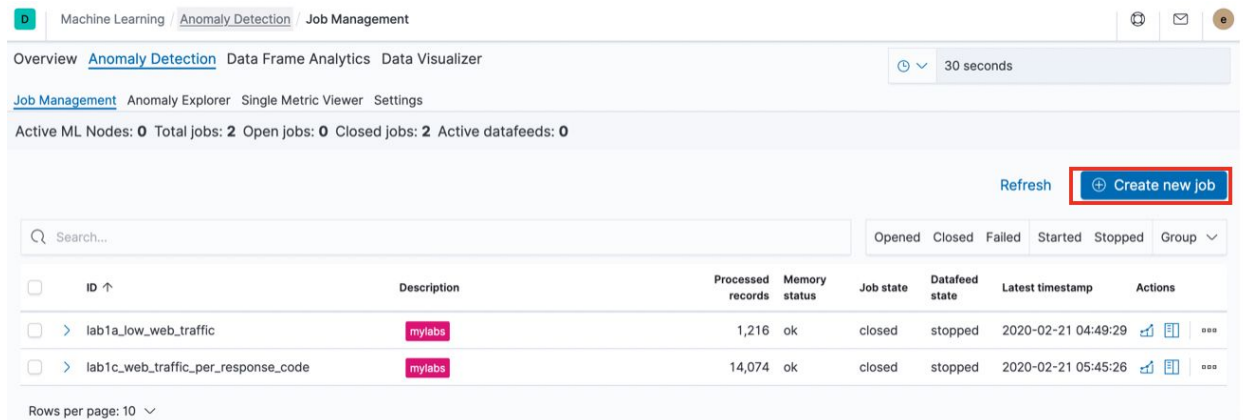# Lab 3 - Population Analysis, Data Frame & Outlier Detection

In this lab, we will be performing the following on sample eCommerce data:
   a. Set up a population analysis job to find unusual customers
   b. Use Dataframe Transform to create a customer entity-centric index
   c. Use Machine Learning's Outlier Detection to find unusual customers

# A - Set Up a Population Analysis Job

1. Click on Kibana > Machine Learning > "Anomaly Detection"
2. Click on the link to "Create new job"



3. Select the "kibana_sample_data_ecommerce" index.

4. Click on the link to create a Population Analysis job



5. Click on the "Use full kibana_sample_data_ecommerce data" button, then click on "Next"



6. Use "customer_full_name.keyword" as the Population field
Add Metric: use "sum" on the field "total_taxful_price"

7. Use 15m bucket span and click on "Next"



8. Name the job "lab3a_unusually_big_orders" and place it under "mylabs" group

9. Click "Next" after passing the job validation



10. Review the Job Settings and click on the "Create job" link to start the job

11. Wait around half a minute while the job completes. Click on the "View Results" link to review the results



12. Customer "Wagdi Shaw" is the highest anomalous customer



13. Explore Anomalous User -
    Click on red tile for user "Wagdi Shaw"

Note:

- Wagdi Shaw had a purchase of $2250, much more than the typical user purchase of $62.45
- (Optional) You could repeat the process of creating custom URL to raw data (as performed in Lab 1d) to see what Wagdi Shaw ordered

# B - Data Frame Transformation

*Use Data Frame analysis on Sample e-commerce logs to find users that spend more money than others*

Note: This is similar to the Population Analysis, except that this is not time-series oriented. In the Population demo, the oddest user ("Wagdi Shaw") was deemed anomalous because the order placed at a particular time was much bigger than the typical user. But, what if Wagdi Shaw had spread out their spending over many orders over a long period of time. We may still desire to know who our biggest spenders are, regardless of what they do on a moment by moment basis. Therefore, we need to summarize their behavior over all time (or a longer span of time)
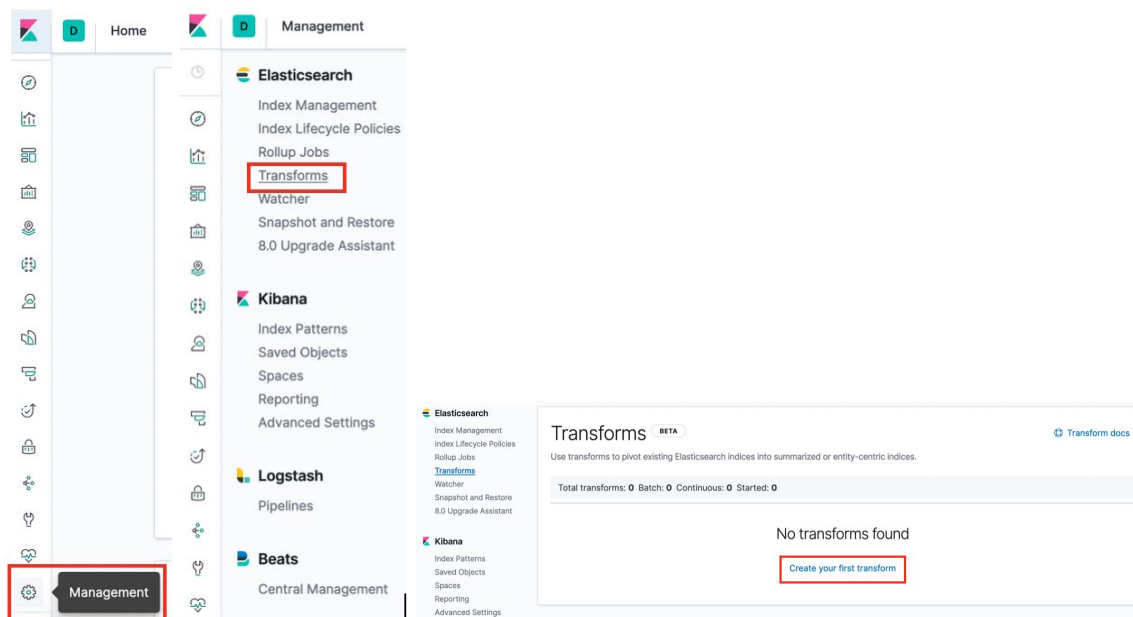
1. Transform the Data -
   We need to Transform the data from the time-series domain to an entity-centric index.
   Click on Kibana > Management > Transforms > Create your first transform

2. Choose "kibana_sample_data_ecommerce" as the source

New transform / Choose a source

| Search... | Sort ⌄ | Types **2** ⌄ |

🔍 [eCommerce] Orders

🔍 [Flights] Flight Log

🖧 kibana_sample_data_ecommerce

🖧 kibana_sample_data_flights

🖧 kibana_sample_data_logs

3. Group By - Pivot on ("Group by") the "customer_full_name.keyword"

**Group by**

customer_full_name.keyword       ✎   ✕

4. Create the following column aggregations
   a. order_id.cardinality
   b. taxful_total_price.sum
   c. products.quantity.sum

Note: This just means that we are aggregating for each customer name:
   a. The number of unique orders purchased to date
   b. Total amount they've ever spent on the e-commerce store, and
   c. The total number of products they've ever bought (by summing those fields over all time)

**Aggregations**

order_id.cardinality       ✎   ✕

taxful_total_price.sum       ✎   ✕

products.quantity.sum       ✎   ✕

Click on the "Next" button to continue

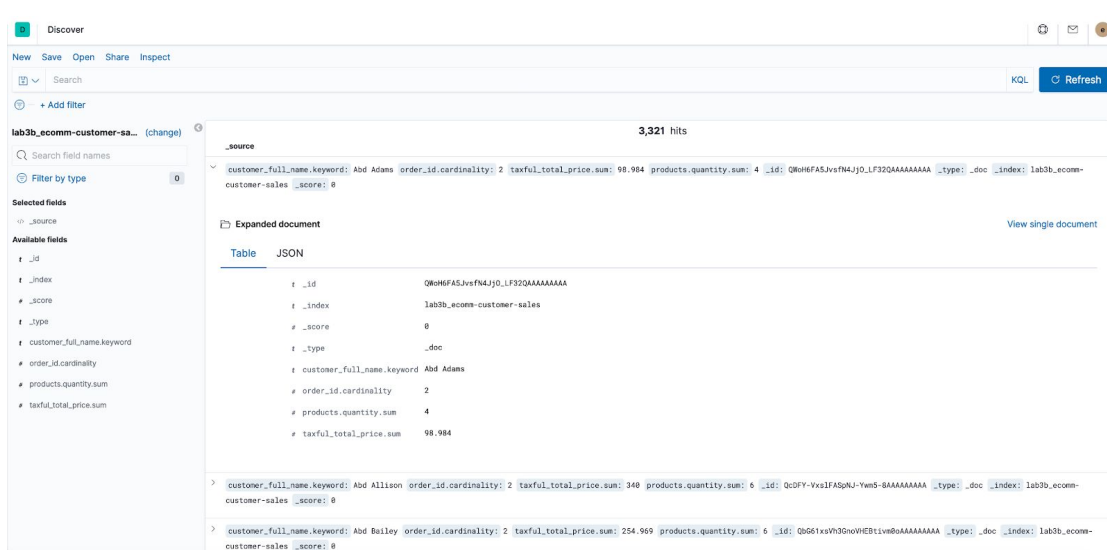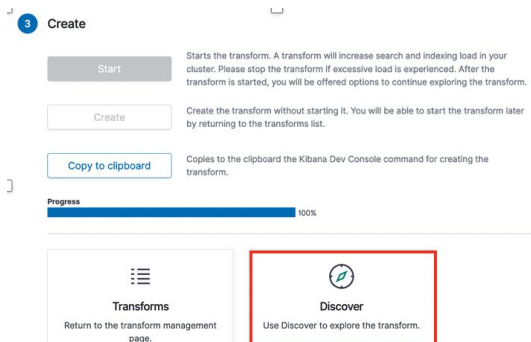5. Name both the Transform ID and Destination Index as : "lab3b_ecomm_customer_sales" and click on "Next"

6. Go ahead and click on the "Create and Start" button



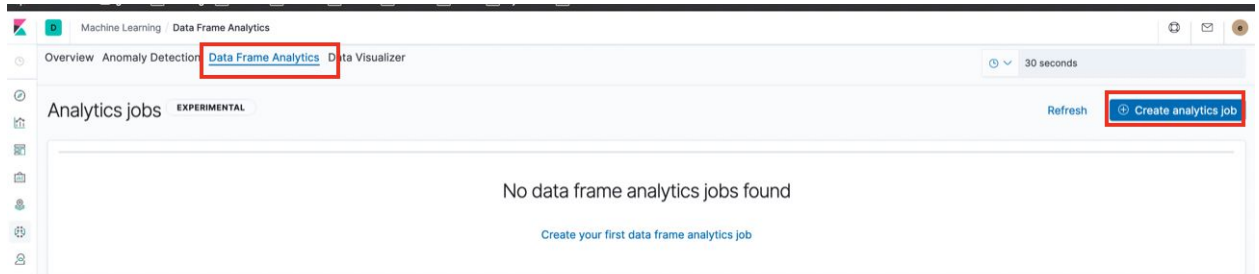7. After the job has completed, click Discover to see the new entity centric (transformed and summarized) index!



This was the first step, to get the data in this form - we can now analyze it further! Let's find the unusual shoppers!

# C - Outlier Detection

1.  Click on Kibana > Machine Learning > Data Frame Analytics link.  Then click on "Create analytics job"



2.  Choose Job type as "outlier detection".
    Enter the index from lab3b as the source index and name this job as well as destination index as **lab3c_ecomm_outliers.** Go ahead and click on the "Create" and "Start" buttons.

3. When the status of the job is "stopped" (it takes a few seconds), click on the "View" button to view the results



4. We are shown a data table with customers ranked by their outlier score.
The data frame analytics job creates an index that contains the original data and outlier scores for each document. The outlier score indicates how different each entity is from other entities.
Note that:
   ● The ml.outlier score is a value between 0 and 1. The larger the value, the more likely they are to be an outlier.
   ● In addition to an overall outlier score, each document is annotated with feature influence values for each field. These values add up to 1 and indicate which fields are the most important in deciding whether an entity is an outlier or inlier. For example, the dark shading on the products.taxful_price.sum field for Wagdi Shaw indicates that the sum of the product prices was the most influential feature in determining that Wagdi is an outlier.