

Final problem:

Problem statement:

The goal of the problem is to generate innovative insights from order dataset and also to increase the average order value (AOV). The dataset for this competition is a flat-file describing customers orders over time. The dataset is anonymized and contains a sample of over 6 million grocery orders from more than 30000 Milkbasket users. The only information provided about users is their sequence of orders and the products in those orders. All of the IDs in the dataset are entirely randomized, and cannot be linked back to any other ID

For each user, we provide more than 300 of their orders, with the sequence of products purchased in each order. We also provide the time of the day the product was added to the basket and the order was assumed to be created at that time itself.

File descriptions:

Full data set: [https://blume-hackathon.s3.ap-south-1.amazonaws.com/milkbasket\\_hackathon\\_data.zip](https://blume-hackathon.s3.ap-south-1.amazonaws.com/milkbasket_hackathon_data.zip)

Sample file: [https://blume-hackathon.s3.ap-south-1.amazonaws.com/milkbasket\\_hackathon\\_sample\\_data.csv](https://blume-hackathon.s3.ap-south-1.amazonaws.com/milkbasket_hackathon_sample_data.csv)

**customer\_id,manufacturer\_id,society\_id,city\_id,route\_id,store\_id,order\_id,order\_date,category\_id,subcategory\_id,product\_id,product\_quantity,selling\_price\_per\_unit,total\_cost,subscription,product\_addedtobasket\_on**  
2698080,1122016,1127168,1120112,1121456,1120112,338048928,2018-05-01,1122576,1125264,1120336,4,16.01,64.04,0,2018-04-30 16:23:46  
1134224,1134336,1120224,1120112,1120336,1120112,337533168,2018-05-01,1122576,1125152,1939280,1,8.79,8.79,1,2018-04-19 01:30:02  
3686704,1150128,1126160,1120112,1123472,1120112,338235520,2018-05-01,1123472,1130640,1681456,1,25.0,25.0,0,2018-04-30 22:26:07

description of columns:

customer\_id = id of the user/customer

manufacturer\_id = brand manufacturer id

society\_id = id of society where the order was created

city\_id = id of city from where the product was ordered

route\_id = id of route which the delivery took

store\_id = id of store from where the product was ordered

order\_id = id of order

order\_date = date of order

category\_id = category id (e.g. beverages)

subcategory\_id = subcategory id (e.g. milk)

product\_id = id of product

product\_quantity = total quantity of product added

selling\_price\_per\_unit = selling price of product (after discounts)

total\_cost = total cost (product\_quantity\*selling\_price\_per\_unit)

subscription = boolean flag to denote if the product was a subscribed product (i.e. recurring)

product\_addedtobasket\_on = time when the product was added to customers basket

Bonus: Create a mobile app that is powered on the above order data for a customer and place suggestions to the customer to buy similar or related products.

Evaluation: Submissions will be evaluated based on innovation, logic, the breath of the approach to the problem, teamwork.

Tech stack: Python, Numpy, Pandas, Scikit learn, Keras, Apache spark

Ideal team: 4 members with Data crunching skills, Stats, Programming in python/java, APIs, serverless