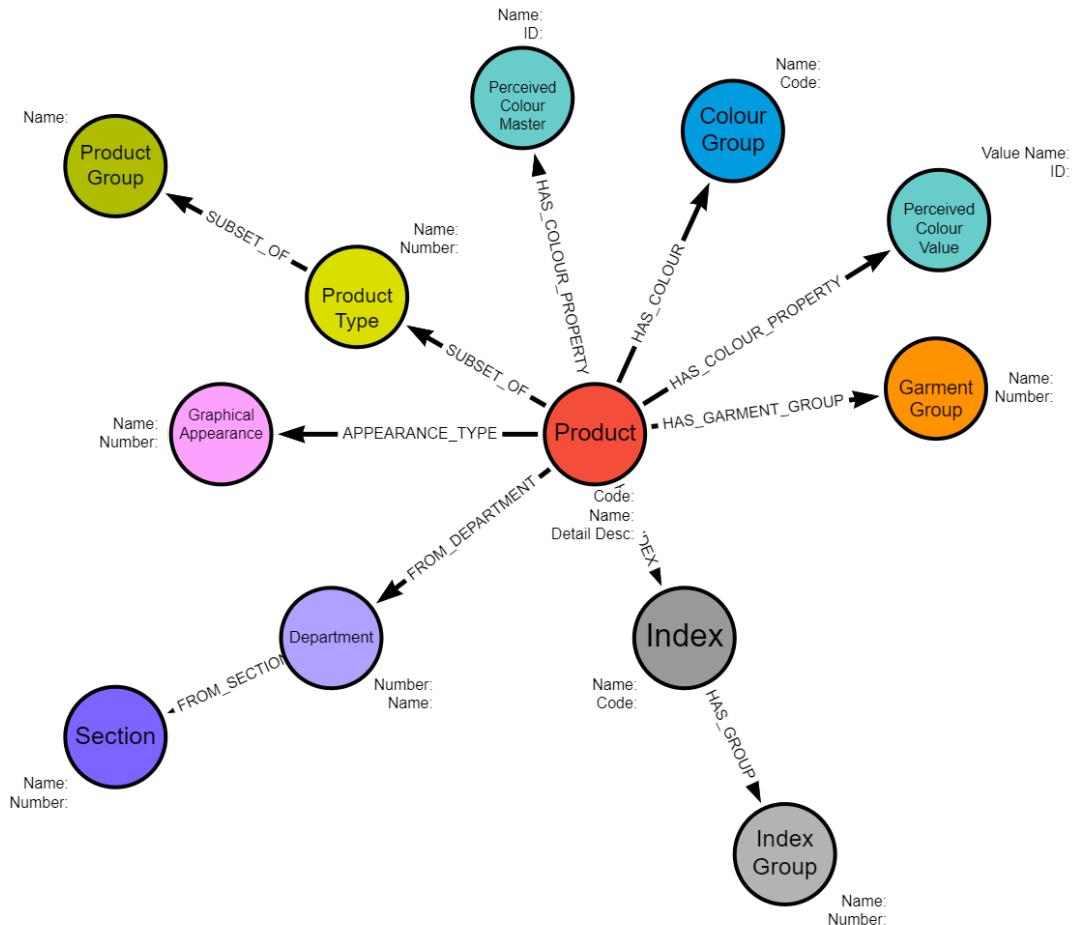
Neo4j Fabric

The initial data model



Neo4j Fabric

What it should have looked like...



Neo4j Fabric

Fabric Setup

```
1 #Fabric config
2 fabric.database.name=fabric
3
4 fabric.graph.0.uri=neo4j://localhost:7687
5 fabric.graph.0.database=neo4j
6 fabric.graph.0.name=neo4j
7
8 fabric.graph.1.uri=neo4j://localhost:7687
9 fabric.graph.1.database=products
10 fabric.graph.1.name=products
```

Neo4j Fabric

Fabric Setup

```
1 # Load the second DB
2 :use system;
3 create database products;
4 use products;
5 CREATE CONSTRAINT uniqueProduct IF NOT EXISTS FOR (n:Product) R
6
7 LOAD CSV WITH HEADERS FROM 'file:///articles.csv' AS row
8 MERGE (prod:Product {code: row.article_id})
```

Warning!

Explore vs exploit

Image Embeddings and Transfer Learning

"Transfer learning (TL) is a research problem in machine learning (ML) that focuses on storing knowledge gained while solving one problem and applying it to a different but related problem. For example, knowledge gained while learning to recognize cars could apply when trying to recognize trucks."

Image Embeddings and Transfer Learning

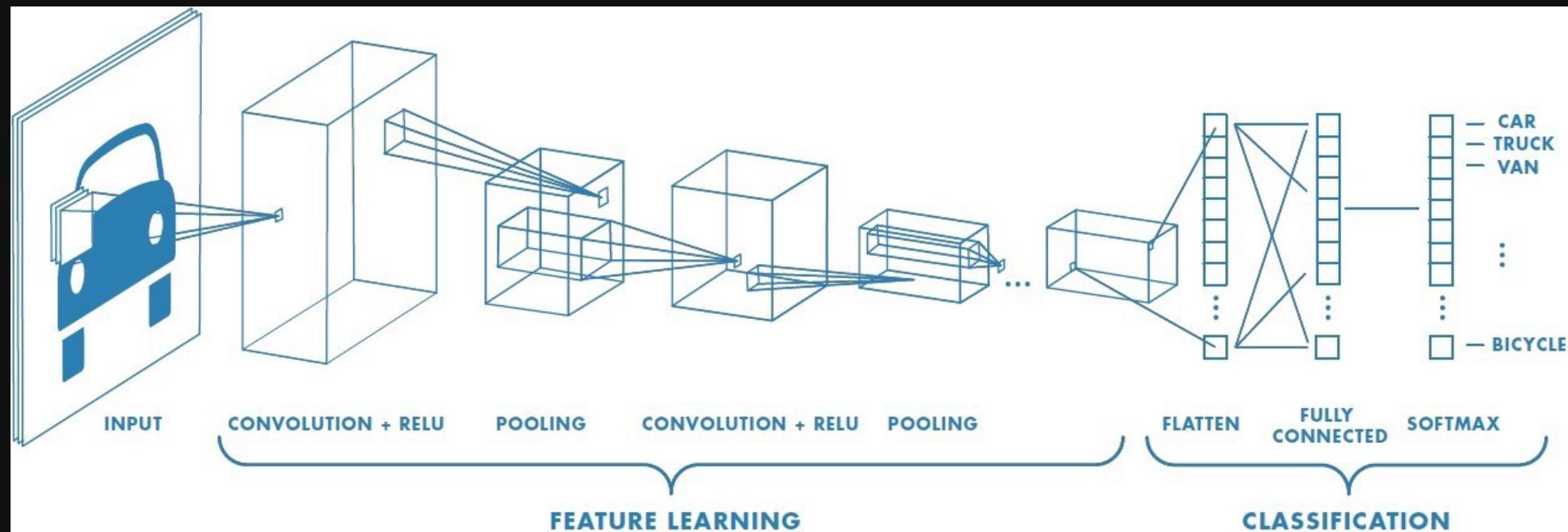


Image Embeddings and Transfer Learning

Cosine Similarity

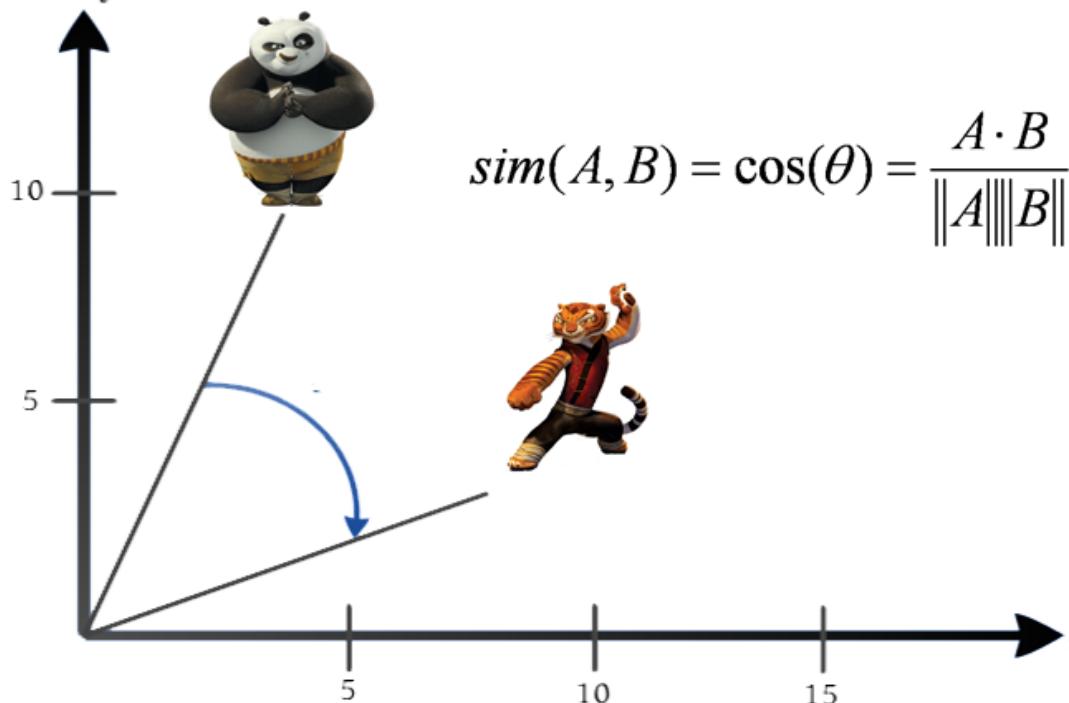


Image Embeddings

Image Embedding Class

```
1 ImageEmbedder():
2     def __init__(self):
3         self._IMAGE_TARGET_SIZE = (224, 224)
4
5         # Create an instance of a pretrained resnet model
6         # include_top to include the 2048x1 embedding layer
7
8         pre_trained_model = resnet_v2.ResNet50V2(
9             input_shape = (224, 224, 3),
10            weights = 'imagenet',
11            include_top=True
12        )
13
14         # Create a model that outputs only the second to last la
15         self.embedding_model = Model(
16             inputs=pre_trained_model.input
```

Image Embeddings

Image Embedding Class

```
3     self._IMAGE_TARGET_SIZE = (224, 224)
4
5     # Create an instance of a pretrained resnet model
6     # include_top to include the 2048x1 embedding layer
7
8     pre_trained_model = resnet_v2.ResNet50V2(
9         input_shape = (224, 224, 3),
10        weights = 'imagenet',
11        include_top=True
12    )
13
14    # Create a model that outputs only the second to last la
15    self.embedding_model = Model(
16        inputs=pre_trained_model.input,
17        outputs=pre_trained_model.get_layer('avg_pool').outp
18    )
```

Image Embeddings

Image Embedding Class

```
10         weights = 'imagenet',
11         include_top=True
12     )
13
14     # Create a model that outputs only the second to last la
15     self.embedding_model = Model(
16         inputs=pre_trained_model.input,
17         outputs=pre_trained_model.get_layer('avg_pool').outp
18     )
19
20     def create_embedding(self, image_path):
21         img = image.load_img(
22             image_path,
23             target_size=self._IMAGE_TARGET_SIZE
24         )
25
```

Image Embeddings

Image Embedding Class

```
13
14     # Create a model that outputs only the second to last la
15     self.embedding_model = Model(
16         inputs=pre_trained_model.input,
17         outputs=pre_trained_model.get_layer('avg_pool').outp
18     )
19
20     def create_embedding(self, image_path):
21         img = image.load_img(
22             image_path,
23             target_size=self._IMAGE_TARGET_SIZE
24         )
25
26         # Convert the image to an array
27         img_vec = image.img_to_array(img)
28         # The model expects batches of images in the shape
```

Image Embeddings

Image Embedding Class

```
22     image_path,
23     target_size=self._IMAGE_TARGET_SIZE
24 )
25
26     # Convert the image to an array
27     img_vec = image.img_to_array(img)
28     # The model expects batches of images in the shape
29     # (n,224,224,3), so add an extra dimension
30     img_vec = np.expand_dims(img_vec, axis=0)
31
32     # resnet_v2 has a built in function to preprocess
33     # images for the resnet model
34     img_vec = resnet_v2.preprocess_input(img_vec)
35
36     # Return the embedding from the function
37     return self.embedding_model.predict(img_vec)
```

Image Embeddings

Neo4j Import Class

```
1 import logging
2 import sys
3 from neo4j import GraphDatabase
4 from neo4j.exceptions import ServiceUnavailable
5
6 class NeoProdEmbeddings:
7     def __init__(self, uri, user, password):
8         self.driver = GraphDatabase.driver(uri, auth=(user, password))
9
10    def close(self):
11        # Don't forget to close the driver connection when you are
12        # done using it
13        self.driver.close()
14
15    @staticmethod
16    def enable_log(level, output_stream):
17        handler = logging.StreamHandler(output_stream)
18        logger = logging.getLogger('neo4j')
19        logger.setLevel(level)
20        logger.addHandler(handler)
```

Image Embeddings

Neo4j Import Class

```
16     handler = logging.StreamHandler(output_stream)
17     handler.setLevel(level)
18     logging.getLogger("neo4j").addHandler(handler)
19     logging.getLogger("neo4j").setLevel(level)
20
21     def add_prod_embedding(self, prod_id, embedding):
22         with self.driver.session(database='products') as session:
23             # Write transactions allow the driver to handle retries
24             result = session.write_transaction(
25                 self._add_and_return_embedding, prod_id, embedding
26             )
27
28     for row in result:
29         print(f"Created embedding for product ID {prod_id}")
30
31     @staticmethod
```

Image Embeddings

Neo4j Import Class

```
28 for row in result:
29     print(f"Created embedding for product ID {prod_id}")
30
31     @staticmethod
32     def _add_and_return_embedding(tx, prod_id, embedding):
33         query = (
34             "MERGE (prod:Product { code: $prod_id }) "
35             "ON MATCH set prod.embedding = $embedding "
36             "RETURN prod.code as prod_code, prod.embedding[0"
37     )
38
39     result = tx.run(query, prod_id=prod_id, embedding=embeddin
40
41     try:
42         return [
43     {
```

Image Embeddings

Neo4j Import Class

```
1 import logging
2 import sys
3 from neo4j import GraphDatabase
4 from neo4j.exceptions import ServiceUnavailable
5
6 class NeoProdEmbeddings:
7     def __init__(self, uri, user, password):
8         self.driver = GraphDatabase.driver(uri, auth=(user, password))
9
10    def close(self):
11        # Don't forget to close the driver connection when you are
12        self.driver.close()
13
14    @staticmethod
15    def enable_log(level, output_stream):
16        handler = logging.StreamHandler(output_stream)
```

Image Embeddings

Running the code

```
1 from embed_functions import image_embedder, neo_helper
2 from pathlib import Path
3
4 if __name__=='__main__':
5     uri = "bolt://localhost:7687"
6     user = 'neo4j'
7     password = 'password'
8
9 # Instantiate each of the classes we made
10 img_emb = image_embedder.ImageEmbedder()
11 neo_emb = neo_helper.NeoProdEmbeddings(uri, user, password)
12
13 IMAGE_DIR = '/Users/grantbeasley/Downloads/images/'
14 images = Path(IMAGE_DIR)
15
16 # Find all of the image files
```

Image Embeddings

Running the code

```
3
4 if __name__=='__main__':
5     uri = "bolt://localhost:7687"
6     user = 'neo4j'
7     password = 'password'
8
9 # Instantiate each of the classes we made
10 img_emb = image_embedder.ImageEmbedder()
11 neo_emb = neo_helper.NeoProdEmbeddings(uri, user, password)
12
13 IMAGE_DIR = '/Users/grantbeasley/Downloads/images/'
14 images = Path(IMAGE_DIR)
15
16 # Find all of the image files
17
18 all_files = [
```

Image Embeddings

Running the code

```
16 # Find all of the image files
17
18 all_files = [
19     image for folder in images.iterdir()
20     if not folder.name == '.DS_Store'
21     for image in folder.iterdir()
22     if image.suffix == '.jpg'
23 ]
24
25 # Iterate through each image and upload to the database
26 for file in all_files:
27     prod_id = file.name.replace('.jpg', '')
28     print(prod_id)
29     embedding = img_emb.create_embedding(file)
30     neo_emb.add_prod_embedding(prod_id, embedding[0].tolist())
```

Finding similar items by image

```
1 // Use appropriate database
2 :use products;
3
4 // Create graph projection with node properties
5 CALL gds.graph.project(
6   'knnGraph', {
7     Product: {
8       label: 'Product',
9     properties:'embedding'
10    }
11  }, '*');
12
13 // Create the top 10 nearest neighbours based on the image em
14 CALL gds.beta.knn.write(
15   'knnGraph', {
16     writeRelationshipType: 'SIMILAR'
```

Finding similar items by image

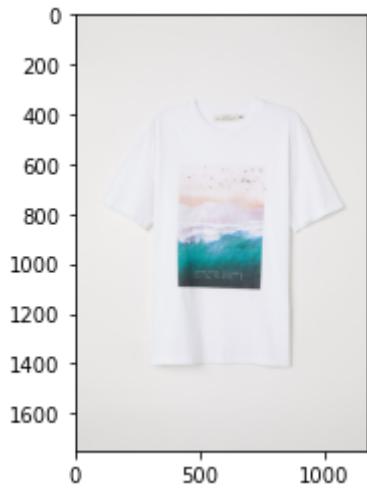
```
1 // Use appropriate database
2 :use products;
3
4 // Create graph projection with node properties
5 CALL gds.graph.project(
6   'knnGraph', {
7     Product: {
8       label: 'Product',
9     properties:'embedding'
10    }
11  }, '*');
12
13 // Create the top 10 nearest neighbours based on the image em
14 CALL gds.beta.knn.write(
15   'knnGraph', {
16     writeRelationshipType: 'SIMILAR'
```

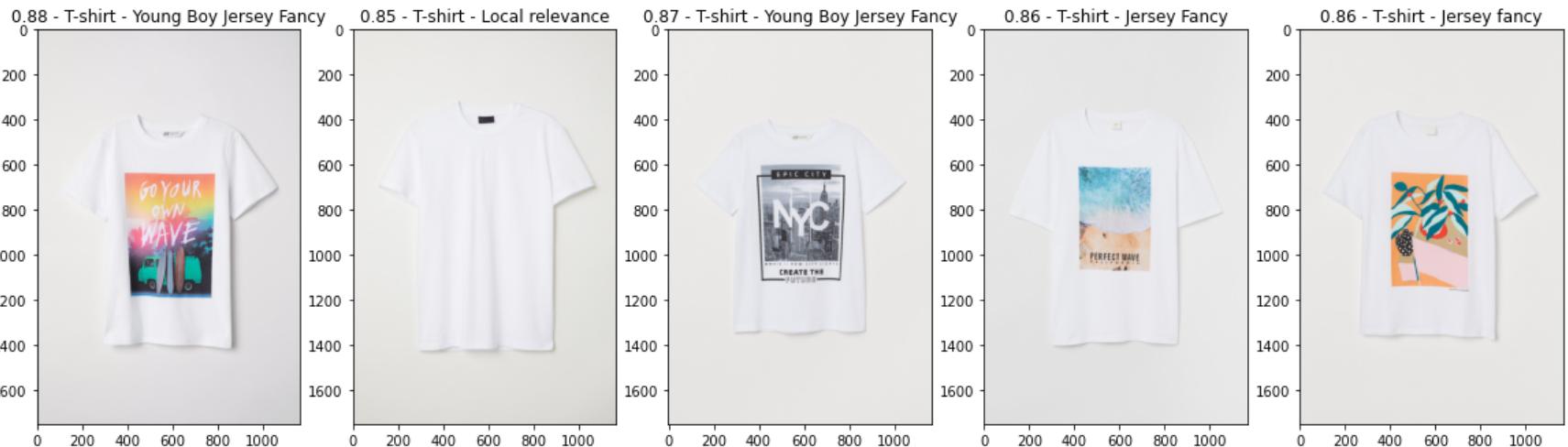
Finding similar items by image

```
10      },
11  }, '*');
12
13 // Create the top 10 nearest neighbours based on the image em
14 CALL gds.beta.knn.write(
15   'knnGraph', {
16     writeRelationshipType: 'SIMILAR',
17     writeProperty: 'score',
18     topK: 10,
19     randomSeed: 42,
20     concurrency: 1,
21     nodeWeightProperty: 'embedding'
22   }
23 )
24
25 VTFID nodesCompared relationshipsWritten;
```

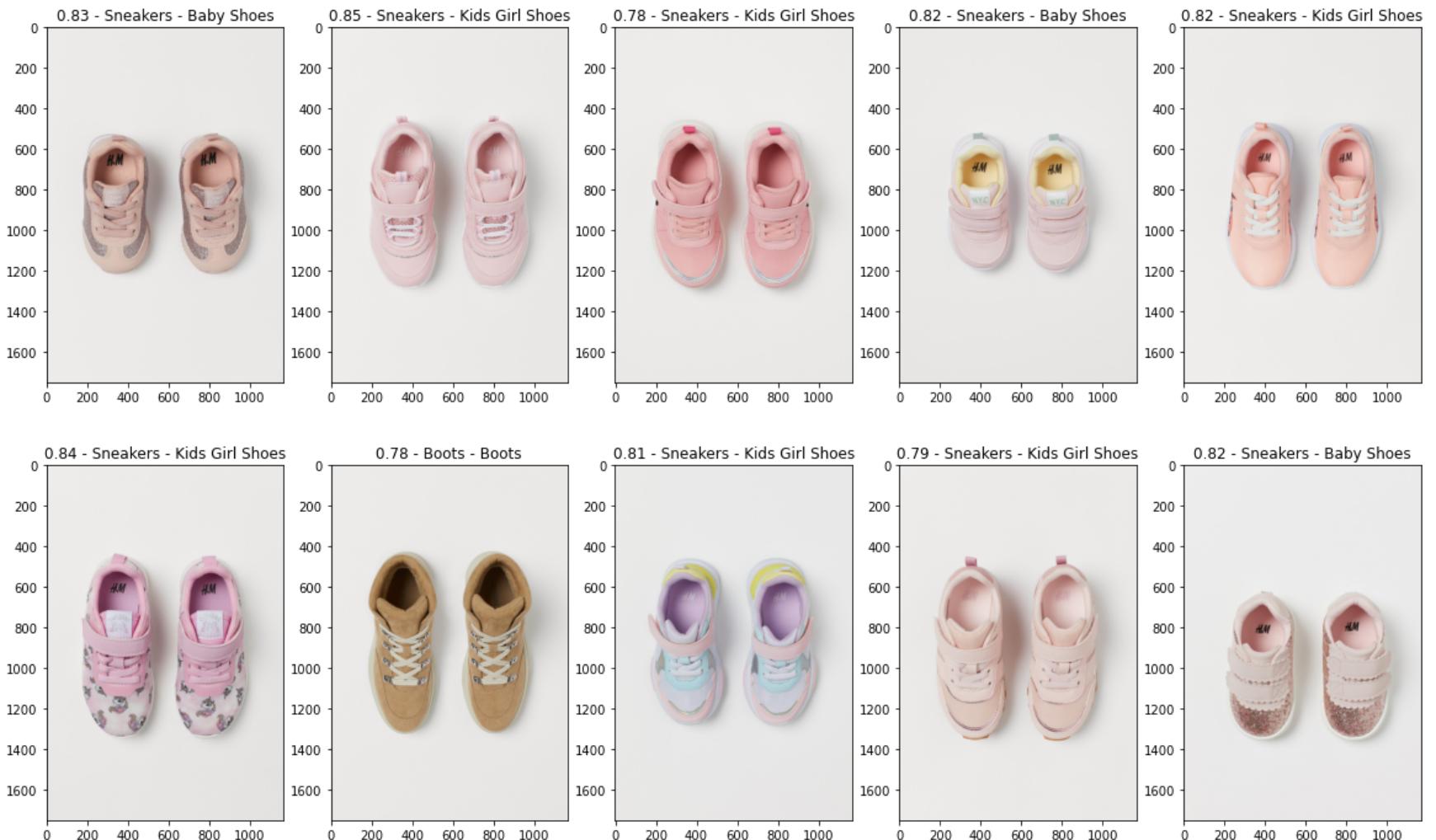
Randomly sampling data for quick validation

```
1 CALL {  
2   USE fabric.products  
3   MATCH (p:Product {code: $prod_code})-[s:SIMILAR]->(p2:Produc  
4   RETURN p2.code as prod_code, s.score as sim_score  
5 }  
6 CALL {  
7   USE fabric.neo4j  
8   WITH prod_code  
9   MATCH (d:Department)<-[:FROM_DEPARTMENT]-(p:Product {code:pr  
10  RETURN type.name as prod_type, d.name as department_name  
11 }  
12  
13 RETURN prod_code, sim_score, prod_type, department_name
```











Collaborative Filtering and Implicit Feedback

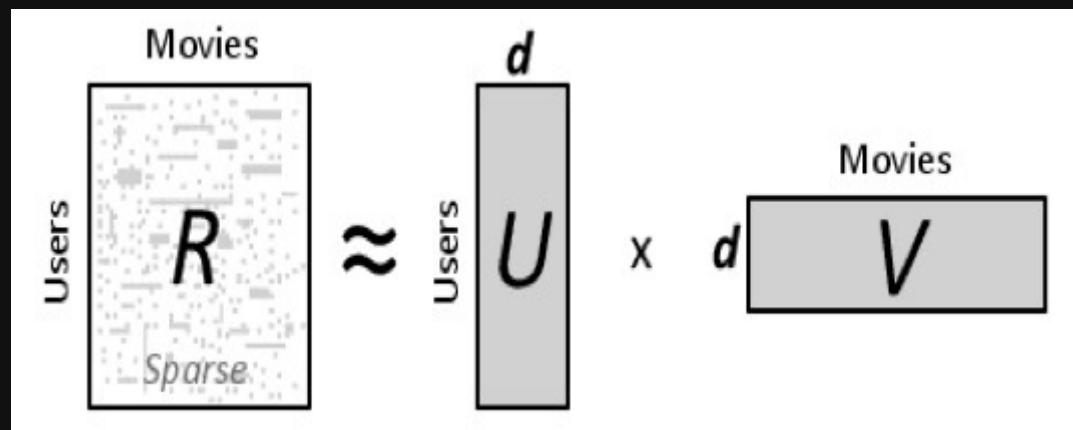


Implicit Feedback

- No negative feedback
- Implicit feedback is noisy
- Explicit feedback indicates preference, implicit feedback indicates confidence
- Need to use different measures

Alternating Least Squares

Python package - implicit



Considerations

- When are we recommending for?
- Sparsity of our dataset

Considerations

- When are we recommending for?
- Sparsity of our dataset

Num of users 1,362,281

Considerations

- When are we recommending for?
- Sparsity of our dataset

Num of users 1,362,281

Number of items 104,547

Considerations

- When are we recommending for?
- Sparsity of our dataset

Num of users 1,362,281

Number of items 104,547

Sparsity 0.022%

Dealing With Sparsity

"While implicit recommendations excel where data is sparse, it can often be helpful to make the interactions matrix a little more dense. We limited our data collection to models that had at least 5 likes. However, it may not be the case that every user has liked at least 5 models. Let's go ahead and knock out users that liked fewer than 5 models. This could possibly mean that some models end up with fewer than 5 likes once these users are knocked out, so we will have to alternate back and forth knocking users and models out until things stabilize."

Decreasing the sparsity of the dataset

```
1 def threshold_likes(df, user_min, article_min):
2     n_users = df['customer_id'].unique().shape[0]
3     n_items = df['article_id'].unique().shape[0]
4     sparsity = float(df.shape[0]) / float(n_users*n_items) * 1
5
6     done = False
7     while not done:
8         starting_shape = df.shape[0]
9         article_counts = df.groupby('customer_id')['article_id']
10        df = df[~df['customer_id'].isin(article_counts[article_c
11        user_counts = df.groupby('article_id')['customer_id'].co
12        df = df[~df['article_id'].isin(user_counts[user_counts <
13        ending_shape = df.shape[0]
14        if starting_shape == ending_shape:
15            done = True
16
```

Decreasing the sparsity of the dataset

```
4     sparsity = float(df.shape[0]) / float(n_users*n_items) * 1
5
6     done = False
7     while not done:
8         starting_shape = df.shape[0]
9         article_counts = df.groupby('customer_id')['article_id']
10        df = df[~df['customer_id'].isin(article_counts[article_c
11        user_counts = df.groupby('article_id')['customer_id'].co
12        df = df[~df['article_id'].isin(user_counts[user_counts <
13        ending_shape = df.shape[0]
14        if starting_shape == ending_shape:
15            done = True
16
17        assert(df.groupby('customer_id')['article_id'].count().min
18        assert(df.groupby('article_id')['customer_id'].count().min
19
```

Decreasing the sparsity of the dataset

```
18     assert df.groupby('article_id').customer_id.nunique().min() > 1
19
20 n_users = df['customer_id'].unique().shape[0]
21 n_items = df['article_id'].unique().shape[0]
22 sparsity = float(df.shape[0]) / float(n_users*n_items) * 100
23
24 return df
25
26 >>> Starting likes info
27 >>> Number of users: 1362281
28 >>> Number of items: 104547
29 >>> Sparsity: 0.022%
30 >>> Ending likes info
31 >>> Number of users: 925154
32 >>> Number of items: 91511
33 >>> Sparsity: 0.036%
```

Decreasing the sparsity of the dataset

```
1 def threshold_likes(df, user_min, article_min):
2     n_users = df['customer_id'].unique().shape[0]
3     n_items = df['article_id'].unique().shape[0]
4     sparsity = float(df.shape[0]) / float(n_users*n_items) * 1
5
6     done = False
7     while not done:
8         starting_shape = df.shape[0]
9         article_counts = df.groupby('customer_id')['article_id']
10        df = df[~df['customer_id'].isin(article_counts[article_c
11        user_counts = df.groupby('article_id')['customer_id'].co
12        df = df[~df['article_id'].isin(user_counts[user_counts <
13        ending_shape = df.shape[0]
14        if starting_shape == ending_shape:
15            done = True
16
```

We lose 32% of users and 12% of items!

Splitting the data

- Train test split
- Implications for number of purchases
- Precision @ k

Training the Model

Split the data and train on the training data.

```
als = implicit.als.AlternatingLeastSquares()  
als.fit(train_data)
```

Eyeballing the Recommendations



$$P = \frac{\# \text{ of our recommendations that are relevant}}{\# \text{ of items we recommended}}$$

Precision

$$P = \frac{\text{\# of our recommendations that are relevant}}{\text{\# of items we recommended}}$$

Precision

Test - Precision @ 5 3.3%

$$P = \frac{\text{\# of our recommendations that are relevant}}{\text{\# of items we recommended}}$$

Precision

Test - Precision @ 5 3.3%

Train - Precision @ 5 49.6%

Purchase Confidence - 2 Possible Approaches

Assume most recent purchases are most relevant
and assign more weight

```
def assign_confidence_recent(date):
    if date >= dt.date(2020, 3, 20):
        return 5
    elif (date >= dt.date(2019, 9, 20)) & (date < dt.date(2020, 3, 20))
        return 3
    else:
        return 1
```

Purchase Confidence - 2 Possible Approaches

We could assume that recent purchases are most relevant, and that purchases within the same season may be related (i.e. user preferences during autumn/winter may show different correlations to spring/summer).

```
def assign_confidence(date):
    if date >= dt.date(2020, 6, 20):
        return 5
    elif date.month in (7, 8, 9):
        return 3
    else:
        return 1
```

Results

Results

Binary
Confidence

Confidence 1

Confidence 2

Results

Binary Confidence	Confidence 1	Confidence 2
Train 49.6%	48.7%	48.7%

Results

	Binary Confidence	Confidence 1	Confidence 2
Train	49.6%	48.7%	48.7%
Test	3.3%	10.5%	10.4%

Bringing it all together

What we currently have:

- A graph model
- Image embeddings and similarities
- A matrix factorization based recommendation system

Challenges

- Image Embeddings - powerful but has limitations
- Matrix factorization method - decreasing sparsity of the dataset also introduces limitations

How we'll combine the various methods

How we'll combine the various methods

1. For each customer, find all purchases and retrieve most similar products based on images

How we'll combine the various methods

1. For each customer, find all purchases and retrieve most similar products based on images
2. Use the graph model and Cypher to filter list and ensure it only contains products from a department the customer has previously purchased from

How we'll combine the various methods

1. For each customer, find all purchases and retrieve most similar products based on images
2. Use the graph model and Cypher to filter list and ensure it only contains products from a department the customer has previously purchased from
3. Use the ALS model to rank all these items and take the top 12

How we'll combine the various methods

1. For each customer, find all purchases and retrieve most similar products based on images
2. Use the graph model and Cypher to filter list and ensure it only contains products from a department the customer has previously purchased from
3. Use the ALS model to rank all these items and take the top 12
4. For those people with less than 12 recommended items, take a random sample from the items from the Cypher query

How we'll combine the various methods

1. For each customer, find all purchases and retrieve most similar products based on images
2. Use the graph model and Cypher to filter list and ensure it only contains products from a department the customer has previously purchased from
3. Use the ALS model to rank all these items and take the top 12
4. For those people with less than 12 recommended items, take a random sample from the items from the Cypher query
5. For those still without 12 recommendations, use the graph model to return up to 12 most popular products purchased by other users who also purchased products by the current customer

Creating the Sparse Matrix

```
1 # Create mappings for sparse matrix
2 all_product_to_idx = {}
3 all_idx_to_product = {}
4 for (idx, prod_id) in enumerate(products):
5     all_product_to_idx[prod_id] = idx
6     all_idx_to_product[idx] = prod_id
7
8 all_customer_to_idx = {}
9 all_idx_to_customer = {}
10 for (idx, cust_id) in enumerate(customers):
11     all_customer_to_idx[cust_id] = idx
12     all_idx_to_customer[idx] = cust_id
13
14 # Create lil matrix to incrementally build sparse matrix of re
15 n_users = len(customers)
16 n_products = len(products)
```

The Recommendation Code

```
1 CALL {  
2 // Return all purchases for a single customer  
3 USE fabric.neo4j  
4 MATCH (c:Customer)-[pur:PURCHASED]->(p:Product)  
5 WHERE c.id = $customer_id  
6 RETURN p.code as prod_code  
7 }  
8 CALL {  
9 // For each of those purchases, return the 10 most similar  
10 // by image embedding  
11 USE fabric.products  
12 WITH prod_code  
13 MATCH (p1:Product)-[:SIMILAR]->(rec:Product)  
14 WHERE p1.code = prod_code  
15 RETURN rec.code as rec_code  
16 }
```

The Recommendation Code

```
1 CALL {  
2 // Return all purchases for a single customer  
3 USE fabric.neo4j  
4 MATCH (c:Customer)-[pur:PURCHASED]->(p:Product)  
5 WHERE c.id = $customer_id  
6 RETURN p.code as prod_code  
7 }  
8 CALL {  
9 // For each of those purchases, return the 10 most similar  
10 // by image embedding  
11 USE fabric.products  
12 WITH prod_code  
13 MATCH (p1:Product)-[:SIMILAR]->(rec:Product)  
14 WHERE p1.code = prod_code  
15 RETURN rec.code as rec_code  
16 }
```

The Recommendation Code

```
6   RETURN p.code as prod_code
7 }
8 CALL {
9 // For each of those purchases, return the 10 most similar
10 // by image embedding
11 USE fabric.products
12 WITH prod_code
13 MATCH (p1:Product)-[:SIMILAR]->(rec:Product)
14 WHERE p1.code = prod_code
15 RETURN rec.code as rec_code
16 }
17 CALL {
18 // For each of these recommendations, filter out the ones
19 // that do not come from a department or index previously
20 // purchases from by the customer
```

The Recommendation Code

```
18 // For each of these recommendations, filter out the ones
19 // that do not come from a department or index previously
20 // purchases from by the customer

21 USE fabric.neo4j
22 with rec_code
23 MATCH (c:Customer)
24 WHERE c.id = $customer_id
25 MATCH (ind:Index)<- [:HAS_INDEX] - (rec:Product) - [:FROM_DEPARTMENT] -
26 WHERE rec.code = rec_code
27 AND EXISTS ((c)-[:PURCHASED]->(:Product) - [:FROM_DEPARTMENT] - )
28 AND EXISTS ((c)-[:PURCHASED]->(:Product) - [:HAS_INDEX]->(ind))
29 RETURN rec.code AS recommended_item
30 }
31
32 RETURN collect(recommended_item) as recommended_items
```

The Recommendation Code

Write all recommendations to a sparse matrix

```
1 for i, cust in enumerate(customers):
2     customer_index = customer_to_idx[cust]
3     recs = neo_recs.get_recommended_items_for_user(cust)
4     rec_indices = [product_to_idx[rec] for rec in recs]
5     rec_matrix[customer_index, rec_indices] = 1
```









MakeAG

The Final Part

The Final Part

1. For each customer, get the recommendations based off of image similarity and filtered with the graph model

The Final Part

1. For each customer, get the recommendations based off of image similarity and filtered with the graph model
2. Check they were part of our ALS model dataset - if they were, then rank all of the recs selected from part 1

The Final Part

1. For each customer, get the recommendations based off of image similarity and filtered with the graph model
2. Check they were part of our ALS model dataset - if they were, then rank all of the recs selected from part 1
3. If they were not, or if less than 12 items returned, randomly sample our recommendations from part 1

The Final Part

1. For each customer, get the recommendations based off of image similarity and filtered with the graph model
2. Check they were part of our ALS model dataset - if they were, then rank all of the recs selected from part 1
3. If they were not, or if less than 12 items returned, randomly sample our recommendations from part 1
4. Finally, if a customer has zero recommendations, use pure Cypher query

If images/ALS don't have the answer...

Cypher query for those with few items from previous query. Only used as backup query.

```
1 MATCH (c:Customer {id: $customer_id})-[:PURCHASED]->(p:Product)
2 [:PURCHASED]-(:Customer)-[:PURCHASED]->(rec:Product)
3 WHERE id(p) <> id(rec)
4 AND NOT EXISTS ((c)-[:PURCHASED]->(rec))
5 WITH c.id as customer_id, rec, COUNT(rec) as score
6 ORDER BY COUNT(rec) DESC LIMIT 12
7 RETURN collect(rec.code) as recommended_items
```

If images/ALS don't have the answer...

Cypher query for those with few items from previous query. Only used as backup query.

```
1 MATCH (c:Customer {id: $customer_id})-[:PURCHASED]->(p:Product)
2 [:PURCHASED]-(:Customer)-[:PURCHASED]->(rec:Product)
3 WHERE id(p) <> id(rec)
4 AND NOT EXISTS ((c)-[:PURCHASED]->(rec))
5 WITH c.id as customer_id, rec, COUNT(rec) as score
6 ORDER BY COUNT(rec) DESC LIMIT 12
7 RETURN collect(rec.code) as recommended_items
```

If images/ALS don't have the answer...

Cypher query for those with few items from previous query. Only used as backup query.

```
1 MATCH (c:Customer {id: $customer_id})-[:PURCHASED]->(p:Product)
2 [:PURCHASED]-(:Customer)-[:PURCHASED]->(rec:Product)
3 WHERE id(p) <> id(rec)
4 AND NOT EXISTS ((c)-[:PURCHASED]->(rec))
5 WITH c.id as customer_id, rec, COUNT(rec) as score
6 ORDER BY COUNT(rec) DESC LIMIT 12
7 RETURN collect(rec.code) as recommended_items
```

If images/ALS don't have the answer...

Cypher query for those with few items from previous query. Only used as backup query.

```
1 MATCH (c:Customer {id: $customer_id})-[:PURCHASED]->(p:Product)
2 [:PURCHASED]-(:Customer)-[:PURCHASED]->(rec:Product)
3 WHERE id(p) <> id(rec)
4 AND NOT EXISTS ((c)-[:PURCHASED]->(rec))
5 WITH c.id as customer_id, rec, COUNT(rec) as score
6 ORDER BY COUNT(rec) DESC LIMIT 12
7 RETURN collect(rec.code) as recommended_items
```

Eyeballing The Recommendations Part 2



What would I do differently?

- Change the graph model/fabric setup
- Optimize ALS model
- Optimize the inference process
- Utilise more of the customer data



© 2000