

- The two dataframes, user table and usage summary table, need to be merged. In the user\_table the column name 'object\_id' was renamed user\_id to make merging the two dataframes easier. There they were merged on user\_id.

Code: **df = data2.merge(data1,on='user\_id',how='outer')**

- Time\_stamp column was renamed login\_time for easier analysis
- Visited column had missing values after merging for customers who did not login in, here nans were filled with 0

Code: **df.visited = df.visited.fillna(0)**

- The response variable needs to be created which is “adopted user” who is a user who has logged into the product on three separate days in at least one seven day period, to do this a few steps were taken

- Create a grouped dataframe, grouped by user id so

**group = df.groupby(['user\_id'])['visited'].count().reset\_index()**

- Then get users with visited value more than 1 because only 1 means they just created the account but did not login

**active\_users = group[group.visited > 1]**

- A new column with 1 for active and 0 for not active is created by using a mask of the list of users

**df\_active = df[df.user\_id.isin(active\_users.user\_id)]**

- Calculate 7 days from login time. If the next login time was in less than 7 days then true for active. Then the trues were summed per user and only users with 3 or more were kept

**df\_active['seven\_days'] = df\_active['login\_time']+timedelta(7)**

**df\_active['active'] = df\_active.login\_time <**

**df\_active.seven\_days.shift(-1)**

**active\_users =**

**df\_active.groupby('user\_id')['active'].sum().reset\_index()**

**adopted\_users = active\_users[active\_users.active>2]**

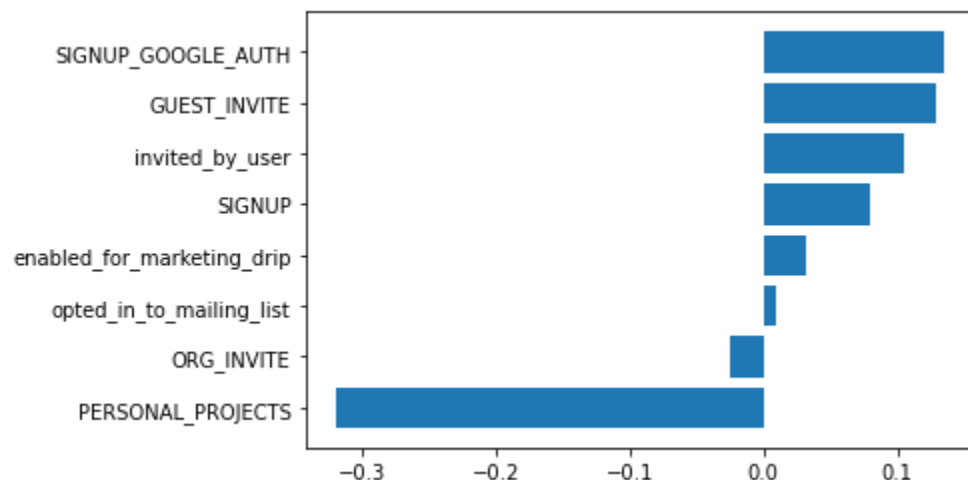
- A new feature was created Adopted\_users that has 1 for adopted and 0 for not. Out of 12,000 users 2,106 are considered adopted, then this list was used to add a new feature to the original user data set with whether they were adopted or not.

```
adopted_user = np.where(df.user_id.isin(adopted_users.user_id),1,0)
df['adopted_user'] = adopted_user
```

- The data has a 17% adoption rate.

The user dataframe is modified:

- Features kept are: creation\_source, opted\_in\_to\_mailing\_list, enabled\_for\_marketing\_drip, org\_id, invited\_by\_user\_id, adopted\_user
- invited\_by\_user\_id has nans which I am assuming to be customers not invited, so it was replaced with 1s if invited and 0 if nan
- X and y were created, y is adopted
- Then creation\_source was onehotencoded
- Train and test, then run logistic regression and run .coef\_ for importance and this is the result I got



The most important feature is creation\_source, followed by invited by user