ausikChattapadhyay_DSC540_Week7-8Ex SC 540 Week 7-8 ausik Chattapadhyay tivity: For this assignment you need to complete 8 of the following exercises against this data.	
apter 7 ilter out missing data ill in missing data emove duplicates ransform data using either mapping or a function eplace values	
iscretization and Binning Ianipulate Strings apter 8 reate hierarchical index ombine and Merge Datasets (you will have to either create a new dataset from your existing data or create a relationship between the data I have provided) eshape	
eshape ivot the data apter 10 rouping with Dicts/Series rouping with Functions rouping with Index Levels	
plit/Apply/Combine ross Tabs apter 11 onvert between string and date time enerate date range	
requencies and date offsets onvert timestamps to periods and back eriod Frequency conversions Load the necessary libraries. port numpy as np port pandas as pd	
II) GENDER AGE COUNTRY COUNTY Grand and Grande SIZED COMMENTS DRESS 113 DAY IDAIN	Coordin
OUT? GENDER AGE COUNTRY COUNTY, ETC Bar wrappers\t(a.k.a. bar wrap	(84
90272840 No Male 40 us or MEH DESPAIR JOY MEH NAN Raisins can go to hell and gold NAN Sunday NAN 1.0 NAN NAN 90272841 No Male 23 usa exton pa JOY DESPAIR JOY DESPAIR NAN NAN NAN Gold NAN Friday NAN 1.0 NAN NAN ws x 120 columns	
dex(['Internal ID', 'Q1: GOING OUT?', 'Q2: GENDER', 'Q3: AGE', 'Q4: COUNTRY',	
<pre>andyDF = candyDF.rename(columns = {'Q1: GOING OUT?' : 'going_out', 'Q2: GENDER' : 'gender', 'Q3: AGE': 'age',</pre>	
<pre>indyDF.drop(columns = ['Internal ID', 'Unnamed: 113', 'Click Coordinates (x, y)'], inplace = True) indyDF.shape 160, 117) Handling null values indyDF.dropna(subset = ['going_out', 'gender', 'age', 'country', 'area'], how = 'all', inplace = True) indyDF.reset_index(drop = True, inplace = True)</pre>	
ndyDF.shape 135, 117) Formatting Columns Going Out Column ndyDF.going_out = candyDF.going_out.fillna('Not Sure') ndyDF.going_out.unique() cay(['No', 'Not Sure', 'Yes'], dtype=object)	
Gender Column hdyDF.gender.value_counts() Le 1467 hale 839 th rather not say 83 her 30 he: gender, dtype: int64 Adding NaN genders to type 3 - I'd rather not say, as it seems to be similar to unknown or NA	
hecking for spaces in text - found none hedyDF.gender = candyDF.gender.fillna("I'd rather not say") hedyDF.gender.value_counts() Le 1467 hale 839 hi rather not say 99 her 30 he: gender, dtype: int64	
Tay(['USA', 'USA', 'us', 'usa', nan, 'canada', 'Canada', 'Us', 'Us', 'Wurica', 'United States', 'uk', 'United Kingdom', 'united states', 'Usa', 'United States', 'Injustration', 'United States of America', 'UAE', 'England', 'UK', 'canada', 'Mexico', 'United States', 'us.a.', 'USAUSAUSA', 'america', 35, 'france', 'United States of America', 'U.S.A.', 'finland', 'unhinged states', 'The United States of America', 'United States', 'The United States', 'North Carolina', 'Unied States', 'Netherlands', 'germany', 'Europe', 'Earth', 'U S', 'u.s.', 'U.K.', 'Costa Rica', 'The United States of America', 'unite states', 'U.S.', '46, 'cascadia', 'Australia', 'insanity lately', 'Greece', 'USA? Hard to tell anymore', '"merica", 'usas', 'Pittsburgh', 45, 'United State', 32, 'France', 'australia', 'A', 'Can', 'Canae', 'New York', 'Trumpistan', 'Ireland', 'United States', 'Korea', 'California', 'Japan', 'USa', 'South africa', 'United States', 'America', 'Inserend to be from Canada, but I am really from the United States.'.	
'I pretend to be from Canada, but I am really from the United States.', 'Usa', 'Uk', 'Iceland', 'Germany', 'Canada', 'Scotland', 'UK', 'Denmark', 'United Stated', 'France', 'Switzerland', 'AhemAmerca', 'UD', 'Scotland', 'South Korea', 'New Jersey', 'CANADA', 'Indonesia', 'United ststes', 'America', 'The Netherlands', 'United Statss', 'endland', 'Atlantis', 'murrika', 'USAI USAI USAI', 'UsAA', 'Alaska', 'United States', 'soviet canuckistan', 'N. America', 'Singapore', 'USSA', 'China', 'Taiwan', 'Ireland', 'hong kong', 'spain', 'Sweden', 'Hong Kong', 'U.S.', 'Narnia', 'u s a', 'United Statea', 'united ststes', 1, 'subscribe to dm4uz3 on youtube', 'United kingdom', 'USA USA USAI!!!', "I don't know anymore", 'Fear and Loathing'], dtype=object) A 699 ited States 497	
ited States 497 217 anda 179 a 139 ited kingdom 1 A USA USA!!!! 1 don't know anymore 1 ance 1 ar and Loathing 1 ne: country, Length: 129, dtype: int64	
Australia', Europe', I pretend to be from Canada, but I am really from the United States.', Europe', South Korea', Fouth Korea', Fouth africa', Frumpistan', Europistan', Eustralia', Europistan', Eustralia', Eus	
as a', a.s.', a.s.a.', a.s.a.a.', a.s.a.a.', a.s.a.a.', a.s.a.a.', a.s.a.a.a.', a.s.a.a.a.a.a.a.a.a.a.a.a.a.a.a.a.a.a.a	
<pre>is', isa', isas') A = [x for x in candyDF.country if (('u' in str(x) or 'U' in str(x)) and</pre>	
'canada ', 'Mexico', 'america', 35, 'france', 'finland', 'Canada ', 'North Carolina ', 'Netherlands', 'germany', 'Europe', 'Earth', 'Costa Rica', 46, 'cascadia', 'Australia', 'insanity lately', 'Greece', "'merica", 45, 32, 'France', 'australia', 'A', 'Can', 'Canae', 'New York', 'Ireland', 'Korea', 'California', 'Japan', 'South africa', 'Iceland', 'Germany', 'Canada'', 'Scotland', 'Denmark', 'France ', 'Switzerland', 'AhemAmerca', 'Scotland', 'South Korea', 'New Jersey', 'CANADA', 'Indonesia', 'America', 'The Netherlands', 'endland', 'Atlantis', 'Alaska', 'N. America', 'Singapore', 'China', 'Taiwan', 'Ireland ', 'hong kong', 'spain', 'Sweden', 'Hong Kong', 'Narnia', 1, 'United kingdom', "I don't know anymore", 'Fear and Loathing'], dtype=object)	
<pre>removing duplicates , wrongly or differently names and updating ndyDF.country = candyDF.country.replace(to_replace = ['america','AhemAmerca',"'merica",'North Carolina ',</pre>	
<pre>ndyDF.country = candyDF.country.replace(to_replace = 'The Netherlands', value = 'Netherlands') ndyDF.country = candyDF.country.replace(to_replace = 'australia', value = 'Australia') ndyDF.country = candyDF.country.replace(to_replace = [1,"I don't know anymore",32,45,35,46,'Fear and Loathing','insanity lately'],</pre>	
stralia 7 pan 5 pan 6 potland 4 potland 3 pland 3 pland 3 pland 2 pland 2 pland 2 pland 3 plan	
agapore 1 Pland 1 Plan	
rea 1 ance 1 sece 1 sta Rica 1 rth 1 rope 1 nland 1 ance 1 rnia 1 ne: country, dtype: int64 Grouping Dataset to 3 Countries - USA, Canada, Others	
creating a dataset Other from existing one which we will merge later ner = [x for x in candyDF.country.unique()] ner.remove('USA') ner.remove('Canada') ner United Kingdom', Mexico', Others', France',	
irance', finland', tetherlands', Germany', Surope', Sarth', Sosta Rica', Australia', Greece', France', Greece',	
Scotland ', South Korea', Indonesia', Atlantis', Singapore', China', Caiwan', Ireland ', Iong kong', Spain', Sweden', Hong Kong', Hong Kon	
'Q6 Anonymous brown globs that come in black and orange wrappers\t(a.k.a. Mary Janes)', 'Q6 Any full-sized candy bar', 'Q6 Black Jacks', 'Q6 Bonkers (the candy)', 'Q6 York Peppermint Patties', 'Q7: JOY OTHER', 'Q8: DESPAIR OTHER', 'Q9: OTHER COMMENTS', 'dress', 'day', 'media_DailyDish', 'media_Science', 'media_ESPN', 'media_Yahoo'], dtype='object', length=117) **Converting Datatype** ndyDF = candyDF.astype({'going_out': 'category', 'gender': 'category', 'country': 'category', 'dress': 'category',	
day': category day	
Method to Convert 4 Columns into one f meltl(row): for c in data.columns: if row[c] == 1: return c Checking Media columns which we will merge into one	
ta = candyDF[candyDF.columns[-4:]] media_DailyDish	
4 NaN 1.0 NaN NaN 30 NaN NaN NaN NaN 31 NaN 1.0 NaN NaN 32 NaN 1.0 NaN NaN 33 NaN NaN NaN NaN 34 1.0 NaN NaN NaN	
5 rows × 4 columns w_col = data.apply(melt1, axis = 1) Adding newly created column ndyDF['media_preference'] = new_col ropping old columns ndyDF.drop(columns = ['media_DailyDish', 'media_Science', 'media_ESPN', 'media_Yahoo'], inplace = True)	
ndyDF.media_preference.value_counts(dropna = False) dia_Science	
rsonal_info_cols = candyDF.columns[:6] estionare_cols = candyDF.columns[5:] ndyDF.columns dex(['going_out', 'gender', 'age', 'country', 'area', 'Q6 100 Grand Bar',	
and pdf ['responses'] = responses and pdf head(3) going_out gender age country area area area and orange wrappers\t(a.k.a. candy) Mary Janes) Q6 Anonymous Q6 Any G7 Q8 Q6 York Q7 Q8 Q7 Q8 Q8 York Q7 Q8 Q9 OTHER COMMENTS Area and orange wrappers\t(a.k.a. candy) Mary Janes) Bottom line is Twix is White	
No Male 44 USA NM MEH DESPAIR JOY MEH DESPAIR DESPAIR DESPAIR DESPAIR Mounds NaN really the and sondy candy gold w Not Sure Male 49 USA Virginia NaN NaN NaN NaN NaN NaN NaN NaN NaN N	ne
ray([44, 49, 40, 23, nan, 53, 33, 43, 56, 64, 37, 59, 48, 54, 36, 45, 25, 34, 35, 38, 58, 50, 47, 16, 52, 63, 65, 41, 27, 31, 61, 46, 42, 62, 29, 39, 32, 28, 69, 67, 30, 22, 26, 51, 70, 24, 18, 19, 57, 60, 66, 12, 55, 72, 21, 11, 9, 68, 20, 6, 10, 71, 90, 13, 99, 7, 88, 39.4, 74, 102, 17, 15, 8, 75, 'See question 2', 14, 100, 76, 77, 73, 70.5, 1, 4], dtype=object) Data Type Conversion adyDF.drop_duplicates(inplace=True) = pd.to_numeric(candyDF['age'],downcast='float',errors='ignore')	
<pre>pd.to_numeric(candyDF['age'],downcast='float',errors='ignore') pd.to_numeric(s,downcast='float',errors='coerce') ndyDF['age'].unique() ndyDF.replace(candyDF['age'],s,inplace=True) ndyDF['age'].replace(['old enough', '45-55','24-50','?','no','Many','hahahahaha','older than dirt','Enough',</pre>	
44.0 49.0 40.0 23.0 NaN ne: age, dtype: float64 m=candyDF.columns n dex(['going_out', 'gender', 'age', 'country', 'area', 'Q6 100 Grand Bar',	
'Q6 Anonymous brown globs that come in black and orange wrappers\t(a.k.a. Mary Janes)', 'Q6 Any full-sized candy bar', 'Q6 Black Jacks', 'Q6 Bonkers (the candy)', 'Q6 White Bread', 'Q6 Whole Wheat anything', 'Q6 York Peppermint Patties', 'Q7: JOY OTHER', 'Q8: DESPAIR OTHER', 'Q9: OTHER COMMENTS', 'dress', 'day', 'media_preference', 'responses'], dtype='object', length=115)	