In [130]: import pandas as pd import numpy as np In [148]: | filePath = 'G:/PRCTICE/Python/' inputFileName = 'SalesOrderHeader.xlsx' outputFileName = 'SalesOrderHeader_Analysis_Result.xlsx' inputFile = filePath + inputFileName outputFile = filePath + outputFileName In [149]: data = pd.read_excel(inputFile) Chage index name and serial: data.index.name = "Index" In [150]: data.index = data.index + 1**Substring:** In [151]: SalesOrderCode = data['SalesOrderNumber'].astype(str).str[:3] Insert new column in a specific position: In [152]: data.insert(loc=7, column='SalesOrderCode', value=SalesOrderCode) In [153]: data.head(3) Out[153]: SalesOrderID RevisionNumber OrderDate DueDate ShipDate Status OnlineOrderFlag SalesOrderCode SalesOrderNumber **PurchaseOrderNumbe** Index 2011-06-2011-05-2011-06-SO4 43659 5 0 SO43659 PO522145787 1 2011-05-2011-06-2011-06-5 SO4 PO18850127500 2 43660 0 SO43660 2011-05-2011-06-2011-06-5 PO18473189620 3 43661 0 SO4 SO43661 31 12 07 3 rows × 27 columns Replace: data['DueDate']=data['DueDate'].replace('1970-01-01', np.nan) Apply IN, Between and & with where clause: In [155]: condition1 = data['SalesOrderCode'].isin(['S04','S05']) condition2 = data['OrderDate'].between('2011-05-31','2013-05-31') data[condition1 & condition2].head(2) Out[155]: SalesOrderID RevisionNumber OrderDate DueDate ShipDate Status OnlineOrderFlag SalesOrderCode SalesOrderNumber **PurchaseOrderNumbe** Index 2011-05-2011-06-2011-06-5 0 SO4 43659 SO43659 PO522145787 1 2011-05-2011-06-2011-06-2 43660 5 0 SO4 SO43660 PO1885012750(2 rows × 27 columns data[condition1 & condition2].tail(2) In [156]: Out[156]: SalesOrderID OrderDate DueDate ShipDate Status RevisionNumber OnlineOrderFlag SalesOrderCode SalesOrderNumber **PurchaseOrderNumbe** Index 2013-05-2013-06-2013-06-7543 51201 SO5 SO51201 NaN 12 2013-05-2013-06-2013-06-5 SO5 7544 51202 SO51202 Na۱ 12 2 rows × 27 columns data.where(condition1 & condition2).head(2) In [211]: Out[211]: RevisionNumber OrderDate DueDate ShipDate OnlineOrderFlag CreditCardID SL# SalesOrderID Status SalesOrderCode SalesOrderNumber NaN NaN NaN NaT NaT NaT NaN NaN NaN NaN NaN 2011-05-2011-06-2011-06-2.0 8.0 5.0 0.0 SO4 SO43660 43660.0 5618.0 12 31 2 rows × 28 columns data.query("RevisionNumber in [8,7]").head(3) In [212]: Out[212]: SalesOrderID OrderDate **DueDate** ShipDate OnlineOrderFlag SalesOrderCode SalesOrderNumber CreditCardID SL# RevisionNumber Status 2011-05-2011-06-2011-06-8 SO4 SO43659 43659 16281.0 31 12 2011-05-2011-06-2011-06-8 43660 SO4 SO43660 5618.0 31 12 2011-05-2011-06-2011-06-8 SO4 SO43661 2 3 43661 0 1346.0 31 12 3 rows × 28 columns group by with and without meging index column: data_groupedby = data.groupby(['SalesOrderCode'], as_index=False).agg('count') In [157]: data_groupedby.head(3) In [158]: Out[158]: OnlineOrderFlag DueDate ShipDate SalesOrderCode SalesOrderID RevisionNumber OrderDate Status SalesOrderNumber **PurchaseOrderNumber** SO4 6341 6341 6341 6341 6341 6341 6341 1758 6341 SO5 10000 10000 10000 10000 10000 10000 10000 10000 1147 SO6 10000 10000 10000 10000 10000 10000 10000 10000 720 3 rows × 27 columns In [159]: data_groupedby = data.groupby(['SalesOrderCode']).agg('count') In [160]: data_groupedby.head() Out[160]: ShipDate Status SalesOrderID RevisionNumber OrderDate DueDate OnlineOrderFlag SalesOrderNumber PurchaseOrderNumber SalesOrderCode 6341 6341 **SO4** 6341 6341 6341 6341 6341 6341 1758 SO₅ 10000 10000 10000 10000 10000 10000 10000 1147 10000 **SO6** 10000 10000 10000 10000 10000 10000 10000 10000 720 **SO7** 5124 5124 5124 5124 5124 5124 5124 5124 181 4 rows × 26 columns **Export every column with distinct value to excel:** writer = pd.ExcelWriter(outputFile, engine='xlsxwriter') In [175]: In [176]: data_groupedby.to_excel(writer, 'Tally') In [177]: total = len(data.index) data.insert(loc=0, column='SL#', value=data.index) In [178]: **for** column **in** data: if(data[column].name == 'SalesOrderCode' or data[column].name == 'SL#'): continue else: col_names = [data[column].name] my_df = pd.DataFrame(columns = col_names) my_df[data[column].name] = data[column].unique().tolist() if((my_df[column].nunique()/total) <= 0.02):</pre> grouped_data = data[['SalesOrderCode',data[column].name,'SL#']] grouped_data = grouped_data.groupby(['SalesOrderCode',data[column].name], as_index=False).agg('count') #Rename index column here. grouped_data = grouped_data.rename(columns={'SL#': 'Count'}) if(len(grouped_data) > 0): grouped_data.to_excel(writer, sheet_name = data[column].name) In [179]: writer.save() writer.close() Set, Reset index: data.set_index("SL#",inplace= True) In [184]: In [185]: data.head() Out[185]: SalesOrderID RevisionNumber OrderDate DueDate ShipDate **Status** OnlineOrderFlag SalesOrderCode SalesOrderNumber PurchaseOrderNumber SL# 2011-06-2011-05-2011-06-5 SO4 PO522145787 43659 SO43659 12 2011-05-2011-06-2011-06-SO4 PO18850127500 2 43660 SO43660 2011-05-2011-06-2011-06-3 43661 5 SO4 SO43661 PO18473189620 2011-05-2011-06-2011-06-5 43662 SO4 SO43662 PO18444174044 31 2011-05-2011-06-2011-06-5 43663 SO4 SO43663 PO18009186470 12 31 5 rows × 27 columns data.reset index(drop=False, inplace=True) In [187]: In [194]: data.head(2) Out[194]: DueDate ShipDate SalesOrderID RevisionNumber OrderDate Status OnlineOrderFlag SalesOrderCode SalesOrderNumber CreditCardID 2011-05-2011-06-2011-06-8 43659 0 SO4 SO43659 16281.0 31 12 2011-06-2011-06-2011-05-8 SO4 SO43660 43660 5618.0 31 12 2 rows × 28 columns Show data by loc & iloc: data.iloc[0:1] In [198]: Out[198]: OnlineOrderFlag SL# SalesOrderID OrderDate DueDate ShipDate Status SalesOrderCode SalesOrderNumber CreditCardID RevisionNumber 2011-05-2011-06-2011-06-8 43659 SO4 SO43659 16281.0 31 12 1 rows × 28 columns In [207]: | data.iloc[0,6:8] Out[207]: Status OnlineOrderFlag Name: 0, dtype: object Rename column: data.rename(columns= {'ShipDate':'ShipingDate'},inplace= False).head(3) Out[210]: SL# SalesOrderID RevisionNumber OrderDate DueDate ShipingDate **Status** OnlineOrderFlag SalesOrderCode SalesOrderNumber **CreditCardIE** 2011-05-2011-06-8 43659 SO4 2011-06-07 SO43659 16281.0 31 2011-05-2011-06-8 43660 2011-06-07 0 SO4 SO43660 5618.0 2011-05-2011-06-8 2 3 43661 2011-06-07 0 SO4 SO43661 1346.0 31 3 rows × 28 columns In [247]: import pandas as pd import numpy as np %%HTML In [252]: <style type="text/css"> table.dataframe td, table.dataframe th { border-style: inset; </style> In [234]: filePath = 'G:/PRCTICE/Python/' inputFileName = 'BusinessEntity.xlsx' outputFileName = 'SalesOrderHeader_Analysis_Result.xlsx' inputFile = filePath + inputFileName outputFile = filePath + outputFileName data = pd.read_excel(inputFile) In [239]: data.head() In [240]: Out[240]: **BusinessEntityID** Title **FirstName MiddleName** LastName Suffix **PhoneNumber** PhoneNumberType **EmailPromotion** AddressT **EmailAddress** david22@adventure-1699 Mr. David R. Robinett NaN 238-555-0100 Home Ho works.com rebecca3@adventure-Cell Н 1700 Ms. 648-555-0100 Rebecca Α. Robinson NaN works.com dorothy3@adventure-Dorothy 2 2 1701 Ms. B. Robinson NaN 423-555-0100 Cell H works.com carolann0@adventure-3 1702 Ms. F. 439-555-0100 Cell 0 Н Carol Ann NaN Rockne works.com scott10@adventure-1703 Mr. Scott M. 989-555-0100 Cell 0 H Rodgers NaN works.com data.set_index(keys=['CountryRegionName','StateProvinceName','City']). Out[246]: BusinessEntityID Title FirstName MiddleName LastName Suffix **PhoneNumber** PhoneNumber^{*} CountryRegionName **StateProvinceName** City Germany Nordrhein-Solingen 1699 Mr. David R. Robinett NaN 238-555-0100 Westfalen **Australia** Seaford Victoria 1700 648-555-0100 Robinson NaN Ms. Rebecca Geelong 1701 Ms. Dorothy В. Robinson NaN 423-555-0100 **United Kingdom England** Lancaster F. 1702 Ms. Carol Ann Rockne 439-555-0100 NaN **East Brisbane Australia** Queensland 1703 Mr. Scott Μ. 989-555-0100 Rodgers NaN **United Kingdom** England Esher-1704 Mr. NaN Rodman NaN 899-555-0100 Jim Molesey **United States** California Concord Rothenberg 1705 326-555-0100 Mr. Eric NaN NaN **Australia New South Wales** St. Leonards 1706 Mr. Michael Rothkugel NaN 358-555-0100 **East Brisbane** Queensland 1707 Rovira Diez 786-555-0100 Mr. Pablo NaN NaN **Victoria** Seaford 1708 Ms. Linda R. Rousey NaN 369-555-0100 **United States** California **Beverly Hills** 1709 583-555-0100 Mr. Luke Roy NaN **Australia** Rockhampton Queensland 1710 Ms. Lisa K. 953-555-0100 Roy NaN France **Seine Saint Denis** Saint-Denis 1711 Mr. Michael NaN NaN 227-555-0100 Ruggiero **Australia Victoria** Seaford 1712 633-555-0100 Ms. Pearlie Rusek NaN **France Pantin Seine Saint Denis** 1713 Ms. Andrea NaN Rusko NaN 587-555-0100 Germany Hessen Hamburg 1714 856-555-0100 Mr. NaN Ruth Andy NaN **Australia** Victoria Warrnambool 1715 Ms. Justine Ryan NaN 498-555-0100 Germany Hessen Salzgitter 1716 Ms. N. Sabella NaN 529-555-0100 Deanna France Seine (Paris) **Paris** 1717 Sacksteder 114-555-0100 Mr. Lane NaN NaN Pas de Calais Boulogne-1718 Mr. Peter NaN Saddow NaN 132-555-0100 sur-Mer **Australia** Queensland **East Brisbane** 1719 Ms. 109-555-0100 Sharon NaN Salavaria NaN **New South Wales Coffs Harbour** 111-555-0100 1720 Mr. Irving W Schmidt NaN Germany Saarland Stuttgart 1721 313-555-0100 Mr. Raymond Sam NaN NaN **United Kingdom England** London 1722 NaN Mandar NaN Samant NaN 129-555-0100 Canada **British Columbia** Cliffside 1723 Mr. Sandstone NaN 114-555-0100 John **Australia** Lavender Bay **New South Wales** Madalena Sanchez 749-555-0100 1724 Ms. NaN **United States** California Colma 1725 Mr. T. Jr. 127-555-0100 Thomas Sanchez **France** Les Ulis Essonne 1726 Mr. 967-555-0100 Ken NaN Sánchez NaN Germany Nordrhein-Solingen 1727 Mr. Mikael NaN Sandberg NaN 854-555-0100 Westfalen Australia **New South Wales Milsons Point** 1728 Ms. R. Sandidge 113-555-0100 Mary NaN Nordrhein-1 (11) 500 555-Germany **Paderborn** 20749 NaN Crystal NaN Liu NaN Westfalen **United States** Oregon Lake Oswego 20750 NaN Isabella M 126-555-0181 Stewart NaN 1 (11) 500 555-**New South Wales Australia Newcastle** 20751 Yang NaN Crystal NaN Α 0112 **United States** California Colma Sanchez 20752 NaN Isabella D NaN 917-555-0121 El Cajon 20753 NaN M Morris 693-555-0146 Isabella NaN **Australia** 1 (11) 500 555-Queensland **Gold Coast** Huang 20754 NaN Crystal L NaN 0127 1 (11) 500 555-**United Kingdom Kirkby** England 20755 NaN W Wu NaN Crystal 0151 France **Seine Saint Denis Drancy** 1 (11) 500 555-20756 NaN Crystal NaN Lin NaN 0113 1 (11) 500 555-Australia East Brisbane Queensland 20757 NaN L NaN Carol Zhou 0141 **United States** Washington Seattle 20758 NaN Isabella NaN NaN 199-555-0140 Rogers Renton 20759 NaN Isabella NaN Reed NaN 754-555-0118 20759 Renton NaN Isabella NaN 754-555-0118 Reed NaN 20760 D 248-555-0151 Lynnwood NaN Isabella Cook NaN Canada **British Columbia** Haney 20761 NaN Morgan 118-555-0127 Isabella NaN 1 (11) 500 555-**Australia South Australia Findon** 20762 NaN Crystal NaN Zhao NaN 0161 1 (11) 500 555-Queensland Rockhampton 20763 NaN Isabella NaN Bell NaN 0157 1 (11) 500 555-Victoria Geelong 20764 NaN Crystal Α Lu NaN 0199 1 (11) 500 555-**New South Wales** Silverwater 20765 D NaN Ruiz NaN Francis 0190 1 (11) 500 555-North Ryde 20766 Ε NaN Crystal Xu NaN 0127 1 (11) 500 555-Germany Brandenburg Eilenburg 20767 NaN Ε Crystal Sun NaN 0119 1 (11) 500 555-France Moselle Metz 20768 NaN Isabella NaN Murphy NaN 0117 Germany 1 (11) 500 555-Hamburg Hamburg 20769 NaN Crystal С Zhu NaN 0180 1 (11) 500 555-Hessen Duesseldorf 20770 NaN NaN Crystal Gao **United States** California **Beverly Hills** 20771 NaN Bailey 808-555-0174 Isabella NaN NaN 1 (11) 500 555-**United Kingdom England** London 20772 NaN Crystal Liang NaN 0120 1 (11) 500 555-W. York 20773 NaN Crystal NaN Guo NaN 0171 California **United States Torrance** Richardson 910-555-0166 20774 NaN Isabella NaN Washington Bremerton 20775 NaN Crystal S He NaN 813-555-0148 Versailles **France** Yveline 1 (11) 500 555-20776 NaN Crystal NaN Zheng NaN 0171 1 (11) 500 555-**Australia New South Wales Darlinghurst** 20777 NaN Crystal NaN Hu NaN 0126 18508 rows × 15 columns In [269]: data_count = data.count() In [270]: data count.head() Out[270]: BusinessEntityID 18508 101 Title FirstName 18508 MiddleName 10666 LastName 18508 dtype: int64 In [271]: col_names = ["Column","Value"] my df = pd.DataFrame([data count],columns = col names) In [272]: my_df Out[272]: Value Column NaN NaN In []: