

# Data Exploration

## Data Set Overview

The table below lists each of the files available for analysis with a short description of what is found in each one.

File Name	Description	Fields
<b>ad-clicks.csv</b>	A detailed description of players' clicks on advertisements	<b>timestamp</b> : when the click occurred. <b>txId</b> : a unique id assigned to the click <b>userSessionId</b> : the session id of the session in which the user made the click <b>teamid</b> : the current team id of the user who made the click <b>userid</b> : the id of the user who made the click <b>adId</b> : the id of the ad clicked on <b>adCategory</b> : the category of the ad clicked on
<b>buy-clicks.csv</b>	A detailed record of players' in-app purchases	<b>timestamp</b> : when the purchase was made. <b>txId</b> : a unique id assigned to the purchase <b>userSessionId</b> : the session id of the session in which the user made the purchase <b>team</b> : the current team id of the user who made the purchase <b>userId</b> : the id of the user who made the purchase <b>buyId</b> : the id of the item purchased <b>price</b> : the price of the item purchased
<b>users.csv</b>	This file contains information about the game's players.	<b>timestamp</b> : when user first played the game.

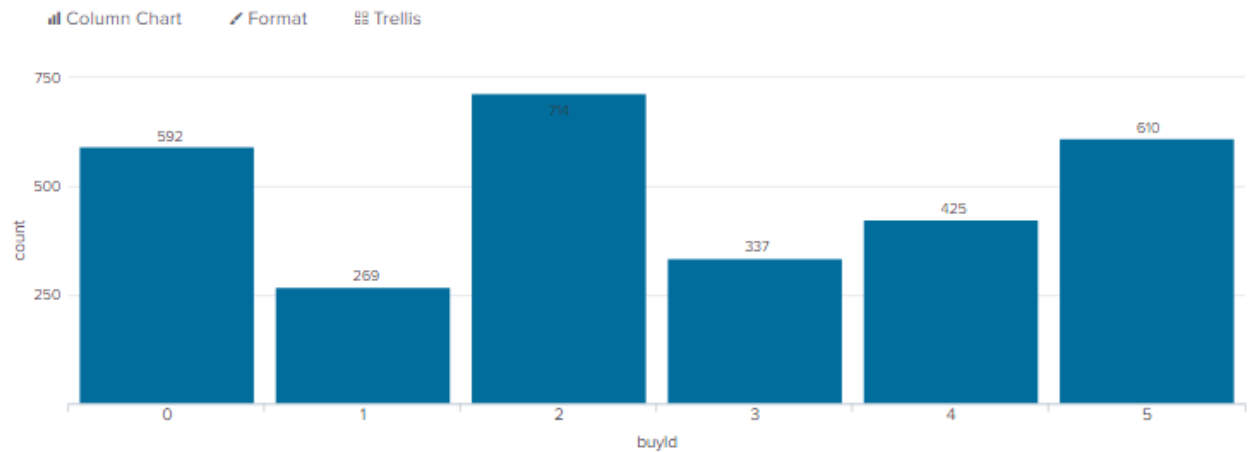
		<p><b>userId</b>: a unique id assigned to the user.</p> <p><b>nick</b>: the nickname chosen by the user.</p> <p><b>twitter</b>: the twitter handle of the user.</p> <p><b>dob</b>: the date of birth of the user.</p> <p><b>country</b>: the two-letter country code where the user lives.</p>
<b>team.csv</b>	A detailed record of terminated teams	<p><b>teamId</b>: the id of the terminated team</p> <p><b>name</b>: the name of the terminated team</p> <p><b>teamCreationTime</b>: when the team was created</p> <p><b>teamEndTime</b>: when the last member left the team</p> <p><b>strength</b>: a measure of team strength at point of termination</p> <p><b>currentLevel</b>: the level of the team at point of termination</p>
<b>team-assignments.csv</b>	A record of players' team-joining events	<p><b>timestamp</b>: when the user joined the team.</p> <p><b>team</b>: the id of the team that the user joined</p> <p><b>userId</b>: the id of the user</p> <p><b>assignmentId</b>: a unique id assigned to this joining event</p>
<b>level-events.csv</b>	A detailed record of every time a team engages with levels	<p><b>timestamp</b>: when the engagement occurred.</p> <p><b>eventId</b>: a unique id assigned to the engagement</p> <p><b>teamId</b>: the id of the team</p> <p><b>teamLevel</b>: the level engaged by the team</p> <p><b>eventType</b>: the type of event, either 'start' or 'end'</p>

<b>user-session.csv</b>	A detailed record of every time a user starts/ends a playing session	<p><b>timestamp</b>: when the event occurred.</p> <p><b>userSessionId</b>: a unique id assigned to the session.</p> <p><b>userId</b>: the current user's ID.</p> <p><b>teamId</b>: the id of the user's current team.</p> <p><b>assignmentId</b>: the team assignment id assigned when the user joined the current team</p> <p><b>sessionType</b>: whether the event type is 'start' or 'end'</p> <p><b>teamLevel</b>: the level of the team during this session.</p> <p><b>platformType</b>: the platform type the user used during the session.</p>
<b>game-clicks.csv</b>	A detailed description of all clicks performed by the user during playing	<p><b>timestamp</b>: when the click occurred.</p> <p><b>clickId</b>: a unique id assigned to the click.</p> <p><b>userId</b>: the id of the user performing the click.</p> <p><b>userSessionId</b>: the id of the session the user was in when the click occurred.</p> <p><b>isHit</b>: Boolean value denotes if the click hit the flamingo (value is 1) or missed the flamingo (value is 0)</p> <p><b>teamId</b>: the id of the user's current team</p> <p><b>teamLevel</b>: the current level of the user's current team</p>

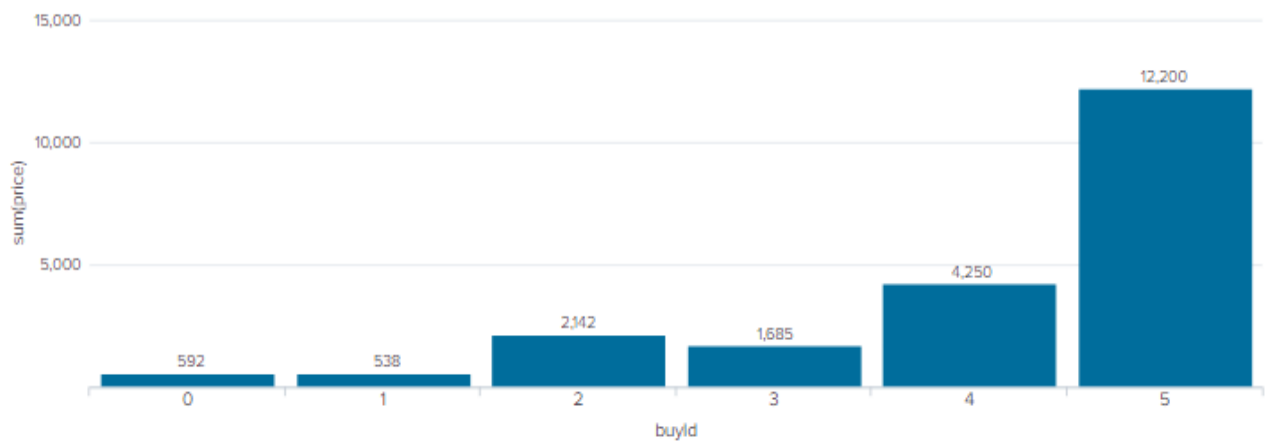
## Aggregation

Amount spent buying items	21407.0
Number of unique items available to be purchased	6

A histogram showing how many times each item is purchased:

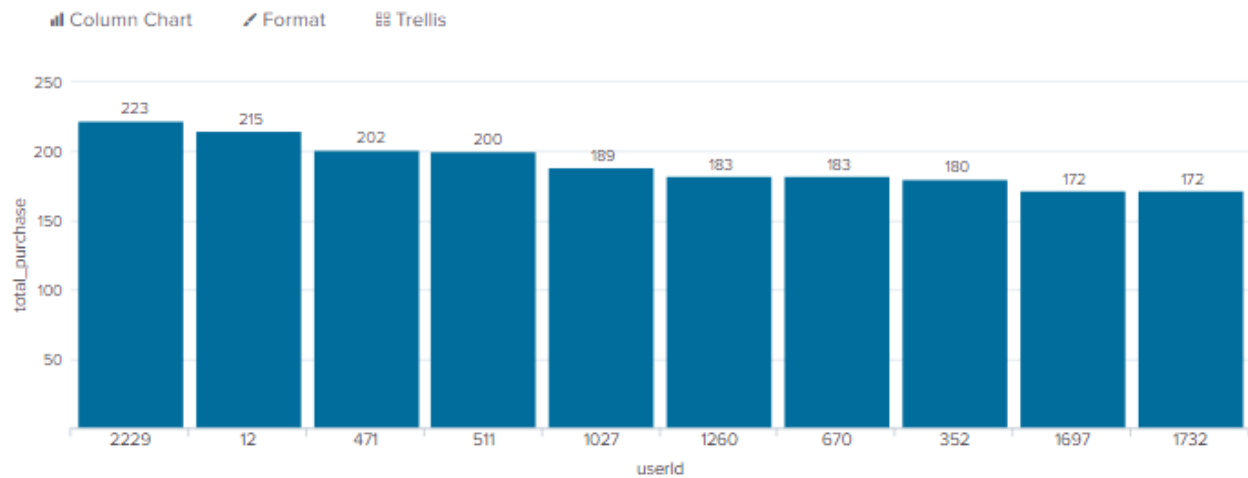


A histogram showing how much money was made from each item:



## Filtering

A histogram showing total amount of money spent by the top ten users (ranked by how much money they spent).



The following table shows the user id, platform, and hit-ratio percentage for the top three buying users:

Rank	User Id	Platform	Hit-Ratio (%)
1	2229	iPhone	11.60
2	12	iPhone	13.07
3	471	iPhone	14.50

## Data Preparation

Analysis of combined\_data.csv

### Sample Selection

Item	Amount
# of Samples	4619
# of Samples with Purchases	1411

\*\*\*Number of samples without purchases = 3208

### Attribute Creation

A new categorical attribute was created to enable analysis of players as broken into 2 categories (HighRollers and PennyPinchers). A screenshot of the attribute follows:

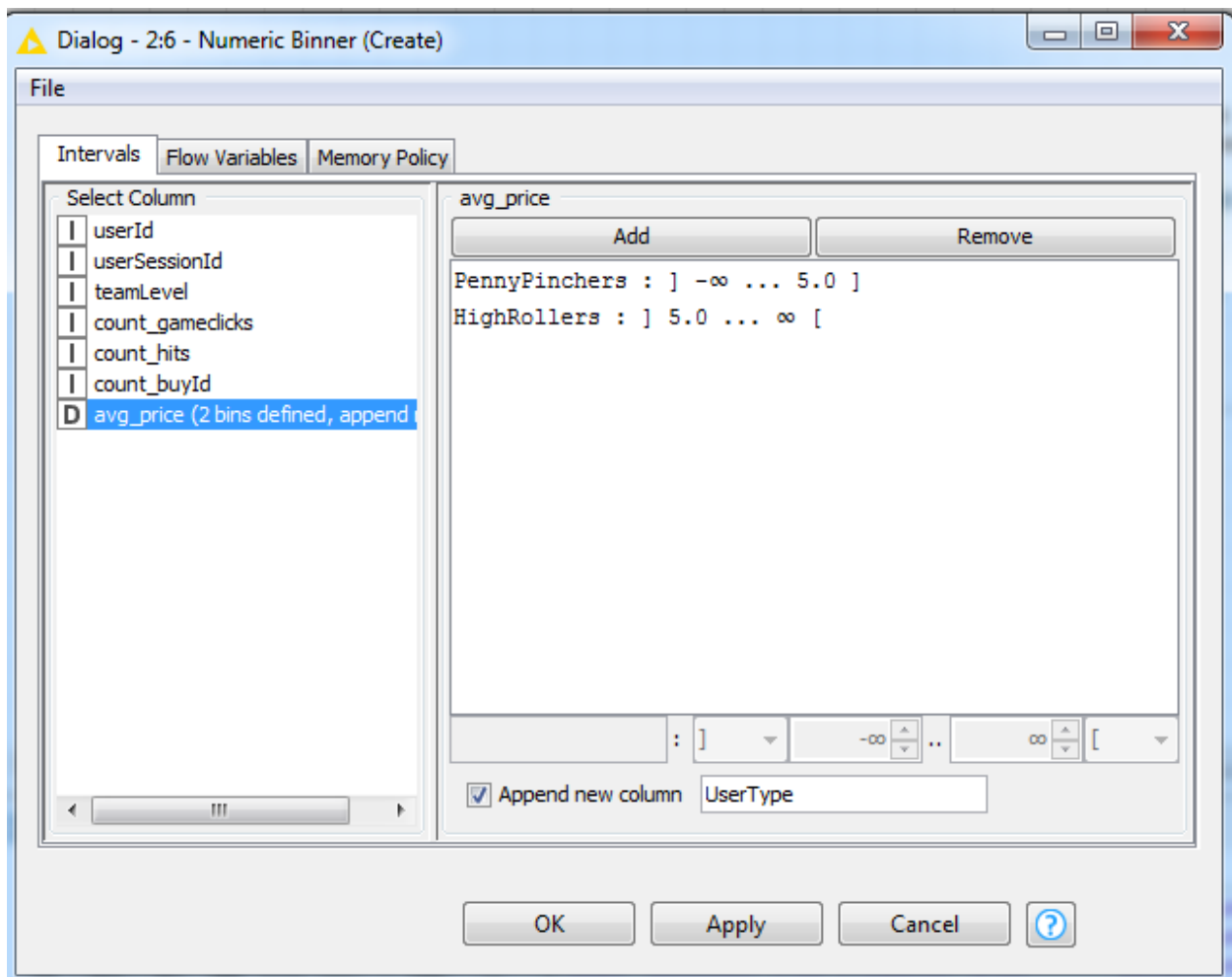


Table "default" - Rows: 1411						
Spec - Columns: 6			Properties	Flow Variables		
Row ID	I teamLevel	S platfor...	I count_...	I count_...	I count_...	S UserType
Row4	1	android	39	0	1	PennyPinchers
Row11	1	iphone	129	9	1	HighRollers
Row13	1	android	102	14	1	PennyPinchers
Row17	1	android	39	4	1	PennyPinchers
Row18	1	android	90	10	1	PennyPinchers
Row31	1	iphone	51	8	1	HighRollers
Row49	1	android	51	6	2	PennyPinchers
Row50	1	android	47	5	2	PennyPinchers
Row58	1	android	46	7	1	PennyPinchers
Row61	1	iphone	41	6	1	HighRollers
Row68	1	android	47	7	1	PennyPinchers
Row72	1	iphone	76	7	1	HighRollers
Row73	1	android	52	2	1	PennyPinchers
Row101	1	android	62	9	1	PennyPinchers
Row122	1	iphone	177	25	2	HighRollers
Row127	1	iphone	54	5	1	HighRollers
Row129	1	android	27	4	2	PennyPinchers
Row131	1	iphone	37	2	1	HighRollers
Row135	1	android	67	5	1	PennyPinchers
Row137	1	iphone	37	5	2	HighRollers

Describe the design of your attribute in 1-3 sentences.

- New column named **"UserType"** was added
- **PennyPinchers** have avg\_price ≤ 5.0\$. Colored in Red (first bin)
- **HighRollers** have avg\_price > 5.0\$. Colored in Green (second bin)

The creation of this new categorical attribute was necessary, because

- This new category is the **target variable** used for **data labelling of a classification task**. A classification task needs discrete categories
- For a **supervised learning classification task** such as our current task, labels are required during **model training**
- The **model score** is also derived from **comparing predicted labels & actual labels** of the **test set**

### Attribute Selection

The following attributes were filtered from the dataset for the following reasons:

Attribute	Rationale for Filtering
UserID	Assigned <b>Randomly</b> by system. Has no relationship to a user's in-game behavior
SessionID	Assigned <b>Randomly</b> by system. Has no relationship to a user's in-game behavior
Avg_price	<p>This is the column we <b>derive target variable from</b> → it has <b>100% correlation</b> to the target variable</p> <p>We get rid of this feature because we already have the <b>UserType column</b> which acts as <b>labels for the classification task</b></p>

\*Remaining features = Team\_level, platformType, count\_clicks, count\_ishits, count\_buyid



## Data Partitioning and Modeling

The data was partitioned into train and test datasets.

The **train data set** was used to create the decision tree model.

The trained model was then applied to **the test dataset**.

This is important because a train model needs to be **evaluated/validated** using **data it has not seen before or touched during training phase**. Such data comes from the test set

When partitioning the data using sampling, it is important to set the random seed because this ensures that the result is **reproduce-able after multiple iterations** of the KNIME workflow and student's result **matches** the model answer

A screenshot of the resulting decision tree can be seen below:



## Evaluation

A screenshot of the confusion matrix can be seen below:

UserType \...	PennyPinc...	HighRollers
PennyPinchers	308	27
HighRollers	38	192

As seen in the screenshot above, the overall accuracy of the model is

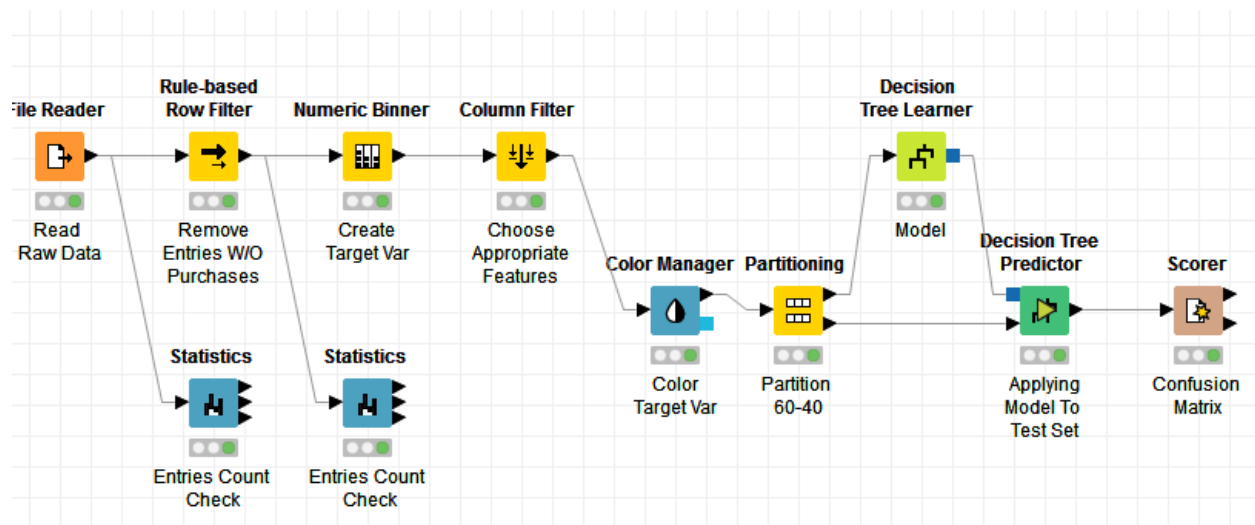
$$\frac{308 + 192}{308 + 27 + 38 + 192} = 88.5\%$$

Write one sentence for each of the values of the confusion matrix indicating what has been correctly or incorrectly predicted.

- 308 of Penny Pinchers in the test set were correctly identified
- 192 of High Rollers in the test set were correctly identified
- 27 of Penny Pinchers in the test set were incorrectly identified as High Rollers
- 38 of High Rollers in the test set were incorrectly identified as Penny Pinchers

## Analysis Conclusions

The final KNIME workflow is shown below:



What makes a HighRoller vs. a PennyPincher?

- High Rollers:
  - Iphone users
- Penny Pinchers
  - The remaining platforms' users

Specific Recommendations to Increase Revenue
1. Promote <b>special promotions</b> of in-app purchases & <b>increase number of ads</b> to users of platforms other than Iphone
2. Send <b>thank you email</b> to Iphone users for their purchases and <b>send vouchers</b> to them to encourage them to keep on spending

## Attribute Selection

Attribute	Rationale for Selection
<b>Purchase_per_adclick</b>	<p>Average purchase amount per ad click. Calculated by total payment divided by total number of ad clicks for that user.</p> <p>This is an important metric (usually called conversion) to measure revenue from each user</p>
<b>Avg_session_duration</b>	<p>Average playing session length. Calculated by taking the time delta between session start and session end for each session ID. The find the average of this duration for each user ID</p> <p>To distinguish between casual players and hardcore players that play the game at length</p>
<b>Hit_rate</b>	<p>Number of hits that hits the flamingo/Total number of hits</p> <p>To distinguish between highly skilled players and lousy ones</p>

## Training Data Set Creation

The training data set used for this analysis is shown below (first 5 lines):

purchase_per_adclick	avg_session_duration	hit_rate
0.000000	2832.500000	0.105535
0.477273	3750.988095	0.134078
0.000000	3834.000000	0.095238
0.000000	4470.000000	0.105960
5.300000	2400.000000	0.100000

Dimensions of the training data set (rows x columns): **1193 x 3** (Not including the index column)

# of clusters created: **4**

## Cluster Centers

Cluster #	Center (purchase_per_adclick, avg_session_length, hit_rate)
1	[-0.21903156, 0.43540912, -0.0538909 ]
2	[-0.32941065, -1.7305019 , -0.42341042]
3	[-0.46824539, -2.23297675, 4.92736511]
4	[2.3048937 , 0.09565116, 0.41877504]

These clusters can be differentiated from each other as follows:

Cluster 1 is different from the others in that these users have

- low conversion &
- **the longest playing time** &
- standard hit rate.
- These are '**lovers**' of the game

Cluster 2 is different from the others in that these users have

- low conversion &
- low hit rate &
- **the worst hit rate.**
- These are the '**least skillful**' of the game

Cluster 3 is different from the others in that these users have

- **the lowest conversion** &
- **the shortest playing time** &
- **the highest hit rate** (extremely high).
- These are '**assassins**' of the game

Cluster 4 is different from the others in that these users have

- **the highest conversion** &
- standard playing time &
- high hit rate.
- These are simply the '**willing spenders**' of the game

## Recommended Actions

Cluster	Action Recommended	Rationale for the action
Cluster 1 – ‘Lovers’	Simply present <b>more advertisements</b> to these players	They spend the longest time playing → should see more ads compared to other players
Cluster 2 – ‘Least skillful’	Present in-app items that are more related to <b>improving hit_rate</b> to these users	These users are least skillful and hence will appreciate these items
Cluster 3 – ‘Assassins’	<p><b>With in-app items:</b> The company can instead</p> <ul style="list-style-type: none"> <li>• Present items related to <b>avatar’s decorations</b></li> <li>• Present more <b>challenging quests</b> to these users and present items (that make the quests easier) inside these quests</li> <li>• Have more <b>discounts/promotions</b></li> </ul> <p><b>With third-party advertisements</b></p> <ul style="list-style-type: none"> <li>• Don’t bother, these don’t have time and is least likely to spend on ads</li> </ul>	<p>Present items that are related to improving hit_rate is useless to these ‘assassins’</p> <p>These players are good so may like avatar decorations to appeal to their ego</p> <p>Harder quests can encourage these players to start using purchased aids</p> <p>Little chance to profit with third party ads from these users</p>
Cluster 4 – ‘Willing spenders’	<p><b>With in-app items:</b></p> <ul style="list-style-type: none"> <li>• Promote <b>more expensive items</b> to these willing spenders</li> </ul> <p><b>With third-party advertisements</b></p> <ul style="list-style-type: none"> <li>• Present <b>more ads</b></li> <li>• <b>Charge the advertisers more</b> money to target these users</li> </ul>	<p>These users are the easiest to convert. Doing so will milk more revenue out of these players</p> <p>These are willing spenders! Advertisers need to pay more to access the gold mine</p>

## Graph Analytics

### Modeling Chat Data using a Graph Data Model

This is a Graph Model for chats created by users of the game “To Catch The Pink Flamingo”. The Graph Model has 4 types of nodes (**Users**, **Teams**, **Team Chat Sessions** and **Chat Items**)

Each node has its own unique ID. There are various types of edges, each represents a different kind of interaction between nodes

Edges such as ‘**CreatesSession**’, ‘Joins’ and ‘Leaves’ describes the interaction Users have on the Chat Sessions

Edges such as ‘**CreatesChat**’ describes the interaction Users have on the Chat Items

Edges such as ‘**PartOf**’ describes which Chat Sessions the Chat Items belong to

Edges such as ‘**OwnedBy**’ describes which Teams the Chat Sessions belong to

Edges such as ‘**Mentions**’ and ‘**ResponsesTo**’ describes the interaction ChatItems have towards Users or other Chat Items respectively

### Creation of the Graph Database for Chats

Describe the steps you took for creating the graph database. As part of these steps

- i) Write the schema of the 6 CSV files

Chat_create_team_chat	<ul style="list-style-type: none"><li>• userid: ID of user: Int</li><li>• teamid: ID of team: Int</li><li>• TeamChatSessionID: ID of chat session: Int</li><li>• timestamp: Timestamp</li></ul> <p>-- for userid <b>CREATING</b> TeamChatSessionID -- for TeamChatSessionID <b>OWNEDBY</b> teamid</p>
Chat_item_team_chat	<ul style="list-style-type: none"><li>• userid: ID of user: Int</li><li>• teamchatsessionid: ID of chat session: Int</li><li>• chatitemid: ID of chat item: Int</li><li>• timestamp: Timestamp</li></ul> <p>-- for userid <b>CREATING</b> chatitemid -- for chatitemid <b>PARTOF</b> teamchatsessionid</p>
Chat_join_team_chat	<ul style="list-style-type: none"><li>• userid: ID of user: Int</li><li>• TeamChatSessionID: ID of chat session: Int</li><li>• timestamp: Timestamp</li></ul> <p>-- for userid <b>JOINING</b> TeamChatSessionID</p>



Chat_leave_team_chat	<ul style="list-style-type: none"> <li>• userd: ID of user: Int</li> <li>• teamchatsessionid: ID of chat session: Int</li> <li>• timestamp: Timestamp</li> </ul> <p>-- for userid <b>LEAVING</b> teamchatsessionid</p>
Chat_mention_team_chat	<ul style="list-style-type: none"> <li>• ChatItem: ID of chat item: Int</li> <li>• userid: ID of user: Int</li> <li>• timestamp: Timestamp</li> </ul> <p>-- for ChatItem <b>MENTIONING</b> userid</p>
Chat_respond_team_chat	<ul style="list-style-type: none"> <li>• chatid1: ID of chat item: Int</li> <li>• chatid2: ID of chat item: Int</li> <li>• timestamp: Timestamp</li> </ul> <p>-- for chatid1 <b>RESPONDINGTO</b> chatid2</p>

ii) Explain the loading process and include a sample LOAD command

- The loading process determines **which columns** in the CSV file should be **used as integer ids** for the **nodes** and **which single column** should be used as **timestamp of the interaction**.
- Nodes are read as (**nodechar: Nodetype {id: toInteger(row[**row order**]})**). Bolded fields are filled in by the Neo4j user
- Edges are read as **-[: Interactiontype {timestamp: row[**order of timestamp column**]}]->**. Bolded fields are filled in by the Neo4j user
- The CSV **headers** reveal what **type of nodes** that is whereas the **type of interaction** has to be **guessed** via the **file name**. There can be **more than one** type of interaction
- If there are many types of interaction, the timestamp column is shared by them

```
1 LOAD CSV FROM "file:/path/to/chat_create_team_chat.csv" AS row
2 MERGE (u:User {id: toInteger(row[0])}) MERGE (t:Team {id: toInteger(row[1])})
3 MERGE (c:TeamChatSession {id: toInteger(row[2])})
4 MERGE (u)-[:CreatesSession{timeStamp: row[3]}]->(c)
5 MERGE (c)-[:OwnedBy{timeStamp: row[3]}]->(t)
```

```
1 LOAD CSV FROM "file:/path/to/chat_item_team_chat.csv" AS row
2 MERGE (u:User {id: toInteger(row[0])})
3 MERGE (c:TeamChatSession {id: toInteger(row[1])})
4 MERGE (m:ChatItem {id: toInteger(row[2])})
5 MERGE (u)-[:CreatesChat {timeStamp: row[3]}]->(m)
6 MERGE (m)-[:PartOf {timeStamp: row[3]}]->(c)
```

```
1 LOAD CSV FROM "file:/path/to/chat_join_team_chat.csv" AS row
2 MERGE (u:User {id: toInteger(row[0])})
3 MERGE (c:TeamChatSession {id: toInteger(row[1])})
4 MERGE (u)-[:Joins {timeStamp: row[2]}]->(c)
```

```

1 LOAD CSV FROM "file:/path/to/chat_leave_team_chat.csv" AS row
2 MERGE (u:User {id: toInteger(row[0])})
3 MERGE (c:TeamChatSession {id: toInteger(row[1])})
4 MERGE (u)-[:Leaves {timeStamp: row[2]}]->(c)

```

```

1 LOAD CSV FROM "file:/path/to/chat_mention_team_chat.csv" AS row
2 MERGE (m:ChatItem {id: toInteger(row[0])})
3 MERGE (u:User {id: toInteger(row[1])})
4 MERGE (m)-[:Mentions {timeStamp: row[2]}]->(u)

```

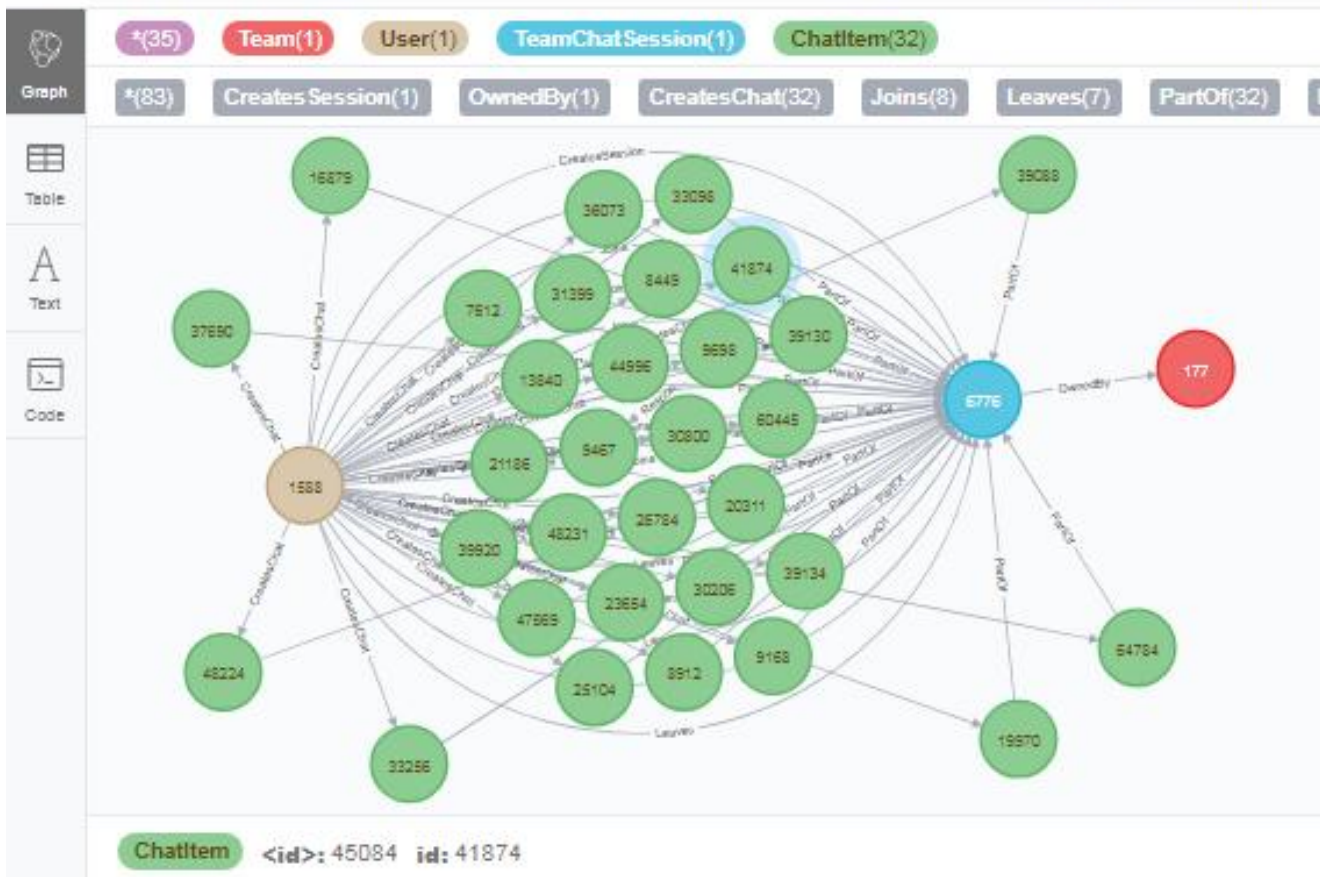
```

1 LOAD CSV FROM "file:/path/to/chat_respond_team_chat.csv" AS row
2 MERGE (m1:ChatItem {id: toInteger(row[0])})
3 MERGE (m2:ChatItem {id: toInteger(row[1])})
4 MERGE (m1)-[:ResponsesTo {timeStamp: row[2]}]->(m2)

```

- iii) Present a screenshot of some part of the graph you have generated. The graphs must include clearly visible examples of most node and edge types. Below are two acceptable examples. The first example is a rendered in the default Neo4j distribution, the second has had some nodes moved to expose the edges more clearly. Both include examples of most node and edge types.

```
$ match (m:ChatItem)-[r3]-(n:User)-[r1]-(c:TeamChatSession)-[r2]-(t:Team) return t,n,c,m,
```



## Finding the longest conversation chain and its participants

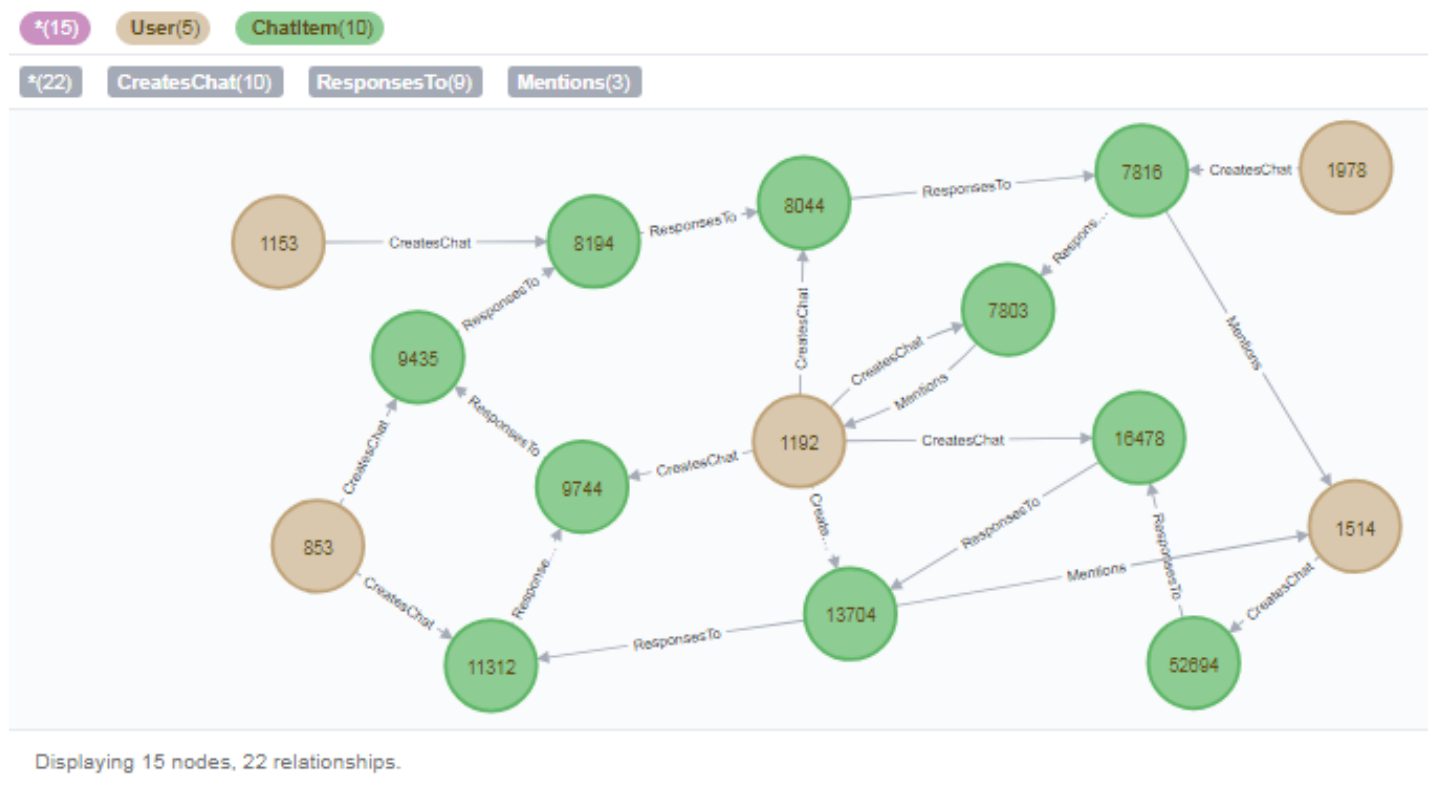
Report the results including the length of the conversation (path length) and how many unique users were part of the conversation chain. Describe your steps. Write the query that produces the correct answer.

```
1 match p=(a:ChatItem)-[e:ResponsesTo*]->(b:ChatItem)
2 return length(p) order by length(p) DESC Limit 1
```

Longest conversation chain's length = 9

```
1 match p=(a:ChatItem)-[e:ResponsesTo*]->(b:ChatItem)
2 where length(p)=9
3 match (n) where n in extract(n in nodes(p))
4 match (u:User)-[r:CreatesChat]->(n)
5 return count(distinct u)
```

Number of unique users in this chain = 5 (1153, 1192, 1514, 1978, 853)



The entire longest conversation chain path, consisting of 5 users and 10 chat items

## Analyzing the relationship between top 10 chattiest users and top 10 chattiest teams

Describe your steps from Question 2. In the process, create the following two tables. You only need to include the top 3 for each table. Identify and report whether any of the chattiest users were part of any of the chattiest teams.

### Chattiest Users

- Find all Users –CreatesChat->ChatItem relationships.
- Then count number of these ChatItems for each User
- Then rank in descending order by this count (limit to top 10)

```
1 match (u:User)-[r:CreatesChat]->(m:ChatItem)
2 return u.id,count(m) as count_chat order by count_chat desc limit 10
```

To enjoy the full Neo4j Browser experience, we advise you to use [Neo4j Browser Sync](#)

```
$ match (u:User)-[r:CreatesChat]->(m:ChatItem) return u.id,count(m) as count_chat order
```

	u.id	count_chat
Table	394	115
A	2067	111
Text	1087	109
Code	209	109
	554	107
	1627	105
	516	105
	999	105
	668	104
	461	104

Users	Number of Chats
394	115
2067	111
1087 ties with 209	109

## Chattiest Teams

- Find all ChatItem –PartOf-> TeamChatSession –OwnedBy-> Team relationships.
- Then count number of these ChatItems for each Team
- Then rank in descending order by this count (limit to top 10)

```
1 match (m:ChatItem)-[r1:PartOf]->(c:TeamChatSession)-[r2:OwnedBy]->(t:Team)
2 return t.id, count(m) as count_chat order by count_chat desc limit 10
```

To enjoy the full Neo4j Browser experience, we advise you to use [Neo4j Browser Sync](#)

```
$ match (m:ChatItem)-[r1:PartOf]->(c:TeamChatSession)-[r2:OwnedBy]->(t:Team) return t.id, cou
```

	t.id	count_chat
Table	82	1324
A	185	1036
Text	112	957
Code	18	844
	194	836
	129	814
	52	788
	136	783
	146	746
	81	736

Teams	Number of Chats
82	1324
185	1036
112	957

Finally, present your answer, i.e. whether or not any of the chattiest users are part of any of the chattiest teams.

- Out of 10 chattiest users, only **user 999** belongs to the chattiest **team 52**
- **No other relationship observed**

## How Active Are Groups of Users?

Describe your steps for performing this analysis. Be as clear, concise, and as brief as possible. Finally, report the top 3 most active users in the table below.

There are two ways to consider connection between neighbors' nodes:

- 1) Consider the edges as **directed edges**. That means if there are 2 nodes interacting with each other that connection should be **counted twice** & maximum number of connections =  $k*(k-1)$

```
1 match (n:User)-[r:InteractsWith]-(m) where n.id in [394,2067,1087,209,554,1627,516,999,668,461]
2 with n.id as nodeid, collect(distinct m.id) as neighbors, count(distinct m) as degree
3 match (u1:User),(u2:User) where u1.id in neighbors and u2.id in neighbors
4 with nodeid, neighbors, degree, sum(case when (u1)-[:InteractsWith]->(u2) then 1 else 0 end) as nb_edges
5 return nodeid, neighbors, degree, toFloat(nb_edges)/toFloat(degree*(degree-1)) as clustering_coeff
6 order by clustering_coeff desc
```

nodeid	neighbors	degree	clustering_coeff
461	[1482, 1675, 482]	3	0.8333333333333334
209	[1672, 63, 516, 2067, 1627, 1265, 2096]	7	0.8333333333333334
516	[209, 2067, 63, 1627, 1672, 1265, 2096]	7	0.8095238095238095
999	[778, 1398, 1606, 1554, 1056, 1587, 1839, 1506, 1601, 909]	10	0.7888888888888889
394	[1012, 2011, 1997, 1782]	4	0.75
1087	[929, 426, 1311, 772, 1879, 1098]	6	0.7333333333333333
554	[2018, 1959, 1687, 1096, 1010, 1412, 610]	7	0.6904761904761905
668	[698, 2034, 648, 458, 1563]	5	0.65
2067	[63, 209, 1672, 516, 1265, 1627, 697, 2096]	8	0.6428571428571429
1627	[516, 2067, 63, 209, 1672, 1265, 697, 2096]	8	0.6428571428571429

THIS METHOD SEEMS TO BE THE ONE IMPLIED BY THE ASSIGNMENT INSTRUCTION. HOWEVER THE RESULTING COEFFICIENT VALUES **DON'T MATCH THE GRADING RUBRIC**

### Most Active Users (based on Cluster Coefficients) – METHOD 1

User ID	Coefficient
461	0.8333
209	0.8333
516	0.8095

- 2) Consider the edges as **un-directed edges**. That means if there are 2 nodes interacting with each other that connection should be **counted exactly once** & maximum number of connections =  $k*(k-1)/2$

```
1 match (n:User)-[r:InteractsWith]-(m) where n.id in [394,2067,1087,209,554,1627,516,999,668,461]
2 with n.id as nodeid, collect(distinct m.id) as neighbors, count(distinct m) as degree
3 match (u1:User),(u2:User) where u1.id in neighbors and u2.id in neighbors and u1.id<u2.id
4 with nodeid, neighbors, degree, sum(case when (u1)-[:InteractsWith]-(u2) then 1 else 0 end) as nb_edges
5 return nodeid, neighbors, degree, toFloat(nb_edges)/toFloat(degree*(degree-1)/2) as clustering_coeff
6 order by clustering_coeff desc
```

The  $u1.id < u2.id$  is to prevent double counting

nodeid	neighbors	degree	clustering_coeff
394	[1012, 2011, 1997, 1782]	4	1.0
461	[1482, 1675, 482]	3	1.0
516	[209, 2067, 63, 1627, 1672, 1265, 2096]	7	0.9523809523809523
209	[1672, 63, 516, 2067, 1627, 1265, 2096]	7	0.9523809523809523
554	[2018, 1959, 1687, 1096, 1010, 1412, 610]	7	0.9047619047619048
999	[778, 1398, 1606, 1554, 1056, 1587, 1839, 1506, 1601, 909]	10	0.8666666666666667
1087	[929, 426, 1311, 772, 1879, 1098]	6	0.8
2067	[63, 209, 1672, 516, 1265, 1627, 697, 2096]	8	0.7857142857142857
1627	[516, 2067, 63, 209, 1672, 1265, 697, 2096]	8	0.7857142857142857
668	[698, 2034, 648, 458, 1563]	5	0.7

THIS METHOD **DOESN'T** SEEM TO BE THE ONE **IMPLIED BY THE ASSIGNMENT INSTRUCTION**. HOWEVER THE RESULTING COEFFICIENT VALUES **MATCH THE GRADING RUBRIC**. BUT **GOD KNOWS WHY** THEY ARE REGARDED AS **THE TOP 3**

#### Most Active Users (based on Cluster Coefficients)

User ID	Coefficient
394	1.0
461	1.0
516	0.9524