

# Coding test for Data Engineer

The objective of this coding test is to get the insights from the data. We'd like you implement the test using Spark or MapReduce APIs in any programming language(Java, Scala, Python...) of your choice. Datasets are available in the data directory. And the schema is as following :

Product (product\_id, product\_name, product\_type, product\_version, product\_price)  
Customer (customer\_id, customer\_first\_name, customer\_last\_name, phone\_number )  
Sales (transaction\_id, customer\_id, product\_id, timestamp, total\_amount, total\_quantity )  
Refund (refund\_id, original\_transaction\_id, customer\_id, product\_id, timestamp, refund\_amount, refund\_quantity)

## Questions :

1. Write down the data quality issues with the datasets provided and steps performed to clean (if any).
2. Display the distribution of sales by product name and product type.
3. Calculate the total amount of all transactions that happened in year 2013 and have not been refunded as of today.
4. Display the customer name who made the second most purchases in the month of May 2013. Refunds should be excluded.
5. Find a product that has not been sold at least once (if any).
6. Calculate the total number of users who purchased the same product consecutively at least 2 times on a given day.

## Extra question (Please skip this if you don't have enough time)

Dataset for Customer\_Extended is in the data directory. And schema is as following :

Customer\_Extended (id, first\_name, last\_name, home\_phone, mobile\_phone, gender, current\_street\_address, current\_city, current\_state, current\_country, current\_zip, permanent\_street\_address, permanent\_city, permanent\_state, permanent\_country, permanent\_zip, office\_street, office\_city, office\_state, office\_country, office\_zip, personal\_email\_address, work\_email\_address, twitter\_id, facebook\_id, linkedin\_id )

## Question :

1. Display all the details of a customer who is currently living at 1154 Winters Blvd.

Please make sure to include the readMe file with clear instructions when you return the coding test.