Ben Deatsman Module 4 Lab Part 1 import pandas as pd In [515... df = pd.read_excel('mbta.xlsx', header=1, index_col=0) In [516... df Out [516... 2007-2007-2007-2007-2007-2007-2007-2007-2007-07 2011-01 2011-02 2011-03 2011-04 2011-05 mode 80 01 02 05 03 04 All Modes NaN NaN 1187.653 1245.959 1256.571 ... 1223.452 1302.414 NaN NaN NaN NaN NaN NaN NaN NaN by Qtr 2 4.000 3.600 40.000 4.300 4.900 5.800 6.521 6.572 5.469 ... 3.674 4.251 5.474 Boat 3.140 3.284 4.431 3 338.675 372.598 ... 352.162 354.367 350.543 357.519 355.479 334.958 346.234 380.399 380.446 385.289 376.317 335.819 339.867 Commuter 142.200 143.000 138.500 137.700 139.500 139.000 142.391 142.364 143.051 ... 128.396 125.463 134.374 134.169 136.140 135.581 13 Rail 5 Heavy Rail 435.294 448.271 458.583 472.201 474.579 477.032 471.735 461.605 499.566 ... 468.418 504.068 516.730 528.631 528.122 529.528 53 243.286 240.262 248.262 246.108 227.935 242.280 Light Rail 227.231 241.444 255.557 234.907 265.748 ... 198.450 219.886 225.776 221.865 Pct Chg / 7 0.020 -0.007 ... -0.040 0.114 -0.002 0.049 0.096 -0.037 0.004 -0.028 0.008 0.050 0.036 0.050 0.054 Private 8 4.772 4.417 4.574 4.542 4.768 4.722 3.936 3.946 4.329 ... 2.213 2.570 2.559 2.762 2.776 2.815 Bus 5.400 5.253 5.308 9 **RIDE** 4.900 5.000 5.500 5.400 5.609 ... 8.387 8.377 5.600 6.735 7.463 8.145 8.059 Trackless 10 12.757 14.393 ... 12.913 13.057 13.479 13.323 13.311 13.142 11.695 12.601 12.599 12.291 12.128 13.444 11.104 **Trolley** 11 TOTAL 1166.974 1191.639 1204.725 1247.105 1244.755 1246.129 1243.952 1223.323 1310.764 ... 1153.413 1220.663 1286.660 1313.283 1302.884 1292.085 128 $11 \text{ rows} \times 59 \text{ columns}$ In [517... df.isnull().sum() Out [517... mode 0 1 2007-01 2007-02 2007-03 2007-04 1 2007-05 1 2007-06 0 2007-07 1 1 2007-08 2007-09 2007-10 1 2007-11 1 0 2007-12 1 2008-01 1 2008-02 0 2008-03 2008-04 1 2008-05 1 2008-06 0 2008-07 1 1 2008-08 2008-09 0 2008-10 1 1 2008-11 2008-12 0 1 2009-01 1 2009-02 0 2009-03 2009-04 1 2009-05 1 0 2009-06 2009-07 1 1 2009-08 2009-09 2009-10 1 1 2009-11 2009-12 0 1 2010-01 1 2010-02 0 2010-03 2010-04 1 2010-05 1 0 2010-06 2010-07 1 2010-08 1 0 2010-09 2010-10 1 2010-11 1 2010-12 0 1 2011-01 1 2011-02 2011-03 0 2011-04 1 1 2011-05 2011-06 2011-07 2011-08 2011-09 2011-10 1 dtype: int64 In [518... print('There are', df.isnull().sum().sum(), 'missing values.') There are 39 missing values. In [519... df = df.drop([1, 7, 11])]print('The df size is', df.shape) The df size is (8, 59)In [520...] df = df.melt(id_vars = 'mode', var_name = 'date', value_name = 'num_riders_1000' df Out [520... mode date num_riders_1000 Boat 2007-01 0 4.000 1 Bus 2007-01 335.819 Commuter Rail 2007-01 142.200 Heavy Rail 2007-01 435.294 3 Light Rail 2007-01 227.231 4 459 Heavy Rail 2011-10 554.932 Light Rail 2011-10 237.768 460 461 Private Bus 2011-10 2.967 462 RIDE 2011-10 8.598 **463** Trackless Trolley 2011-10 12.297 464 rows × 3 columns In [521... print('The new dimensions are', df.shape) The new dimensions are (464, 3) In [522... df[['year', 'month']] = df['date'].str.split('-', expand=True) df.drop(columns=['date'], inplace=True) df.head() Out [522... mode num_riders_1000 year month 4.000 2007 01 Boat 335.819 2007 1 Bus 01 **2** Commuter Rail 142.200 2007 01 435.294 2007 Heavy Rail 01 227.231 2007 4 Light Rail 01 In [523... df[df['num_riders_1000'] == 40] mode num_riders_1000 year month Out [523... 16 Boat 40.0 2007 03 In [524... df['num_riders_1000'].replace(40.000, 4.000, inplace=True) df.groupby('mode', as_index=False).agg({'num_riders_1000': 'mean'}) Out [525... mode num_riders_1000 0 4.447276 Boat Bus 358.590448 1 137.351534 2 Commuter Rail 489.293483 3 Heavy Rail 232.988810 4 Light Rail Private Bus 5 3.352190 6 RIDE 6.603672 7 Trackless Trolley 12.125052 In [526... df['month'].replace({'01': 'Jan', '02': 'Feb', '03': 'Mar', '04': 'Apr', '05': 'May', '06': 'Jun', '07': 'Jul', '08': 'Aug', '09': 'Sep', '10': 'Oct', '11': 'Nov', '12': 'Dec'}, inplace=True) In [527... df_sort = df[df['month'] == 'Jan'] df_sort.groupby(['mode', 'month'], as_index=False).agg({'num_riders_1000': 'mean'}) Out [527... mode month num_riders_1000 0 Jan 3.3072 Boat 342.3718 1 Bus Jan 137.0198 2 Commuter Rail Jan 460.2818 3 Heavy Rail Jan 4 Light Rail 217.3920 Jan 3.4746 5 Private Bus Jan 6 RIDE Jan 5.9354 7 Trackless Trolley 12.4658 Jan In [528... df boat = df[df['mode'] == 'Boat'] df_boat.groupby(['mode', 'year'], as_index=False).agg({'num_riders_1000': 'sum'}).sort_values('num_riders_1000', ascending=False) Out [528... mode year num_riders_1000 Boat 2007 57.055 Boat 2010 52.931 Boat 2009 50.552 Boat 2008 50.349 Boat 2011 47.055 In [529... df_heavyrail = df[df['mode'] == 'Heavy Rail'] df_heavyrail.groupby(['mode', 'month'], as_index=False).agg({'num_riders_1000': 'mean'}).sort_values('num_riders_1000', ascending=False) Out [529... mode month num_riders_1000 11 Heavy Rail 517.70580 Sep 10 Heavy Rail 516.42520 Oct 0 Heavy Rail 501.91240 Apr 6 Heavy Rail 497.65420 Jun 5 Heavy Rail 495.26940 Jul 494.73740 8 Heavy Rail May 9 Heavy Rail Nov 489.65700 7 Heavy Rail 483.99600 Mar 481.39340 3 Heavy Rail Feb 477.42940 1 Heavy Rail Aug 460.28180 4 Heavy Rail Jan 2 Heavy Rail 446.59225 Dec Part 2 from completejourney_py import get_data In [531... cj_data = get_data() import pandas as pd transactions = cj_data['transactions'] demographics = cj_data['demographics'] products = cj_data['products'] cj_data.keys() In [532... Out[532... dict_keys(['campaign_descriptions', 'coupons', 'promotions', 'campaigns', 'demographics', 'transactions', 'coupon_redemptions', 'products']) transactions.columns.intersection(demographics.columns) Out[533... Index(['household_id'], dtype='object') In [534... df1 = transactions.merge(demographics, how='left', on='household_id') df1 Out [534... household_id store_id basket_id product_id quantity sales_value retail_disc coupon_disc coupon_match_disc week transaction_timestamp 35-44 35-4 0 1 900 31198570044 1095275 1 0.50 0.00 0.0 0.0 2017-01-01 11:53:26 2017-01-01 12:10:28 1 900 330 31198570047 9878513 1 0.99 0.10 0.0 0.0 1 45-2 1228 0.0 1 2017-01-01 12:26:30 406 31198655051 1041453 1 1.43 0.15 0.0 12 54 55-Unc 3 2017-01-01 12:30:27 906 31198705046 1020156 1 1.50 0.29 0.0 0.0 1 64 55-Unc 906 2.78 0.80 0.0 0.0 2017-01-01 12:30:27 2 1 4 31198705046 1053875 64 1469302 679 447 41453103606 14025548 1 0.79 0.20 0.0 0.0 2018-01-01 03:50:03 NaN 53 2018-01-01 04:01:20 1469303 2070 311 41453083334 909894 1 1.73 0.17 0.0 0.0 53 2070 2 5.00 2.98 0.0 0.0 53 2018-01-01 04:01:20 1469304 311 41453083334 933067 45-54 2018-01-01 04:01:20 1469305 2070 311 41453083334 1029743 1 2.60 0.29 0.0 0.0 53 2070 0.13 0.0 0.0 53 2018-01-01 04:01:20 1469306 1061220 1 1.19 311 41453083334 1469307 rows × 18 columns In [535... df1.notnull().sum() household_id Out [535... 1469307 store_id 1469307 basket_id 1469307 1469307 product_id quantity 1469307 sales_value 1469307 retail_disc 1469307 coupon_disc 1469307 coupon_match_disc 1469307 1469307 week transaction_timestamp 1469307 828850 age income 828850 home_ownership 607943 marital_status 687418 household_size 828850 household_comp 828850 kids_count 828850 dtype: int64 In [536... (transactions .merge(demographics, how='left', on='household_id') .groupby('age', as_index=False) .agg({'sales_value': 'sum'}) .sort_values(by='sales_value', ascending=False) Out [536... age sales_value 971822.19 **3** 45-54 **2** 35-44 724357.40 **1** 25-34 453372.46 176600.84 65+ **4** 55-64 173153.70 **0** 19-24 125673.05 pizza_filter = products['product_type'].str.contains('pizza', case=False, na=False) (products[pizza_filter] .merge(cj_data['transactions'] , how='inner', on='product_id') .groupby(['product_id', 'product_type', 'product_category'], as_index=False) .agg({'sales_value': 'sum'}) .sort_values(by='sales_value', ascending=False) Out [595... product_id product_type product_category sales_value PIZZA/TRADITIONAL 97 944139 FROZEN PIZZA 1344.50 70 906838 PIZZA/PREMIUM FROZEN PIZZA 1263.78 357 12648296 PIZZA/TRADITIONAL FROZEN PIZZA 1230.36 113 PIZZA/TRADITIONAL 1125.18 969568 FROZEN PIZZA 85 925626 PIZZA/ECONOMY FROZEN PIZZA 1017.60 PIZZA/ECONOMY 29 610240 FROZEN PIZZA 0.99 1231548 PIZZA/ECONOMY 203 FROZEN PIZZA 0.99 232 2063332 0.67 PIZZA/ECONOMY FROZEN PIZZA 266 6876820 PIZZA/ECONOMY FROZEN PIZZA 0.59 10341942 QSR FD: COLD PIZZA 334 PREPARED FOOD 0.00 394 rows × 4 columns pizza_cat_filter = products['product_category'].str.contains('pizza', case=False, na=False) & products['product_type'].str.contains('snack|appetizer', In [597... (products[pizza_cat_filter] .merge(cj_data['transactions'] , how='inner', on='product_id') .groupby(['product_id', 'product_type', 'product_category'], as_index=False) .agg({'sales_value': 'sum'}) .sort_values(by='sales_value', ascending=False) Out [597... product_id product_type product_category sales_value 37 1990.50 907631 SNACKS/APPETIZERS FROZEN PIZZA 34 890695 SNACKS/APPETIZERS 1686.28 FROZEN PIZZA 18 845193 SNACKS/APPETIZERS FROZEN PIZZA 1507.54 22 856252 SNACKS/APPETIZERS FROZEN PIZZA 1133.90 829291 SNACKS/APPETIZERS 15 FROZEN PIZZA 1020.61 FROZEN PIZZA 2215639 SNACKS/APPETIZERS 113 1.89 186 15831485 SNACKS/APPETIZERS FROZEN PIZZA 1.79 99 1335834 SNACKS/APPETIZERS FROZEN PIZZA 1.25 2017119 SNACKS/APPETIZERS 110 FROZEN PIZZA 0.48 108 1917329 SNACKS/APPETIZERS FROZEN PIZZA 0.38 188 rows × 4 columns In [601... | pb_filter = products['product_type'].str.contains('peanut butter', case=False, na=False) pb_products = products[pb_filter] print('There are', pb_products['product_id'].nunique(), 'unique peanut butter products') There are 144 unique peanut butter products pb_merge = pb_products.merge(cj_data['transactions'] , how='inner', on='product_id') pb_merge['transaction_month'] = pd.to_datetime(pb_merge['transaction_timestamp']).dt.month (pb_merge.groupby('transaction_month', as_index=False) .agg({'sales value': 'sum'}) .sort_values(by='sales_value', ascending=False)) Out [588... transaction_month sales_value 11 12 1019.08 9 8 905.36 7 8 895.10 11 882.32 10 10 9 878.89 5 6 865.15 7 6 863.38 0 1 863.36 5 850.37 4 3 4 790.18 2 1 755.76 3 2 733.60 January has the highest sales value for products containing peanut butter in the product type