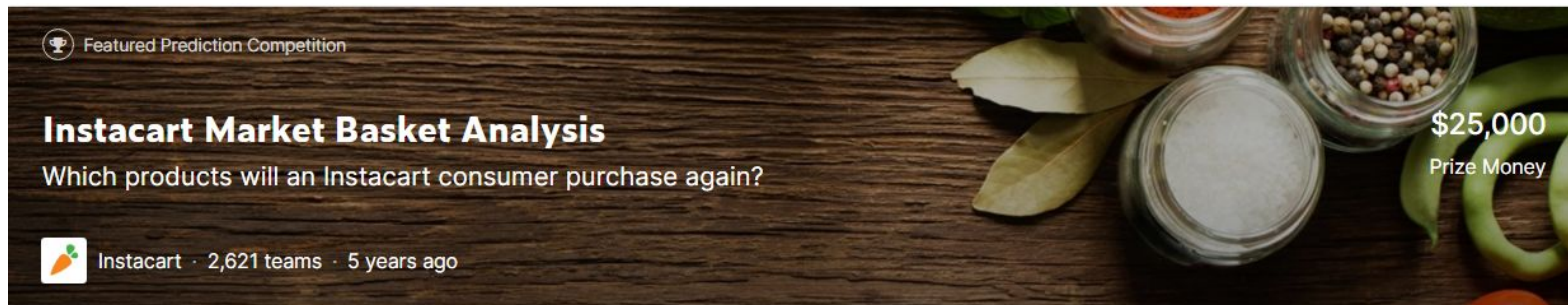# Instacart market basket analysis

Karen  2022/07/11
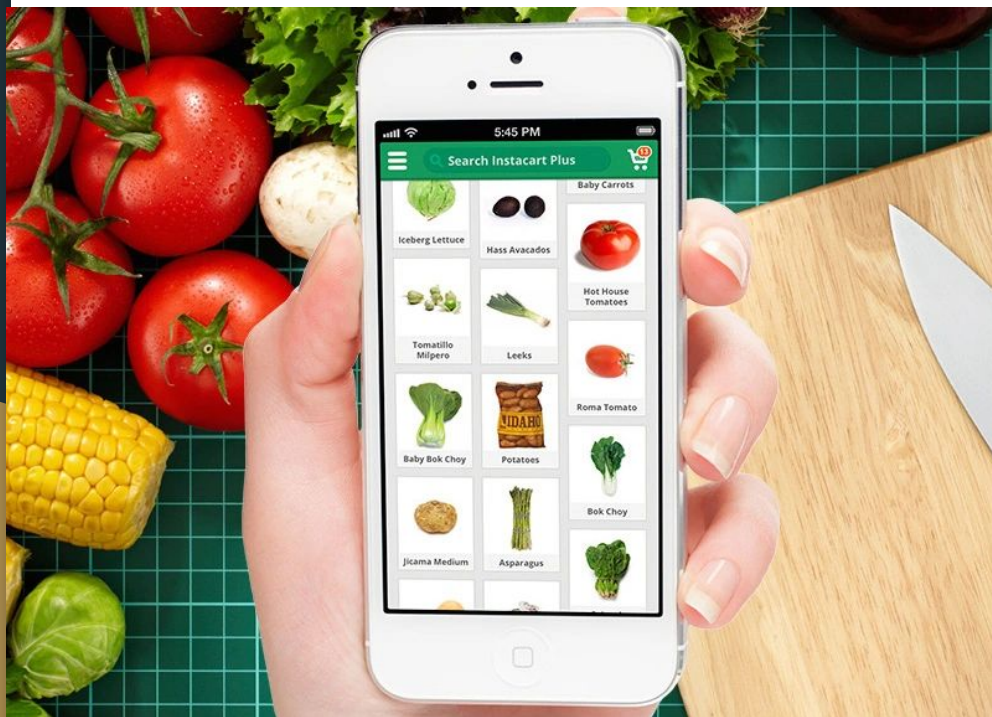
# Outline

- Background
- Exploratory Data Analysis (EDA)
- Modeling (Xgboost)
- Apyori association analysis
- Summary

# Background



Featured Prediction Competition

**Instacart Market Basket Analysis**
Which products will an Instacart consumer purchase again?

$25,000
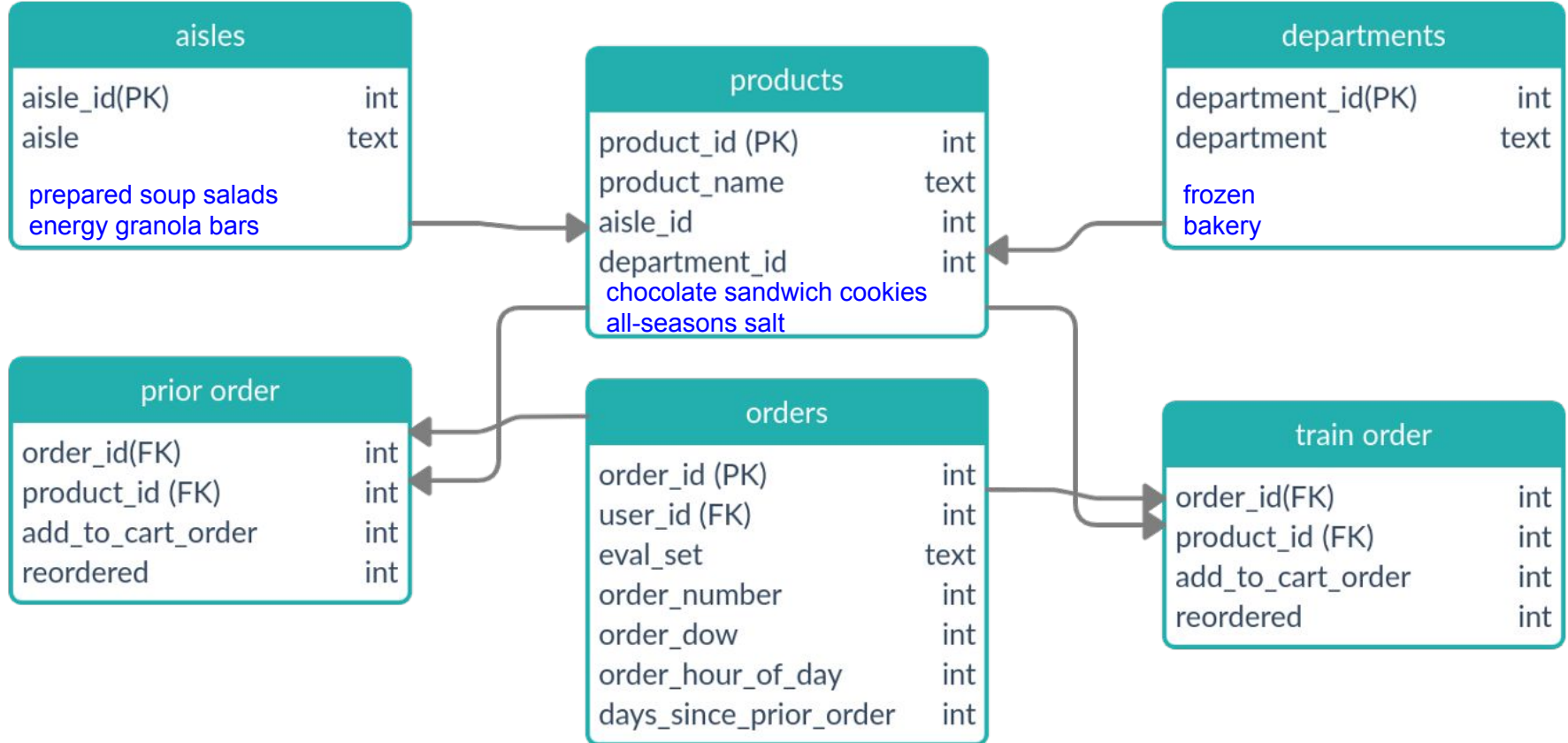Prize Money

Instacart · 2,621 teams · 5 years ago

- 公開資料集 :Instacart Online Grocery Shopping Dataset 2017
- 根據之前的購買資料計算客戶再度訪問平台時可能再次購買的商品

- 成立於 2012 年
- 提供 O2O 生鮮雜貨代買代送服務
  - O2O (Online to Offline)
  - 團購服務
- 提供使用者比價資訊
- 提供購物專家購買的最佳路徑，依照天氣、交通調整運費

# EDA - Briefly review

- 來自約 20 萬名 Instacart 用戶
- 約 340 萬的訂單數量
- 將近 5 萬件商品項目
- 這些商品的類別, 分佈 21 種
- 商品擺放的位置, 約有 134 個商品陳列走道位置
- 對於每個用戶提供 4 ~ 100 個訂單資料

|  | order_number | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| User A | p | p | p | p | p | tr |  |  |  |  |
| User B | p | p | p | p | p | p | p | p | te |  |
| User C | p | p | p | p | p | p | p | tr |  |  |
| User D | p | p | p | tr |  |  |  |  |  |  |

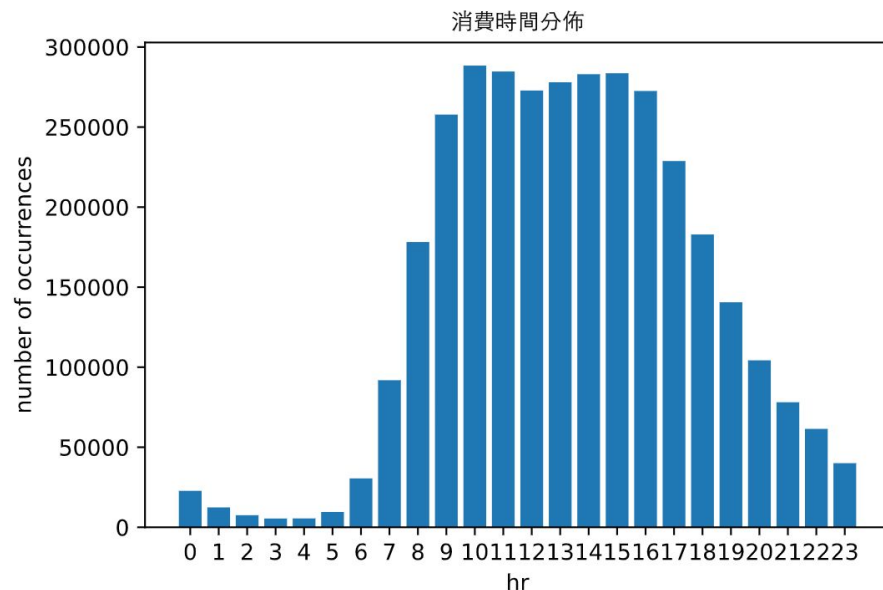Prior (p)
Train (tr)
Test (te)

# EDA- Dataset review

- 拆分三個資料集
  - Prior dataset
  - Train dataset
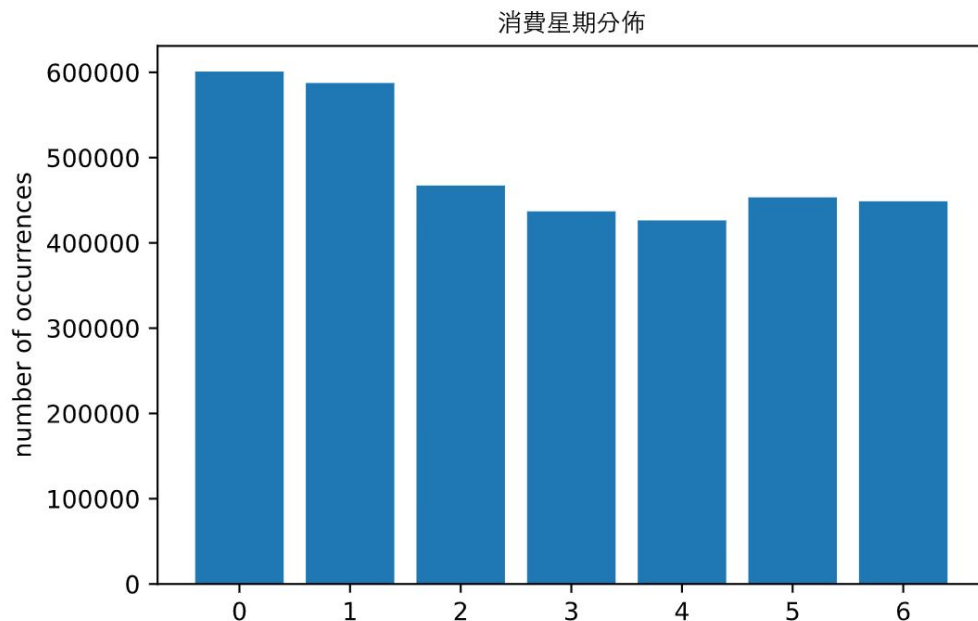  - Test dataset



count of eval_set

# EDA- Distribution of purchase time

- 消費集中在10am~4pm

# EDA- Distribution of purchase weekday

- 0為週六、1為週日以此類推
- 週六最愛買！
- 週三小週末反而消費少

消費星期分佈

# EDA- Distribution of order interval



距上次購買間隔天數分佈

# EDA- Purchase frequency

- 最低購買次數為4次
- 最高為100次



購買次數分佈

# EDA- Distribution of numbers of items

# EDA- Purchase frequency by product

- 分析過往訂單記錄, 找出 TOP 10 明星商品
- Banana, the king of the fruit?



Top 10 product

- 分析過往訂單記錄, 找出 Last 10 不受歡迎商品
- 特殊食品、酒？！

# EDA- Distribution of the department which product is only purchased one

# EDA- Distribution of product department

# Define reorder

- reordered: 1 if this product has been ordered by this user in the past, 0 otherwise

| first purchase | |
|---|---|
| product_name | reordered |
| A | 0 |
| B | 0 |

| second purchase | |
|---|---|
| product_name | reordered |
| A | 1 |
| C | 0 |
| D | 0 |

| third purchase | |
|---|---|
| product_name | reordered |
| A | 1 |
| B | 1 |
| E | 0 |

# EDA - Reorder ratio

# EDA - Distibution of department  by reorder

### reordered

| | |
|---|---|
| ■ | produce - 32.23 % |
| ■ | dairy eggs - 18.91 % |
| ■ | beverages - 9.19 % |
| ■ | snacks - 8.65 % |
| ■ | frozen - 6.35 % |
| ■ | bakery - 3.86 % |
| ■ | pantry - 3.41 % |
| ■ | deli - 3.34 % |
| ■ | canned goods - 2.56 % |
| ■ | meat seafood - 2.11 % |
| ■ | dry goods pasta - 2.10 % |
| ■ | breakfast - 2.08 % |
| ■ | household - 1.57 % |
| ■ | babies - 1.27 % |
| ■ | personal care - 0.76 % |
| ■ | international - 0.52 % |
| ■ | alcohol - 0.46 % |
| ■ | pets - 0.31 % |
| ■ | missing - 0.15 % |
| ■ | bulk - 0.10 % |
| ■ | other - 0.08 % |

### non_reordered

| | |
|---|---|
| ■ | produce - 24.93 % |
| ■ | dairy eggs - 13.40 % |
| ■ | snacks - 9.23 % |
| ■ | pantry - 9.21 % |
| ■ | frozen - 7.71 % |
| ■ | beverages - 7.01 % |
| ■ | canned goods - 4.35 % |
| ■ | dry goods pasta - 3.51 % |
| ■ | household - 3.33 % |
| ■ | bakery - 3.28 % |
| ■ | deli - 3.10 % |
| ■ | breakfast - 2.34 % |
| ■ | meat seafood - 2.30 % |
| ■ | personal care - 2.29 % |
| ■ | babies - 1.34 % |
| ■ | international - 1.28 % |
| ■ | alcohol - 0.49 % |
| ■ | missing - 0.34 % |
| ■ | pets - 0.29 % |
| ■ | other - 0.16 % |
| ■ | bulk - 0.11 % |

| | product_id | product_name | aisle_id | department_id | department | aisle |
|---|---|---|---|---|---|---|
| 42767 | 39812 | Organic Thyme | 16 | 4 | produce | fresh herbs |
| 42755 | 31717 | Organic Cilantro | 16 | 4 | produce | fresh herbs |
| 41883 | 11165 | Tuscan Kale | 83 | 4 | produce | fresh vegetables |
| 42464 | 19881 | Bartlett Pear | 24 | 4 | produce | fresh fruits |
| 41867 | 10358 | Organic White Mushrooms | 83 | 4 | produce | fresh vegetables |
| 42634 | 42411 | Young Coconut | 24 | 4 | produce | fresh fruits |
| 41821 | 6773 | Onions | 83 | 4 | produce | fresh vegetables |
| 41187 | 4539 | Santa Fe Caesar Complete Salad Kit | 123 | 4 | produce | packaged vegetables fruits |
| 42796 | 15772 | Pineapple Spears | 32 | 4 | produce | packaged produce |
| 41677 | 43787 | Bolthouse Farms Baby Cut Carrots | 123 | 4 | produce | packaged vegetables fruits |

# EDA - Distibution of aisles by reorder

## reordered



- fresh fruits - 13.66 %
- fresh vegetables - 10.64 %
- packaged vegetables fruits - 5.91 %
- yogurt - 5.19 %
- milk - 3.62 %
- water seltzer sparkling water - 3.21 %
- packaged cheese - 3.00 %
- soy lactosefree - 2.31 %
- chips pretzels - 2.23 %
- bread - 2.04 %
- refrigerated - 1.99 %
- eggs - 1.67 %
- frozen produce - 1.48 %
- energy granola bars - 1.42 %
- crackers - 1.37 %
- ice cream ice - 1.28 %
- lunch meat - 1.25 %
- soft drinks - 1.20 %
- baby food formula - 1.17 %
- frozen meals - 1.14 %
- cream - 1.13 %
- fresh dips tapenades - 1.13 %
- cereal - 1.13 %
- juice nectars - 1.03 %
- fresh herbs - 1.02 %
- packaged produce - 1.00 %

## non_reordered



- fresh vegetables - 10.42 %
- fresh fruits - 7.69 %
- packaged vegetables fruits - 4.80 %
- yogurt - 3.41 %
- packaged cheese - 3.05 %
- chips pretzels - 2.23 %
- ice cream ice - 1.91 %
- frozen produce - 1.80 %
- water seltzer sparkling water - 1.71 %
- baking ingredients - 1.70 %
- soup broth bouillon - 1.51 %
- crackers - 1.48 %
- soy lactosefree - 1.47 %
- refrigerated - 1.46 %
- milk - 1.45 %
- bread - 1.45 %
- energy granola bars - 1.37 %
- fresh herbs - 1.37 %
- spices seasonings - 1.35 %
- frozen meals - 1.30 %
- canned jarred vegetables - 1.29 %
- cereal - 1.22 %
- oils vinegars - 1.20 %
- baby food formula - 1.17 %
- lunch meat - 1.17 %
- condiments - 1.15 %
- spreads - 1.11 %
- nuts seeds dried fruit - 1.11 %
- dry pasta - 1.09 %
- fresh dips tapenades - 1.05 %
- juice nectars - 1.02 %
- hot dogs bacon sausage - 1.02 %
- canned meals beans - 1.01 %
- other creams cheeses - 1.00 %

```
fresh fruits                    2726251
fresh vegetables                2123540
packaged vegetables fruits      1178700
yogurt                          1034957
5  milk                          722128
water seltzer sparkling water    640988
packaged cheese                  598280
soy lactosefree                  460069
chips pretzels                   444036
bread                            408010
refrigerated                     397213
eggs                             333408
frozen produce                   295616
energy granola bars              283351
crackers                         272645
ice cream ice                    256194
lunch meat                       249963
soft drinks                      238981
baby food formula                233042
frozen meals                     228222
Name: aisle, dtype: int64
```

```
fresh vegetables                1445090
fresh fruits                    1066410
packaged vegetables fruits       665106
yogurt                           472626
packaged cheese                  423182
chips pretzels                   309703
ice cream ice                    264907
frozen produce                   249491
water seltzer sparkling water    237162
baking ingredients               235996
soup broth bouillon              208858
crackers                         205785
soy lactosefree                  204424
refrigerated                     201896
15  milk                         201531
bread                            200459
energy granola bars              190484
fresh herbs                      190007
spices seasonings                187516
frozen meals                     180298
Name: aisle, dtype: int64
```

| 主鍵Primary key (根據prior訂單購買的產品) | | 參數X (根據prior訂單建立) | | | train/test | future orders | 欲預測的變數Y |
| --- | --- | --- | --- | --- | --- | --- | --- |
| user_id | product_id | ... | ... | ... | eval_set | order_id | reordered |
| 1 | 196 | | | | train | 1187899 | 1 |
| 1 | 10258 | | | | train | 1187899 | 0 |
| 1 | 10486 | | | | train | 1187899 | 1 |
| 1 | 10686 | | | | train | 1187899 | 1 |
| 1 | 15435 | | | ... | train | 1187899 | 0 |
| 1 | 12376 | | | | test | 1187968 | |
| 2 | 11698 | | | | train | 1256788 | 1 |
| 2 | 12495 | | | | train | 1256788 | 1 |
| 2 | 14571 | | | | test | 1257530 | |

ref:https://medium.com/@PTLin0519/kaggle%E7%AB%B6%E8%B3%BD-instacart-market-basket-analysis-%E4%B8%80-%E7%AB%B6%E8%B3%BD%E7%B0%A1%E4%BB%8B%E8%88%87%E6%8E%A2%E7%B4%A2%E6%80%A7%E6%95%B8%E6%93%9A%E5%88%86%E6%9E%90-972183f2a19b

# Feature Engineering

- 商品面
  - 被購買的次數
  - 被重複購買比例
  - 被第一次購買的次數
  - 被第二次購買的次數
- 客戶面
  - 距離上一次購買天數的總和和平均
  - 平均單次購買的 產品數量
  - 平均將該商品放入購物車的順序
  - 平均購買該商品的次數
  - 連續沒有購買該商品的次數

# Modeling

$$F1\ score = 2 * \frac{Precision * Recall}{Precision + Recall}$$

**XGBoost**

- Baseline model score　　: 0.37625
- Grid Search score　: 0.37968 (38.3%)

- 高準確度
- 適合作為 baseline model
- 準確度較不受無效特徵干擾

| # | △ | Team | Members | | Score | Entries | Last |
|---|---|---|---|---|---|---|---|
| 1 | — | 胡萝卜 | | | 0.40914 | 62 | 5Y |
| 2 | — | ===== KEEP OUT 🔒===== | | | 0.40820 | 138 | 5Y |
| 3 | — | sjv | | | 0.40810 | 76 | 5Y |

# Features importance analysis



- 重複購買率
  (user_reorder_ratio)
- 被購買的次數 (prod_order)
- 距離上一次購買天數的平均
  (mean_days_prior)
- 平均單次購買的產品數量
  (user_avg_basket)

# Apyori association analysis

- 關聯分析: 找尋資料彼此之間的關聯, 透過兩種主要的方式來進行分析
  - 頻繁項集: 經常一起出現的物品集合
  - 關聯規則: 表達數據之間的可能存在很強關聯性

- 支持度(Support) : 表示為 item-set 在整個 AllSamples 中出現的頻率

  Support(X) = number(X) / number(AllSamples)

- 信心度(Confidence) : 表示當事件X發生的情況下, 同時會發生Y的可能性

  Confidence(X→Y) = P(Y|X) , = P(X∩Y) / P(X)

| 交易編號 | 商品 |
|:---:|:---:|
| 0 | 豆漿、萵苣 |
| 1 | 萵苣、尿布、葡萄酒、甜菜 |
| 2 | 豆漿、尿布、葡萄酒、橙汁 |
| 3 | 萵苣、豆漿、尿布、葡萄酒 |
| 4 | 萵苣、豆漿、尿布、橙汁 |

# Summary

- 藉由過去的訂單數據, 利用 model 能預測下次顧客是否再度購買商品
- 觀察 organic 在資料集中, 可能成為重要的特徵
- 通過關聯分析發現訂單中的商品有交互關係, 例如購買A產品經常購買B

# Future works

- data augmentation: 加入用戶最近的 3 - 5 個訂單提高訓練數據量
- product feature engineering: organic  feature, alternative item
- none prediction model: 有可能用戶下一次訂單中不回購任何商品
- other models: RNN, CNN, LGBMClassifier
- design new training flow: 由於產品間有交互作用, 可以設計新的訓練流程, 不只是將每一個產品當成獨立的分類問題

Thank you