# INFO7374 Algorithmic Digital Marketing

| Summary | In this codelab, we analyzed Elo dunnhumby_The-Complete-Journey dataset to summarize the insights. |
| --- | --- |
| Author | Yifu Liu and Pengfei He |

# About Datasets

## product.csv

All product are included in this csv file
Each product is unique

contains :

- PRODUCT_ID
- MANUFACTURER
- DEPARTMENT

## transaction.csv

Only products with purchasing quantity are included in the dataset.
Product id are consistent between product.csv

contains:

- household_key
- PRODUCT_ID
- QUANTITY
- COUPON_DISC

## coupon.csv

contains:

- COUPON_UPC
- PRODUCT_ID

● CAMPAIGN

# Coupon.csv

contains:

- household_key
- Coupon
- Days
- CAMPAIGN

# campaign_table.csv

Identifiers that can be used to link to other sources of movie data are contained in the file links.csv.

contains :

- CAMPAIGN
- household_key
- End date
- Description

# Campaign_desc.csv

Contains:

- Start date
- End date
- CAMPAIGN
- Description

# causal_data.csv

contains:

- Store id
- mailer
- display

# Data Wrangling

**Data wrangling**, sometimes referred to as **data** munging, is the process of transforming and mapping **data** from one "raw" **data** form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.

## Trifacta

We have used Trifacta for joining the tables
And we could aggregate data and filter

Add Datasets

campaign_table – 2.csv

campaign_table – 3

campaign_table – 4

campaign_table – 5

campaign_
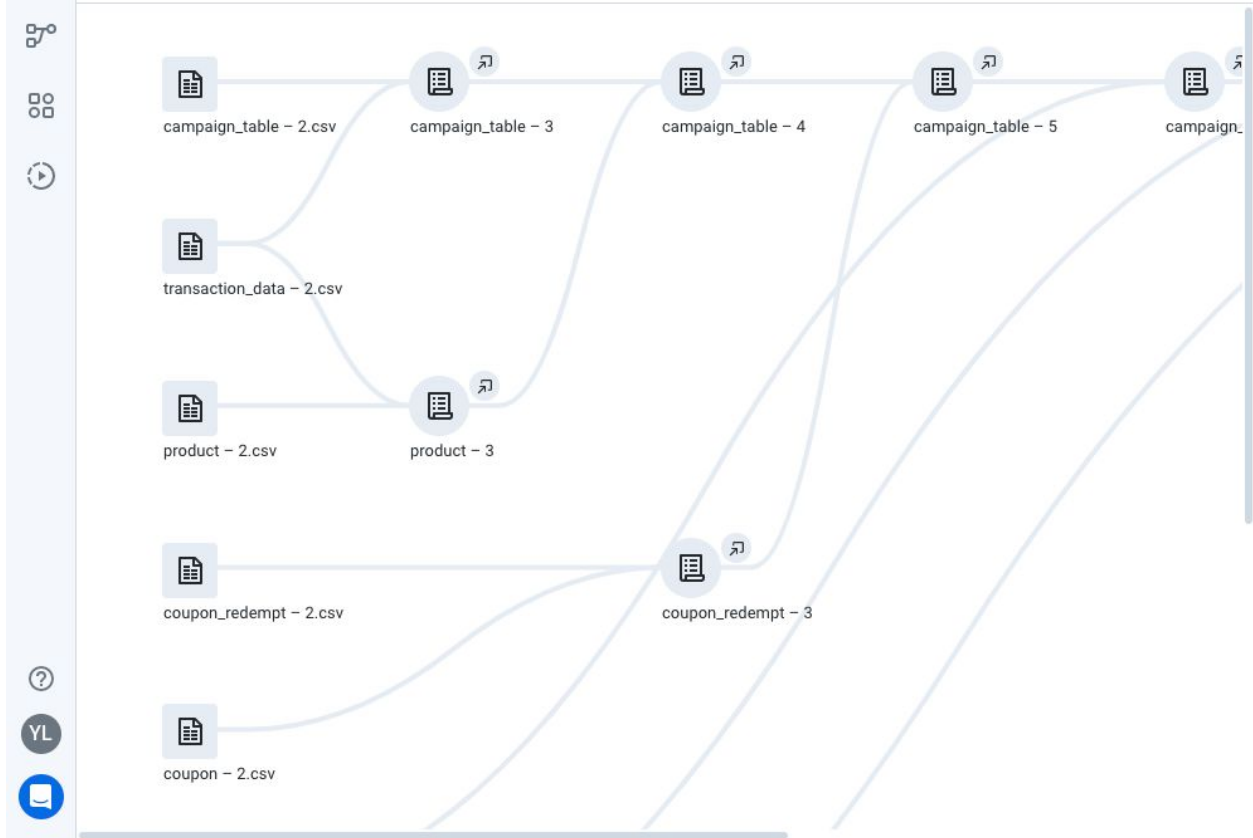
transaction_data – 2.csv

product – 2.csv

product – 3

coupon_redempt – 2.csv

coupon_redempt – 3

coupon – 2.csv

ASSIGNMENT >
campaign_table – 4 ⌄
Initial Sample

Run Job

| # household_key | # household_key1 | # household_key2 | ᴬᴮᶜ DESCRIPTION | # CAMPAIGN | # BASKET_ |
|---|---|---|---|---|---|
| 1 - 2.49k | 1 - 2.49k | 1 - 2.49k | 3 Categories | 1 - 30 | 26.98B - 27.67B |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 17 | 17 | 17 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 27 | 27 | 27 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 212 | 212 | 212 | TypeA | 26 | 2 |
| 208 | 208 | 208 | TypeA | 26 | 2 |

34 Columns    28,791 Rows    3 Data Types

Run Job

**Group by** ✕

Group by

| # PRODUCT_ID | ✕ | ✕ ⌄ |

Values

Value or formula

Type                                                    required

✓    Group by as new table

Group by as new column(s)

Cancel    **Add**

**Row Labels**

| # PRODUCT_ID | # PRODUCT_ID1 | # CAMPAIGN | # CAMPAIGN1 |
|---|---|---|---|
| 869.39k - 879.88k | 869.39k - 879.88k | 2 - 30 | 2 - 30 |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 871158 | 871158 | 26 | |
| 878234 | 878234 | 26 | |

**Preview**

| # PRODUCT_ID |
|---|
| 869.39k - 879.88k |
| 871158 |
| 878234 |
| 871337 |
| 875118 |
| 873847 |
| 873203 |
| 877017 |
| 869688 |
| 873654 |
| 878683 |

1 Column    38 Rows    1 Data Type

# Data Integration, Profiling and Cleaning

## Report - Job 171661
genome_stats Flow - genome_stats

All Data     7 columns   30,424 rows   3 data types

● 100% valid values    ● 0% mismatching values    ● 0% missing values

### tag

| | |
|---|---|
| Type | String |
| Valid | 30,424 |
| Mismatched | 0 |
| Empty | 0 |

**Top 20 values**

| | |
|---|---|
| original | 584 |
| comedy | 513 |
| action | 468 |
| based on a book | 294 |
| imdb top 250 | 232 |
| funny | 228 |
| chase | 221 |
| adapted from:book | 214 |
| mentor | 209 |
| franchise | 198 |
| animation | 197 |
| relationships | 182 |
| romantic | 180 |
| horror | 179 |
| sci-fi | 177 |
| drama | 171 |
| good action | 169 |
| romantic comedy | 168 |
| love story | 168 |
| hilarious | 166 |

### ratings_count

| | |
|---|---|
| Type | Integer |
| Valid | 30,424 |
| Mismatched | 0 |
| Empty | 0 |

25,680

2,936   952   432   184   104   80   16   40

0k 10k 20k 30k 40k 50k 60k 70k 80k 90k

| | |
|---|---|
| Minimum | 1,001 |
| Lower quartile | 1,588 |
| Median | 2,863 |
| Upper quartile | 6,592 |
| Maximum | 81,491 |

### ratings_mean

| | |
|---|---|
| Type | String |
| Valid | 30,424 |
| Mismatched | 0 |
| Empty | 0 |

**Top 20 values**

| | |
|---|---|
| 3.1381082310000004 | 32 |
| 3.5809299589999997 | 32 |
| 3.9244774110000002 | 16 |
| 3.833333333 | 16 |
| 3.281270465 | 16 |
| 3.698573305 | 16 |
| 3.7076640660000004 | 16 |
| 3.397994769 | 16 |
| 3.4579283110000003 | 16 |
| 3.61858006 | 16 |
| 2.8323678489999997 | 16 |
| 4.103381267 | 8 |
| 3.2007952289999997 | 8 |
| 4.040050528 | 8 |
| 3.86651446 | 8 |
| 2.9085053989999996 | 8 |
| 3.575661765 | 8 |
| 3.325702713 | 8 |
| 3.3910329989999997 | 8 |
| 3.227533119 | 8 |

### relevance

| | |
|---|---|
| Type | Decimal |
| Valid | 30,424 |
| Mismatched | 0 |
| Empty | 0 |

21,351

6,464

2,049

56   504

0.4   0.5   0.6   0.8   0.9   1

| | |
|---|---|
| Minimum | 0.49 |
| Lower quartile | 0.88 |
| Median | 0.94 |
| Upper quartile | 0.97 |
| Maximum | 1.00 |

Assignment › product – 3
▶ Job 167149
Finished 05/23/2020

Download results   ⋯

Overview    Output Destinations    Profile    Dependencies

All data

Download as PDF    Download as JSON

19 columns    1.049M rows    3 data types

● 99.9% valid values    ● 0.1% mismatching values    ● 0% missing values

Results profile by column

| # PRODUCT_ID | | # PRODUCT_ID1 | | # MANUFACTURER | | ABC DEPARTMENT | | ABC BRAND | |
|---|---|---|---|---|---|---|---|---|---|
| Valid | 1,048,575 | Valid | 1,048,575 | Valid | 1,048,575 | Valid | 1,048,575 | Valid | 1,04 |
| Mismatched | 0 | Mismatched | 0 | Mismatched | 0 | Mismatched | 0 | Mismatched | |
| Empty | 0 | Empty | 0 | Empty | 0 | Empty | 0 | Empty | |

DEPARTMENT Top 20 values:

| | |
|---|---|
| GROCERY | 667,426 |
| DRUG GM | 113,314 |
| PRODUCE | 102,475 |
| MEAT-PCKGD | 45,749 |
| MEAT | 36,029 |
| DELI | 25,366 |
| PASTRY | 15,361 |
| NUTRITION | 12,004 |
| KIOSK-GAS | 8,366 |
| SEAFOOD-PCKGD | 4,265 |
| SALAD BAR | 3,691 |
| COSMETICS | 3,006 |
| MISC SALES TRAN | 2,452 |
| FLORAL | 1,792 |
| SEAFOOD | 1,578 |
| MISC. TRANS. | 884 |
| SPIRITS | 842 |
| COUP/STR & MFG | 396 |
| TRAVEL & LEISUR | 363 |
| GARDEN CENTER | 232 |

BRAND Top 2 values:

| | |
|---|---|
| National | 7 |
| Private | 3 |

| | PRODUCT_ID | PRODUCT_ID1 | MANUFACTURER |
|---|---|---|---|
| Minimum | 25,671 | 25,671 | 1 |
| Lower quartile | 914,030 | 914,030 | 69 |
| Median | 1,019,858 | 1,019,858 | 537 |
| Upper quartile | 1,121,095 | 1,121,095 | 1,258 |
| Maximum | 13,512,965 | 13,512,965 | 6,434 |

---

| genres | | title | | movieId | |
|---|---|---|---|---|---|
| Type | String | Type | String | Type | Integer |
| Valid | 30,424 | Valid | 30,424 | Valid | 30,424 |
| Mismatched | 0 | Mismatched | 0 | Mismatched | 0 |
| Empty | 0 | Empty | 0 | Empty | 0 |

genres Top 20 values:

| | |
|---|---|
| Drama | 2,544 |
| Comedy | 2,288 |
| Comedy|Drama | 1,224 |
| Comedy|Romance | 1,168 |
| Drama|Romance | 1,096 |
| Comedy|Drama|Romance | 968 |
| Drama|Thriller | 608 |
| Crime|Drama | 448 |
| Action|Adventure|Sci-Fi | 424 |
| Crime|Drama|Thriller | 424 |
| Horror | 384 |
| Horror|Thriller | 344 |
| Documentary | 328 |
| Comedy|Crime | 320 |
| Drama|War | 304 |
| Action|Crime|Thriller | 296 |
| Action|Adventure|Thriller | 288 |
| Action|Sci-Fi|Thriller | 280 |
| Action|Comedy | 264 |
| Children|Comedy | 256 |

title Top 20 values:

| | |
|---|---|
| War of the Worlds (2005) | 32 |
| Confessions of a Dangerous Mind (2002) | 32 |
| Hamlet (2000) | 16 |
| SLC Punk! (1998) | 16 |
| Journey to the Center of the Earth (2008) | 16 |
| Enron: The Smartest Guys in the Room (2005) | 16 |
| Aladdin (1992) | 16 |
| Dracula (1931) | 16 |
| Hostage (2005) | 16 |
| 9 (2009) | 16 |
| Prisoners (2013) | 8 |
| Attack the Block (2011) | 8 |
| Carnal Knowledge (1971) | 8 |
| Grumpy Old Men (1993) | 8 |
| Reversal of Fortune (1990) | 8 |
| Modern Times (1936) | 8 |
| Papillon (1973) | 8 |
| Dead Man (1995) | 8 |
| Quiz Show (1994) | 8 |
| Band of Brothers (2001) | 8 |

movieId:

| | |
|---|---|
| Minimum | 1 |
| Lower quartile | 1,770 |
| Median | 3,824 |
| Upper quartile | 43,392 |
| Maximum | 201,773 |

## Advantages:

- It is very handy to see all the data in all table,  you can check the column easily before

joining
- Trifacta is a good tool in the big data tool category of a tech stack.

- Trifacta is an open source tool with GitHub stars and GitHub forks. Here's a link to

  Trifacta's open source repository on GitHub

## Disadvantages:
- Dataset cant exceed 100MB
- Dataset load is lower than 100MB
- It run slowing during the job

# Pandas

Data preprocessing

Missing data handling

```
[5]:  transaction_data = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Compl
      coupon = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complete Journe
      coupon_redempt = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complete
      campaign_table = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complet
      product = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complete Journe
      hh_demographic = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complet
      campaign_desc = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complete
      causal_data = pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumby - The Complete J
```

```
[6]:  campaign_table.info()
      campaign_desc.info()
      causal_data.info()
      coupon.info()
      coupon_redempt.info()
      hh_demographic.info()
      product.info()
      transaction_data.info()
      campaign_desc = campaign_desc.dropna()
      campaign_table = campaign_table.dropna()
      causal_data = causal_data.dropna()
      coupon = coupon.dropna()
      coupon_redempt = coupon_redempt.dropna()
      hh_demographic = hh_demographic.dropna()
      product = product.dropna()
      transaction_data = transaction_data.dropna()
```

Sampling from casual data which is 600MB（over 100MB）

[25]: `df2.loc[df2['PRODUCT_ID']==826830]`

[25]:

|  | Unnamed: 0 | PRODUCT_ID | STORE_ID | WEEK_NO | display | mailer |
|---|---|---|---|---|---|---|
| 3570016 | 5370012 | 826830 | 286 | 17 | 0 | A |
| 3570017 | 5370013 | 826830 | 286 | 18 | 0 | H |
| 3570018 | 5370014 | 826830 | 286 | 38 | 0 | A |
| 3570019 | 5370015 | 826830 | 286 | 39 | 0 | A |
| 3570020 | 5370016 | 826830 | 286 | 49 | 0 | A |
| ... | ... | ... | ... | ... | ... | ... |
| 3571852 | 5371848 | 826830 | 34280 | 90 | 6 | A |
| 3571853 | 5371849 | 826830 | 34280 | 91 | 6 | 0 |
| 3571854 | 5371850 | 826830 | 34280 | 92 | A | 0 |
| 3571855 | 5371851 | 826830 | 34280 | 93 | 6 | 0 |
| 3571856 | 5371852 | 826830 | 34280 | 96 | 0 | A |

1841 rows × 6 columns

[24]: `df2.loc[df2['PRODUCT_ID']==1018588]`

[24]:

|  | Unnamed: 0 | PRODUCT_ID | STORE_ID | WEEK_NO | display | mailer |
|---|---|---|---|---|---|---|
| 10614868 | 12414864 | 1018588 | 286 | 11 | 0 | A |
| 10614869 | 12414865 | 1018588 | 286 | 37 | 0 | A |
| 10614870 | 12414866 | 1018588 | 288 | 11 | 0 | A |
| 10614871 | 12414867 | 1018588 | 288 | 37 | 0 | A |

```
[60]: df.drop(df.index[1:7000000],inplace=True)
      df.drop(df.index[1000000:35786520],inplace=True)
```

```
[61]: len(df)
```

```
[61]: 1000000
```

```
[62]: df.to_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/dunnhumb

      df2=pd.read_csv('/Users/check4068/Desktop/算法营销作业/dunnhumby_The-Complete-Journey/d
      print(df2)
```

```
         Unnamed: 0  PRODUCT_ID  STORE_ID  WEEK_NO  display  mailer
0                 0       26190       286       70        0       A
1           7000000      869077       366       41        0       A
2           7000001      869077       366       52        0       A
3           7000002      869077       366       56        0       D
4           7000003      869077       366       57        0       A
...             ...         ...       ...      ...      ...     ...
999995      7999994      897158       293       38        0       A
999996      7999995      897158       293       88        0       F
999997      7999996      897158       295       35        0       H
999998      7999997      897158       295       38        0       A
999999      7999998      897158       295       88        0       F

[1000000 rows x 6 columns]
```

We saw the product id between 800000 and 1000000 is very frequent, which is good for joining, matching more % while joining with other tables
So we sample the data between 800000 and 1000000

```
[155]: a1 = transaction_data.groupby(by=['DAY']).agg({'QUANTITY':'sum'})
       a1.head()
       p1 = sns.distplot(a1, kde=False, hist=True)
       p1.set(title='Distribution of days',
               xlabel='# of Day',
               ylabel='# of product sold quant');
```


Distribution of days

We use groupby aggregation function to see the trend of days

```
[150]: a3 = coupon.groupby(by=['CAMPAIGN']).agg({'PRODUCT_ID':'count'})
       a3.head()
       p3 = sns.distplot(a3, kde=False, hist=True)
       p3.set(title='Distribution of Product campaign frequency',
              xlabel='# of PRODUCT',
              ylabel='# of SUM of Campaign');
```



Distribution of Product campaign frequency

We use groupby aggregation function to see the product campaign frequency

```
[131]: merge1 = coupon.merge(coupon_redempt, on=['COUPON_UPC'], how='inner')
       merge1.head()
```

[131]:

| | COUPON_UPC | PRODUCT_ID | CAMPAIGN_x | household_key | DAY | CAMPAIGN_y |
|---|---|---|---|---|---|---|
| 0 | 10000089064 | 27754 | 9 | 321 | 446 | 9 |
| 1 | 10000089064 | 27754 | 9 | 1773 | 439 | 9 |
| 2 | 10000089064 | 243186 | 9 | 321 | 446 | 9 |
| 3 | 10000089064 | 243186 | 9 | 1773 | 439 | 9 |
| 4 | 10000089064 | 872316 | 9 | 321 | 446 | 9 |

We merge two table in pandas, but we need to find out what we should join on

## Advantages:
- It is good for data preprecessing
- It is handy to plot a small part of data

## Disadvantages:

- We need to find out what we should join on during table merging
- Not as handy as trifacta when mutli-table are involved, because it is hard to find out the right coloums in so many coloums

# Snowflake and EA

The data is imported into Snowflake, custom warehouse, schema and table is created.



Create custom SQL query to get information

Connected to EA

## Advantages:

- User-friendly UI, especially for large dataset
- Strongly computing capabilities in handling huge datas

## Disadvantages:

- Dataset has a limitation, so we have to sample the data to use it.
- The free-trial has time-limit.

# Questions to Answer

## 1 Which columns are dimensions, which columns are measures?

Dimensions are columns like: Department, Brand, Commodity_Desc, Sub-Commodity_Desc

Measures are columns like: Curr_size_of_product, Sales_value, Trans_Time, Retail_Disc, Retail_disc, Trans_time,

Columns we choose to drop(missing values or null): Coupon_disc, Coupon_match_disc,

## 2 How would you generate new dimensions? What will you do to summarize measures?

We mainly use map with lambda expressions to generate new dimensions and use built-in methods from pandas to compute the measures.

For example, we will use like:

table[new dimension] = table[old dimension].map(lambda)

Table_mean = table[measure].mean()

## 3 Dashboards:

We face some problems in saving the object to EA after connecting. So we used a sample dataset to demonstrate how to implement a dashboard.

# The Motivator 2 ▾

Data updated: Today at 6:09 PM

| View As | Activity Owner | Account Name | Time Period |
|---|---|---|---|
| All | All | All | All |

**Bruce Kennedy**

**Johnny Green**

**Eric Sanchez**

**Catherine Bro...**

| Total Activities | Completed Activities | | Overdue Activities | |
|---|---|---|---|---|
| 2,061 | No Results Found ◄ ► | 0 | 2,061 ◄ ► | 0 |

Metrics are compared to team average

| Calls | Emails | Events | Tasks |
|---|---|---|---|
| 1,907 | 82 | No Results Found | 72 |

| Calls By Week | Emails By Week | Events By Week | Tasks By Week |
|---|---|---|---|

No results found

No results found

| 1.9к | 82 | | 72 |
|---|---|---|---|

| Inbound Calls | 603 | High Priority | 17 | Complete | No Results Found | Complete | No Results Found |
|---|---|---|---|---|---|---|---|
| Outbound Calls | 615 | Normal Priority | 57 | Open | No Results Found | Open | 72 |