# 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 | ERD TABLE | DESCRIPTION | FIELDS | |
|---|---|---|---|---|
| ad-clicks.csv | AdClicks | A line is added to this file when a player clicks on an advertisement in the Flamingo app. | timestamp | when the click occured |
| | | | txId | unique id for the click (within ad-clicks.log) for the click |
| | | | userSessionid | id of user session for user who made click |
| | | | teamid | current team id of user who made the click |
| | | | userid | user id of user who made the click |
| | | | adId | id of the ad licked on |
| | | | adCategory | category/type of ad clicked on |
| buy-clicks.csv | InAppPurchases | A line is added when a player makes an in-app purchase on Flamingo app | timestamp | when the click occured |
| | | | txId | unique id for the click (within ad-clicks.log) |
| | | | userSessionid | id of user session for user who made click |
| | | | team | current team id of user who made the purchase |
| | | | userid | user id of user who made the click |
| | | | buyId | id of the item purchased |
| | | | price | price of the item purchased |
| users.csv | User | File contains a line for each user playing the game. | timestamp | when the click occured |
| | | | userId | user id of user who made the click |
| | | | nick | nickname chosen by the user |

| | | | twitter | twitter handle of the user |
|---|---|---|---|---|
| | | | dob | date of birth of the user |
| | | | country | 2-letter country code where the user lives |
| team.csv | Team | File contains a line for each team terminated in the game. | teamId | id of the team |
| | | | name | name of the team |
| | | | teamCreationTime | timestamp when team was created |
| | | | teamEndTime | timestamp when last member of the team |
| | | | strength | measure of team strength  roughly corresponding to the success of a time |
| | | | currentLevel | current level of a team |
| team-assignments.csv | TeamAssignment | A line is added each time a user joins a team. A user can be in at most a single team at a time. | timestamp | when the user joined the team |
| | | | team | id of the teamuser |
| | | | userId | id of the user |
| | | | assignmentId | unique id for this assignment |
| level-events.csv | LevelEvent | A line is added each time a team starts or finishes a level in the game. | timestamp | when the click occured |
| | | | eventId | unique id for the event |
| | | | teamid | id of the team |
| | | | teamLevel | level started or completed |
| | | | eventType | type of event (start or end) |

| user-session.csv | User_Sessions | Each line describes a user session, which denotes when a user starts and stops playing the game.<br>When a team goes to next game level, the session is ended for each user in the team and a new one is started. | timestamp | when the click occured |
| --- | --- | --- | --- | --- |
| | | | userSessionid | unique id for the session |
| | | | userId | current user's ID |
| | | | teamid | current user's team |
| | | | assignmentId | team assignment id for the user to the team |
| | | | sessionType | whether the event is the start or end of a session |
| | | | teamLevel | level of team during the session |
| | | | platformType | type of platform of the user during the session |

| game-clicks.csv | GameClicks | A line is added each time a user performs a click in the game. | timestamp | when the click occured |
| --- | --- | --- | --- | --- |
| | | | clickId | unique id for the click |
| | | | userId | click user's ID |
| | | | userSessionId | id of the session of user when click occurs |
| | | | idHit | if click hits flamingo (val=1) or missed (val=0) |
| | | | teamId | id of the team of the user |
| | | | teamId | id of the team of user |
| | | | teamLevel | level of team during the session |

## Aggregation

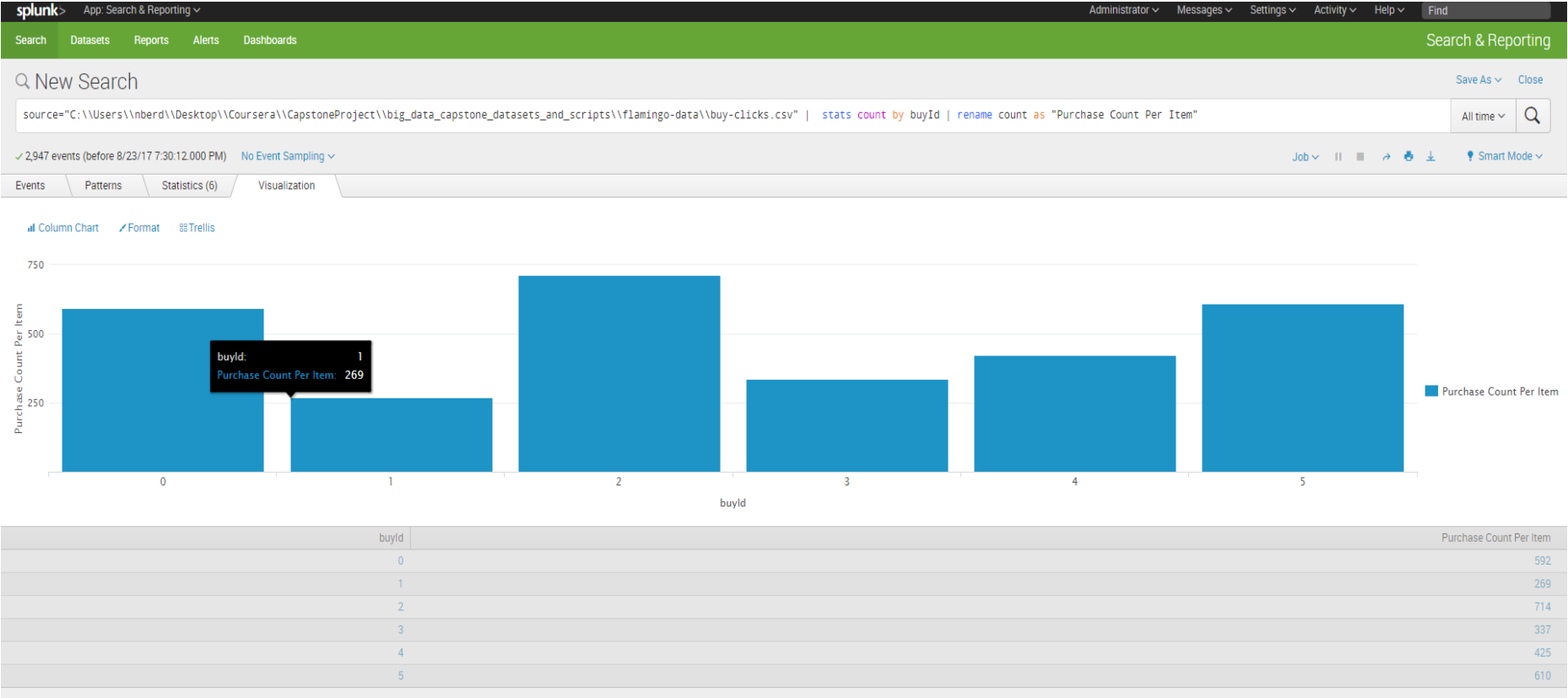| | | |
|---|---|---|
| Amount spent buying items | source="buy-clicks.csv" \| stats sum(price) | 21407.0<br><br>==This is a simple aggregation where we get all the products that gamers bought through the application and then we apply the aggregation function sum() on the price column to calculate the amount that gamers spent.== |
| Number of unique items available to be purchased | source="buy-clicks.csv" \| stats dc(buyId) | 6<br>==To find the available unique items we have to calculate the number of different item categories. Each line/record has an id of the item purchased (buyId), thus we have to fetch all items and make a DISTINCT COUNT on the buyId column.== |
| | | |

# Histograms

1. Purchase Count Per Item

This query is compiled as:
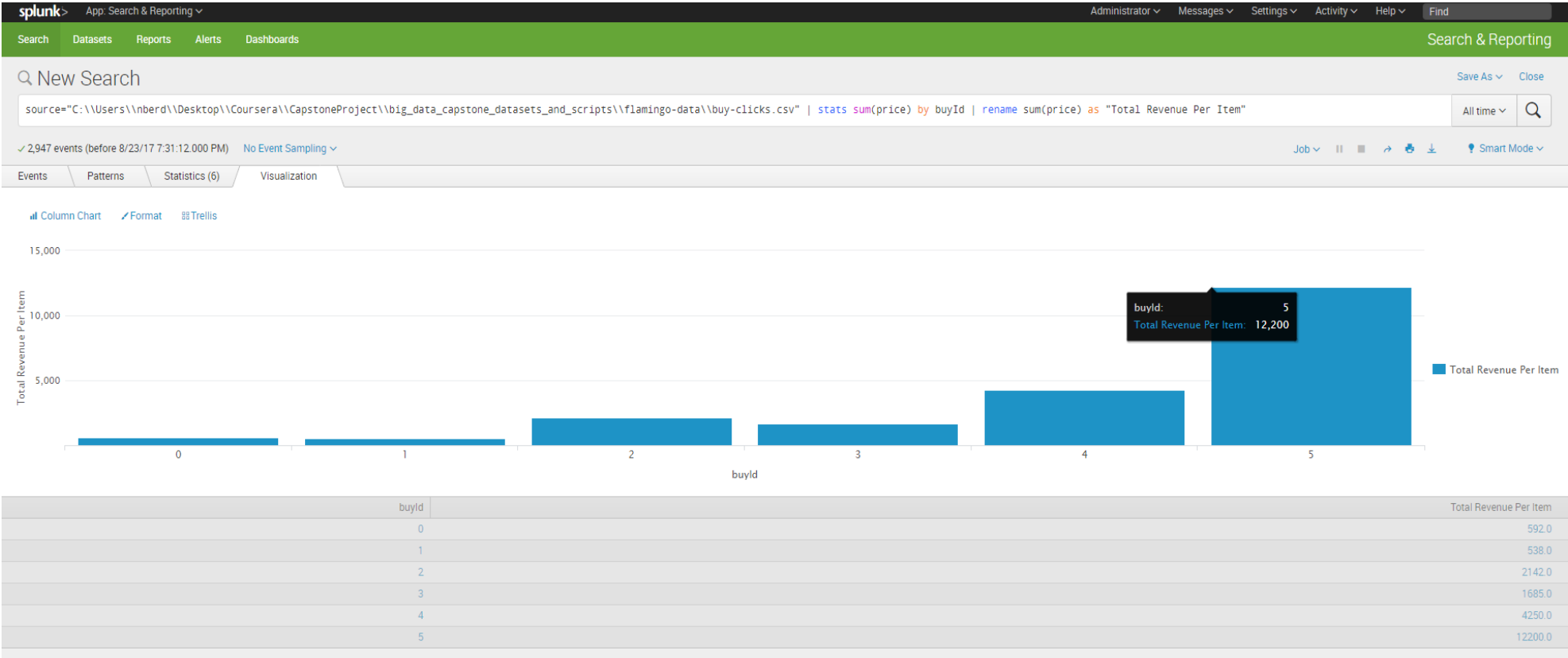source="C:\\...\\buy-clicks.csv" | stats count by buyId | rename count as "Purchase Count Per Item"

The produced histogram is the following (screenshot from Splunk console):

## 2. Total Revenue Per Item

source="buy-clicks.csv" | stats sum(price) by userId |sort sum(price) desc |rename sum(price) as "Total Money Spent" |head 10

All time ⌄   🔍

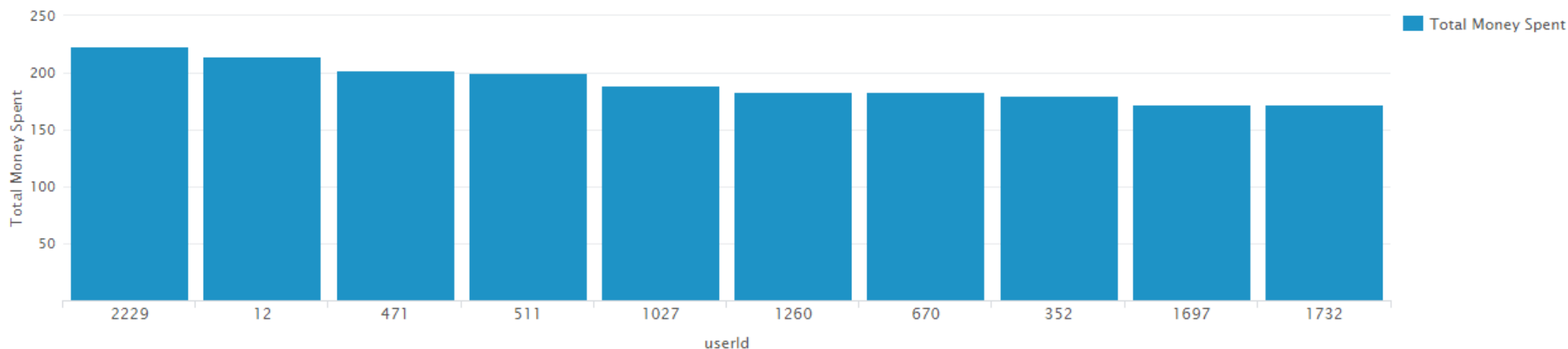✓ 2,947 events (before 6/23/17 8:40:48.000 PM)   No Event Sampling ⌄       Job ⌄  ‖  ▪  ↗  🖶  ⤓      💡 Smart Mode ⌄

Events    Patterns    Statistics (10)    Visualization

📊 Column Chart ⌄    ✎ Format ⌄



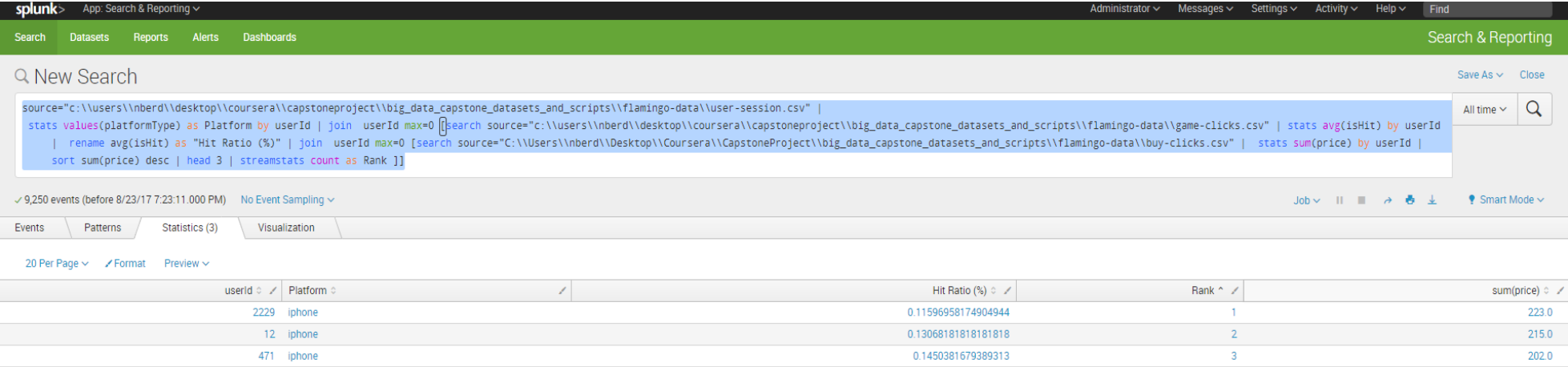| userId | Total Money Spent |
| --- | --- |
| 2229 | 223.0 |
| 12 | 215.0 |
| 471 | 202.0 |
| 511 | 200.0 |
| 1027 | 189.0 |
| 1260 | 183.0 |
| 670 | 183.0 |
| 352 | 180.0 |
| 1697 | 172.0 |
| 1732 | 172.0 |

The above query and the histogram indicate the top 10 buyers among the gamers. Each bar indicates the Total money that a gamer spent to for purchases.

# 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 | 11.6% (61/526) |
|---|------|--------|----------------|
| 2 | 12   | iphone | 13.1% (92/704) |
| 3 | 471  | iphone | 14.5% (76/524) |