us re re <r th="" us<=""><th>rl = 'https://datawebscrapping.vercel.app/super_store.html' sr = dotenv_values('/credential.env') esponse = requests.get(url ,('user_agent':usr['usr']}) esponse Response Response Sing bs4 to extract data from website oup = BeautifulSoup(response.text , 'html.parser') rticles = soup.find_all('div', ('class':'article'))</th></r>	rl = 'https://datawebscrapping.vercel.app/super_store.html' sr = dotenv_values('/credential.env') esponse = requests.get(url ,('user_agent':usr['usr']}) esponse Response Response Sing bs4 to extract data from website oup = BeautifulSoup(response.text , 'html.parser') rticles = soup.find_all('div', ('class':'article'))
San co fo	<pre>rticles = soup.find_all('div', {'class':'article'}) aving Each Articles in new text files (automatic) ount = 0 or i in articles: if i: with open(f"article{count+1}.txt", 'w') as arti:</pre>
Pr Sci ta #t co co co [< < < < < < < < < < < < < < < < < < <	trap Tables ables = soup.find_all('table',('class':'dataframe sample_1000_r'))[0] tables olumns = tables.find_all('th') olumns = tables.find_all('th') tch>Row ID, tch>Corder ID, tch>Corder Date, tch>Ship Date, tch>Ship Mode, tch>Segment, tch>Segment, tch>Segment, tch>Segment, tch>Country, tch>Country
< < < < < < < < < < < < < < < < < < <	Country,
· · · · · · · · · · · · · · · · · · ·	Order Date', Ship Date', Ship Mode', Customer ID', Segment', Country', Country', Country', City', State', Postal Code', Region', Product ID', Category', Sub-Category', Product Name', Sales', Canantity', Profit']
ro ro	<pre>cows = tables.find_all('tr') cows[0] #checking tr style="text-align: right;"></pre>
<t <="" ro<="" t="" td=""><td><pre>th>sub-Category th>roduct Name th>roduct Name th>sales th>ch>Quantity th>ch>quantity th>ch>quantity th>roduct Name th>ch>quantity th>roduct Name th>quantity th>quantit</pre></td></t>	<pre>th>sub-Category th>roduct Name th>roduct Name th>sales th>ch>Quantity th>ch>quantity th>ch>quantity th>roduct Name th>ch>quantity th>roduct Name th>quantity th>quantit</pre>
	'7/20/2017', '7/26/2017', 'Standard Class', 'CS-11950', 'Consumer', 'United States', 'Chicago', 'Illinois', '60610', 'Central', 'OFF-ST-10002974', 'Office Supplies', 'Storage', 'Trav-L-File Heavy-Duty Shuttle II, Black', '\$69.71', '2', '88.71']
df df #	5523 CA-2016-110982 06-05-2016 06-07-2016 First Class CK-12205 Consumer United States Santa Clara California 95051 West OFF-BI-10000829 Office Supplies Binders Avery Non-Stick Binders \$21.55 6 \$7.00 6415 CA-2017-142671 11-10-2017 11/14/2017 Second Class DR-12940 Home Office United States Hollywood Florida 33021 South OFF-BI-10004099 Office Supplies Binders GBC VeloBinder Strips \$11.52 5 \$-7.68 69487 CA-2017-130505 10-12-2017 10-12-2017 Same Day NF-18385 Consumer United States Wheeling West Virginia 26003 East FUR-TA-10001932 Furniture Tables Chromcraft 48" x 96" Racetrack Double Pedestal \$673.34 3 \$-76.95
rer #m df df cc Ra	Data Cleaning/Preprocessing/Cleansing/Transformation smove un_necessary columns/variables making all variables name to lower case with '_' f.columns = df.columns.str.lower().str.replace(" ", "_") f.drop(inplace =True , columns =['row_id', 'postal_code', 'country']) f.info() class 'pandas.core.frame.DataFrame'> snageIndex: 1000 entries, 0 to 6999 state columns (total 16 columns):
0 1 2 3 4 5 6 7 8 9 1 1 1	Column
dt me: Th	bypes: object(16) emory usage: 125.1+ KB ne data is not in the correct format: 1. order_date and ship_date are objects instead of datetime. 2. sales, quantity, and profit should be integers or floats instead of objects. rst Change datetime data format 7/20/2017 ate_cols = df[['order_date', 'ship_date']] or col in date_cols: df[col] = df[col].str.replace(',', '-').str.strip() df[col] = pd.to_datetime(df[col]), format = '%m-8d-%Y')
pr df <c Ra Da #</c 	rint('process done') f.info() class 'pandas.core.frame.DataFrame'> class 'pandas.co
1 1 1 1 dt me:	7 state 1000 non-null object 8 region 1000 non-null object
df df nu	f.sales = df.sales.apply(lambda x :float(re.findall('\d+\.\d+',x)[0])) f.profit = df.profit.apply(lambda x :re.sub(r"[^\d]","",x)) f.profit.replace('-',0,inplace =True) um_col =df[['profit','quantity']] or i in num_col: df[i] = pd.to_numeric(df[i]) f.info()
Ra Da # - 0 1 2 3 4 5 6 7 8 9 1 1	9 product_id 1000 non-null object 10 category 1000 non-null object 11 sub-category 1000 non-null object
1 1 1 the state of	12 product_name 1000 non-null object 13 sales 1000 non-null float64 14 quantity 1000 non-null int64 15 profit 1000 non-null float64 15 profit 1000 non-null float64 16 profit 1000 non-null float64 17 profit 1000 non-null float64 18 profit 1000 non-null float64 19 profit 1000 non-null float64 10 profit 1000 non-null float64 20 profit 1000 non-null float64 21 profit 1000 non-null float64 22 profit 1000 non-null int64 23 profit 1000 non-null int64 24 profit 1000 non-null float64 25 profit 1000 non-null int64 26 profit 1000 non-null float64 26 profit 1000 non-null int64 27 profit 1000 non-null int64 28 profit 1000 non-null int64 29 profit 1000 non-null int64 20 profit 100 non-nul
Rer Exp Sal	Inder HTML: 100% If [00:01&kt00:00, 1.37s/it] port report to file: 100% If [00:00&kt00:00, 35.99it/s] Data Analysis and Visualization ales Trends: How sales trends changed over time? Are there any noticeable patterns or seasonality? # these are custom color code 1 use chatgpt to generate these color code Warm Shades arm palette = ['#ff9999', '#ff6666', '#ff4ddd', '#ff3333', '#ff0000']
# co # ea # su # pa dffdf # # i	<pre>conglates = ['#ff9999', '#ff6666', '#ff4d4d', '#ff3333', '#ff0000'] Cool Blues col_blues palette = ['#cce5ff', '#99c2ff', '#66b3ff', '#3399ff', '#007bff'] Earth Tones arth_tones_palette = ['#d9e3d8', '#b8d1c6', '#9ac9b3', '#6cb4a0', '#4e8c6f'] Sunset Colors unset_palette = ['#ffcc99', '#ff3660', '#ff9933', '#ff6000', '#cc3300'] Pastel Shades astel_palette = ['#f7c6c7', '#f7a6a8', '#f78180', '#f75050', '#d72027'] f('month') = df.order_date.dt.month # extraing month and year from date feature f('year') = df.order_date.dt.year int32 to int 64 f(['month', 'year']] = df[['month', 'year']].astype('int64')</pre>
df <c Ra Da #</c 	Class pandas.core.frame.DataFrame >
1 1 1 1 1 1 dt me:	9 product_id 1000 non-null object 10 category 1000 non-null object 11 sub-category 1000 non-null object 12 product_name 1000 non-null object 13 sales 1000 non-null ifloat64 14 quantity 1000 non-null int64 15 profit 1000 non-null int64 16 month 1000 non-null int64 17 year 1000 non-null int64 17 year 1000 non-null int64 18 privot_table emory usage: 140.8+ KB privot_table f.pivot_table(index ='year', columns= 'month', values= 'sales', aggfunc='sum').reset_index() nonth year 1 2 3 4 5 6 7 8 9 10 11 12 0 2014 98773 321.28 3688.71 1254.39 1982.63 3952.49 2232.41 1272.92 3667.53 282.68 7 4009.14 4557.28
	1 2015 919.45 221.60 2011.20 2987.48 629.76 3314.64 3195.15 3243.56 4824.11 688.05 6429.57 7224.19 2 2016 2358.98 1241.72 3225.47 880.06 3052.49 2978.58 3474.89 4716.15 5303.58 1483.89 5478.64 6335.38 3 2017 2926.11 1763.55 4605.40 3115.98 3286.28 5015.95 4131.55 3916.49 6227.49 3793.36 7557.62 8411.40 ales_by_month = df.groupby('month').agg(Sales = ('sales', 'sum')).reset_index() ales_by_month Nonth Sales 1 7192.57 1 2 3548.15
3 4 5 6 7 8 9 10 11	4 5 8951.16 5 6 15261.66 6 7 13034.00 7 8 13149.12 8 9 20022.71 9 10 8792.17 10 11 23474.97
0 1 2 3 Sal	year Sales 2014 30753.38 2015 35688.76 2016 40529.83 2017 54751.48 sles by Month lt.figure(figsize=(12, 6)) ns.lineplot(data=sales_by_month, x='month', y='Sales', marker='o', palette='#d9e3d8') lt.title('Sales Treads by Month')
pl pl pl	1t.ylabel('Sales') lt.yaticks(rotation=45) lt.grid(True) 1t.show() Sales Trends by Month
Sales	10000
pl sn pl pl	Month lt.figure(figsize=(12, 5)) ns.lineplot(sales_by_year , x = 'year', y ='Sales' , marker ='o', markerfacecolor = 'black') lt.grid(True) lt.title('Sales Trend by Year') lt.xlabel('Year') lt.xlabel('Year') sext(0, 0.5, 'Sales') Sales Trend by Year
Sales	50000 45000 40000
• (d	30000 2014.0 2014.5 2015.0 2015.5 2016.0 2016.5 2017.0 sales Performance by Region: • Which regions have the highest and lowest sales? How does sales performance vary across different regions? df.region.value_counts(normalize=True) * 100).reset_index() region proportion East 29.1
2 3 sa sa 1 3	West 290 Central 258 South 161 ales_by_region = df.groupby('region').sales.sum().reset_index().sort_values(by = 'sales', ascending =False) ales_by_region
pr pr 3 1 0	South 21389.26 refit_by_region = df.groupby('region').profit.sum().reset_index().sort_values(by = 'profit', ascending =False) region profit West 9942.45 East 8116.86 Central 5626.90 South 4166.75 att.figure(figsize = (14,5))
sn pl pl pl sn pl pl	It.subplot(1,2,1) ns.barplot(sales_by_region, x = 'sales' ,y = 'region', hue = 'sales' , palette=cool_blues_palette) It.xlabel('Sales_By_Region') It.xlabel('Region') It.subplot(1,2,2) ns.barplot(profit_by_region, x = 'profit' ,y = 'region', hue = 'profit' , palette=cool_blues_palette) It.xlabel('Profit By Region') It.xlabel('Profit') It.ylabel('Region') ext(0, 0.5, 'Region') Sales By Region
Region	East - West - East - East - Central - Central - Sales 21389.26 39330.36 Profit 4166.75 5626.9
se se	South 49839.57 South 9942.45 10000 20000 30000 40000 50000 0 2000 4000 6000 8000 10000 Sales What are the sales and profit contributions from different customer segments? Which segment is the most profitable? egment_contribution =df.pivot_table(index = 'segment', values = ['sales', 'profit'], aggfunc = 'sum').reset_index() segment profit sales
se pl pl pl	
Total Value	80000 -
df	20000 - Constitute Constitute Segment F. head () Order date while date while mode customer id a segment of the state region of product id a category subcategory and product name sales quantity profit month year.
1 2 3 4 cat	order_ide
0 1 2	category sales Furniture 66927.72 Office Supplies 48086.39 Technology 46709.34 rofit_cat = df.groupby('category').profit.sum().reset_index().sort_values(by ='profit',ascending =False) rategory profit rategory profit
0 pr pr 6 0 3	Office Supplies 1357001 Furniture
12 5 14 11 9 1 10 7 2 15 8	2 Paper 3229.29 5 Chairs 3063.80 4 Storage 2704.90 1 Machines 2186.91 9 Furnishings 1476.56 1 Appliances 1244.22 0 Labels 971.31 7 Envelopes 752.72 2 Art 522.71 5 Supplies 190.21
pl pl sn pl pl pl pl	<pre>A Bookcase -1045.95 6 Tables -4411.24 1t.figure(figsize = (28,18)), 1t.subplot(2,2,1) ns.barplot(x , x = 'category', y = 'sales') 1t.title('Sales By Category') 1t.subplot(2,2,2) ns.barplot(sub_cat_subCat , x = 'sub-category', y='sales' ,palette= cool_blues_palette) 1t.xticks(rotation = 45) 1t.title('Sales By Sub-Category') 1t.subplot(2,2,3)</pre>
sn pl sn pl pl Te	ns.barplot(profit_cat , x = 'category', y='profit' ,palette= cool_blues_palette) lt.subplot(2,2,4) ns.barplot(profit_Scat , x = 'sub-category', y='profit' ,palette= cool_blues_palette) lt.xticks(rotation =45) lt.ticks(rotation =45) lt.ticle('Profit_By Sub-Category') ext(0.5, 1.0, 'Profit_By Sub-Category') Sales By Category Sales By Sub-Category 30000 -
sales	25000 - 40000 - 20000 - 20000 - 20000 - 15000 -
	Profit By Category Profit By Sub-Category Profit By Sub-Category 14000
	12000 - 10000 - 8000 -
	Technology Office Supplies category Furniture Codiffer Supplies Report Chart Scott State Burker Report Chart Scott State Burker State St
	able scapping without bs4 using pandas andas only scrapes/extracts tables from webpages with less control over data