

In [363]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
```

Part 0: Reading the data

- Please find the data (Fytlyff_DS_Interview.csv) and read it as a Dataframe

In [364]:

```
1 data = pd.read_csv(r"C:\Users\Danish\Desktop\Data Science\fltrr\Fyt
2
3 print("Total Number of Rows: ",data.shape[0])
4 print("Total Number of Columns: ",data.shape[1])
5 data.head()
6 # A general view of the data read from csv file
```

Total Number of Rows: 2160

Total Number of Columns: 10

Out[364]:

	Year	Month	MobileWeb_or_Web	Type_of_Customers?	Where_Are_T
0	2019	Jan	Desktop_Website	Existing_Customer	
1	2019	Jan	Desktop_Website	Existing_Customer	
2	2019	Jan	Desktop_Website	Existing_Customer	
3	2019	Jan	Desktop_Website	Existing_Customer	
4	2019	Jan	Desktop_Website	Existing_Customer	

In [365]:

```
1 print("    ", "Name of columns/Features")
2 for ech in range(len(data.columns)):
3     print(ech+1, "    ", data.columns[ech])
4 print("-"*20, "DATA TYPES", "-"*20)
5 print("    ", "Data type of columns/Features")
6 for ech in range(len(data.columns)):
7     print(ech+1, "    ", data.dtypes[ech])
```

Name of columns/Features

```
1 Year
2 Month
3 MobileWeb_or_Web
4 Type_of_Customers?
5 Where_Are_They_comming_from?
6 Which_Place_in_India?
7 How_many_Landed_on_our_Page?
8 How_many_Landed_on_the_our_Page_and_clicked_on_a_but
ton?
9 How_many_Landed_on_the_our_Page_and_clicked_on_a_but
ton_and_started_filling_the_Form?
10 How_many_Landed_on_the_our_Page_and_clicked_on_a_bu
tton_and_started_filling_the_Form_and_Completed_and_submi
ted_the_form?
```

----- DATA TYPES -----

Data type of columns/Features

```
1 int64
2 object
3 object
4 object
5 object
6 object
7 float64
8 float64
9 int64
10 int64
```

Part 1: Data cleaning

Write a function called `data_cleaning()` which, when called, would perform the following activity:

- 1. Replaces the NA values with 0s in the data
- 2. In column 'B' replace Jan with 1, feb with 2, march with 3 and so on...
- 3. In column 'E' Replace "Came_From_Google" with "Google" and "Landed_on_the_page_Directly" with "Direct_traffic"

In [366]:

```
1 monthDict={'Jan':1, 'Feb':2, 'Mar':3, 'Apr':4, 'May':5, 'Jun':6,
2           'Jul':7, 'Aug':8, 'Sep':9, 'Oct':10, 'Nov':11, 'Dec':12}
3
4 # data_cleaning function
5 def data_cleaning(d):
6
7     # 1. Replaces the NA values with 0s in the data
8     d.fillna(0,inplace=True)
9
10    # 2. In column 'B' replace Jan with 1, feb with 2, march with 3
11    d['Month'] = monthDict.get(d['Month'],0)
12
13    # 3. In column 'E' Replace "Came_From_Google" with "Google" and
14    # "Landed_on_the_page_Directly" with "Direct_traffic"
15    if d["Where_Are_They_comming_from?"] == "Came_From_Google":
16        d["Where_Are_They_comming_from?"] = "Google"
17        #for "Came_From_Google": TO "Google"
18
19    if d["Where_Are_They_comming_from?"] == "Landed_on_the_page_Dir
20        d["Where_Are_They_comming_from?"] = "Direct_traffic"
21        # "Landed_on_the_page_Directly": TO "Direct_traffic"
22    return d
23
24 data=data.apply(
25     data_cleaning, # funciton call
26     axis=1,
27     result_type='expand')
28
29
30 # data.Month = data.Month.astype(str)
31 # data.Year = data.Year.astype(str)
```

Part 2: Descriptive statistics

Write a function called `descriptive_stats ()` which, when called, would perform the following activity:

- 1. Generates the summary statistics (Mean, Median, Quartile, standard deviation) of all the numerical columns
- 2. Produce a list of all the unique values & data types present in the non-numeric columns

In [367]:

```
1  stat_names = ["Mean", "Median", "Quartile", "standard deviation"]
2
3  def descriptive_stats(x):
4      column_labels=(x.columns)
5      # each column name is stored in column_labels variable
6      num_dict={}
7      # Numeric_columns variable will contain the name of numeric col
8      non_num_dict={}
9      # Non_Numeric_columns variable will contain the name of non num
10
11     # 1. Generates the summary statistics (Mean, Median, Quartile,
12     # of all the numerical columns
13
14     for label in range(len(column_labels)):
15         # This will Filter the Column Name with their respective va
16         # according to the Type of the data into the respective col
17         if np.issubdtype(x[column_labels[label]].dtype, np.number):
18             num_dict[column_labels[label]] = list(x[column_labels[l
19         else:
20             non_num_dict[column_labels[label]] = list(x[column_labe
21
22     df_dict = {}
23     list_keys = list(num_dict.keys())
24     list_value = list(num_dict.values())
25     df2 = pd.DataFrame(data=non_num_dict)
26
27     # these list are create to make a data frame that can be easily
28     # other wise output will not be clearly understood when printed
29     mn=[]
30     md=[]
31     qtl=[]
32     sd=[]
33
34     for ele in range(len(list_keys)):
35         mn.append(np.mean(list_value[ele]))
36         md.append(np.median(list_value[ele]))
37         qtl.append(np.quantile(list_value[ele],[0.25,0.75]))
38         sd.append(np.std(list_value[ele]))
39     stat_values=[mn,md,qtl,sd]
40
41     for u in range(len(stat_values)):
```

```

42         df_dict[stat_names[u]] = stat_values[u]
43
44     df = pd.DataFrame(data=df_dict,index=list_keys)
45
46     #PART 2
47     # 2. Produce a list of all the unique values &
48     # data types present in the non-numeric columns
49
50     key_list = list(non_num_dict.keys())
51     val_list = list(non_num_dict.values())
52
53     for ech_el in range(len(key_list)):
54         print(ech_el+1,"list of all the unique values in: ",
55               key_list[ech_el],"& Data type(dtype) of Column","\n",
56               df2[key_list[ech_el]].value_counts())
57         print("-"*100)
58     return df
59
60 df=descriptive_stats(data)
61
62 print("Mean", "Median", "Quartile", "standard deviation","informati
63 df.head()
64

```

1 list of all the unique values in: MobileWeb_or_Web & Data type(dtype) of Column

Desktop_Website	1080
Mobile_website	1080

Name: MobileWeb_or_Web, dtype: int64

2 list of all the unique values in: Type_of_Customers? & Data type(dtype) of Column

Existing_Customer	1080
New_Customer	1080

Name: Type_of_Customers?, dtype: int64

3 list of all the unique values in: Where_Are_They_comming_from? & Data type(dtype) of Column

Google	720
Direct_traffic	720
Unidentified_Sources	720

Name: Where_Are_They_comming_from?, dtype: int64

```

-----
4 list of all the unique values in: Which_Place_in_India? & Data type(dtype) of Column
  Bangalore      432
Chennai          432
Dehradun         432
Indore           432
Pune             432
Name: Which_Place_in_India?, dtype: int64
-----

```

Mean Median Quartile standard deviation information

Out[367]:

How_many_Lan

How_many_Landed_on_the_our_Page_and_cli

How_many_Landed_on_the_our_Page_and_clicked_on_a_button_and_starte



Some important points about Part 2: Descriptive statistics (descriptive_stats()) function

- Data passed to the descriptive_stats() should be pd.dataframe type (remaining not tested)
- function generate two out puts 1--(by print() of python) 2--(dataFrame)

- 1 print() give out of "Produce a list of all the unique values data types present in the non-numeric columns" of this part of the question
- 2 is DataFrame give Mean, Median, Quatile, Standard deviation of each column mention in the index (row wise) in data Frame
- Quartile values 25th percentile and 75th percentile are calculated in list
- This function is dependend up on numpy and pandas

Part 3: Prescriptive statistics

Can you write code and present the data which would help us answer (Text in “” are column names):

- 1. “Which_Place_in_India?” has the highest “How_many_Landed_on_the_our_Page?”
- 2. “How_many_Landed_on_the_our_Page_and_clicked_on_a_button_and_started_fi and_Completed_and_submitted_the_form?” divided by “How_many_Landed_on_our_Page?” is highest for “Which_Place_in_India?”

In [368]:

```
1 # 1. “Which_Place_in_India?” has the highest “How_many_Landed_on_th
2
3 result =data[["Which_Place_in_India?","How_many_Landed_on_our_Page?
4 print(result)
```

```
Which_Place_in_India?          Pune
How_many_Landed_on_our_Page?    11274131.0
dtype: object
```


In [369]:

```
1 # 2. "How_many_Landed_on_the_our_Page_and_clicked_on_a_button_and_st
2 # divided by
3 # "How_many_Landed_on_our_Page?"
4 # is highest for "Which_Place_in_India?"
5
6 grouped_data = (
7     data.groupby(
8         ["Which_Place_in_India?"])[
9         "How_many_Landed_on_the_our_Page_and_clicked_on_a_button_an
10     ].sum())/(data.groupby(
11         ["Which_Place_in_India?"])[ "How_many_Landed_on_our_Page?" ].sum(
12
13 print(" place is: ",grouped_data.idxmax(),"\n",
14       "value is: ",grouped_data.max(),)
```

```
place is:  Bangalore
value is:  0.19752416721431196
```

Part 4: Simple Machine learning questions

Write a function called `pred_future()` which, when called, would perform the following activity:

- 1. Predict
"How_many_Landed_on_the_our_Page_and_clicked_on_a_button_and_started_filli
_and_Completed_and_submited_the_form?" for the complete year of 2022
- 2. Generate the overall MAPE of your prediction for the year 2021.

In [370]:

```
1 print(data.shape)
2 # data.head(60)
3
4
5 df=data[['Year', 'Month', 'How_many_Landed_on_the_our_Page_and_clicks']]
6 df=df.rename(columns = {"How_many_Landed_on_the_our_Page_and_clicks": "How_many_Landed_on_the_our_Page_and_clicks"})
7 print(df.shape)
8 # df
```

(2160, 10)

(2160, 3)

In [371]:

```
1 df19 = df[df["Year"]==2019].reset_index(drop=True)
```

In [372]:

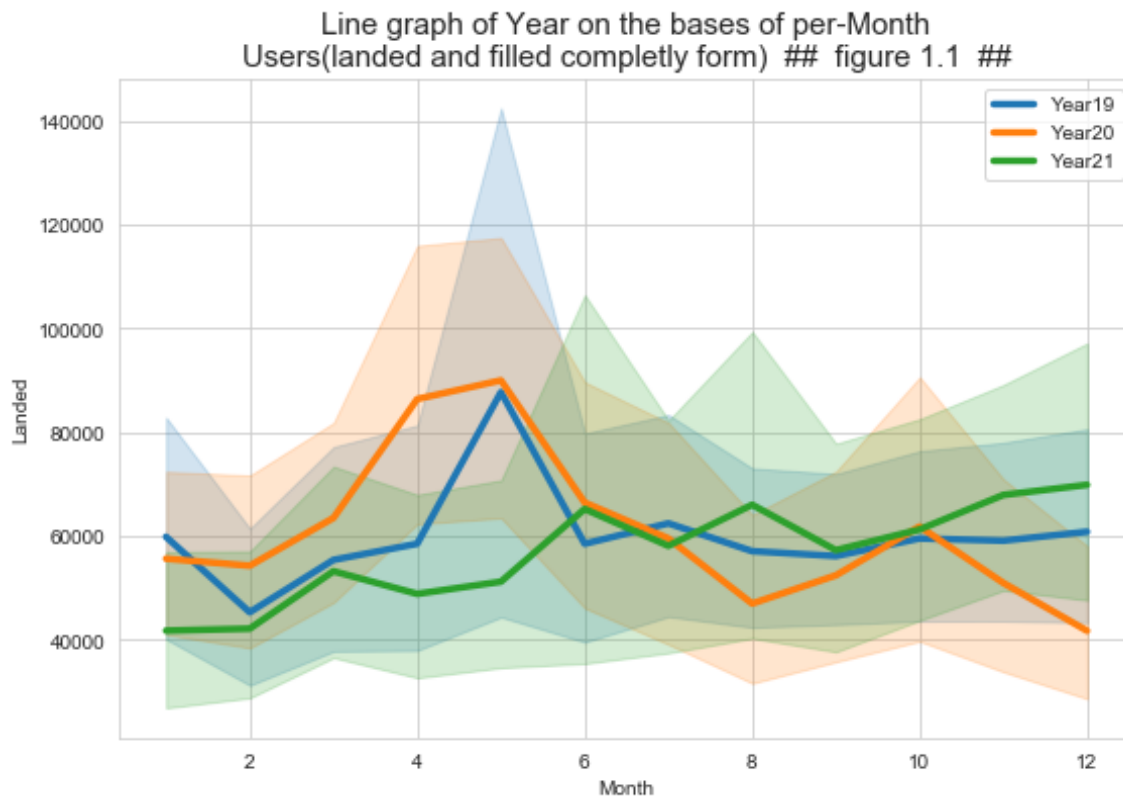
```
1 df20 = df[df["Year"]==2020].reset_index(drop=True)
2 # df20.head()
```

In [373]:

```
1 df21 = df[df["Year"]==2021].reset_index(drop=True)
2 # df21.head()
```

In [374]:

```
1 plt.figure(figsize=(9,6))
2 plt.title("""Line graph of Year on the bases of per-Month
3 Users(landed and filled completely form)  ##  figure 1.1  ##""",font
4 sns.lineplot(x='Month',y="Landed", label = 'Year19', data=df19, pa
5 sns.lineplot(x='Month',y="Landed", label = 'Year20', data=df20,pal
6 sns.lineplot(x='Month',y="Landed", label = 'Year21', data=df21,pal
7
8 plt.show()
9 # df20
```



In [375]:

```
1 def merge_date(x):
2     y=str(x["Year"])+str(x["Month"])
3     return y
```

In [376]:

```
1 df_19=df19.groupby(["Month", "Year"])["Landed"].sum()
2 df_19=df_19.reset_index()
3 df_19["Year_month"]=df_19.apply(merge_date,axis=1)
4 df_19.Year_month=df_19.Year_month .astype(np.datetime64())
```

In [377]:

```
1 df_20=df20.groupby(["Month", "Year"])["Landed"].sum()
2 df_20=df_20.reset_index()
3
4 df_20["Year_month"]=df_20.apply(merge_date,axis=1)
5 df_20.Year_month=df_20.Year_month .astype(np.datetime64())
6
```

In [378]:

```
1 df_21=df21.groupby(["Month", "Year"])["Landed"].sum()
2 df_21=df_21.reset_index()
3
4 df_21["Year_month"]=df_21.apply(merge_date,axis=1)
5 df_21.Year_month=df_21.Year_month .astype(np.datetime64())
6
```

In [379]:

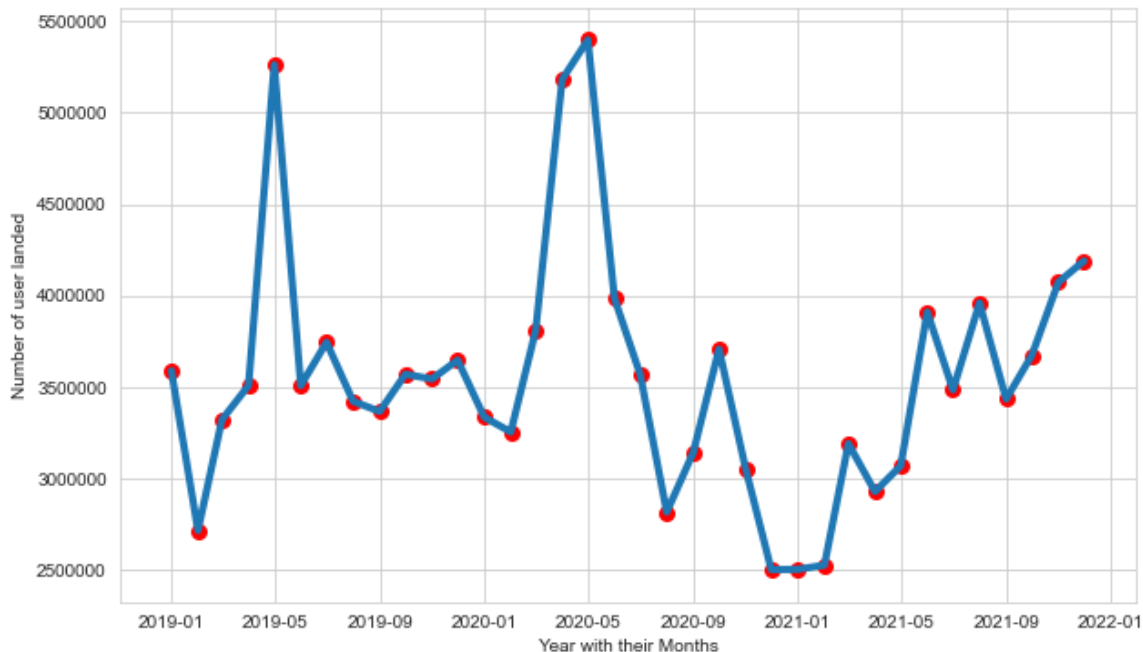
```
1 l=np.array((df_19["Landed"]+df_21["Landed"]+df_21["Landed"])/3, dtype
2 df_22=pd.DataFrame(data={"Year":2022, "Month":df_19["Month"], "Landed
3
4 df_22["Year_month"]=df_22.apply(merge_date,axis=1)
5 df_22.Year_month=df_22.Year_month .astype(np.datetime64())
6
7
8 test_y=df_22["Year_month"]
```

In [380]:

```
1 df_sum = pd.concat([df_19,df_20,df_21])
2 df_sum.columns
3 y=df_sum['Year_month']
4 x=df_sum.drop(columns=['Year_month',"Year","Month"])
```

In [381]:

```
1 plt.figure(figsize=(10,6))
2 sns.lineplot(x='Year_month',y="Landed",data=df_sum,linewidth=4)
3 plt.scatter(df_sum.Year_month,df_sum.Landed,linewidth=3,color="red")
4 plt.xlabel("Year with their Months")
5 plt.ylabel("Number of user landed")
6 plt.show()
7
```



In [382]:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.model_selection import train_test_split
```

In [383]:

```
1
2 X_train,x_test,Y_train,y_test= train_test_split(x,y,test_size=0.30,
3 print(X_train.shape,x_test.shape)
4 print(Y_train.shape,y_test.shape)
5 # y_test.head()
```

(25, 1) (11, 1)

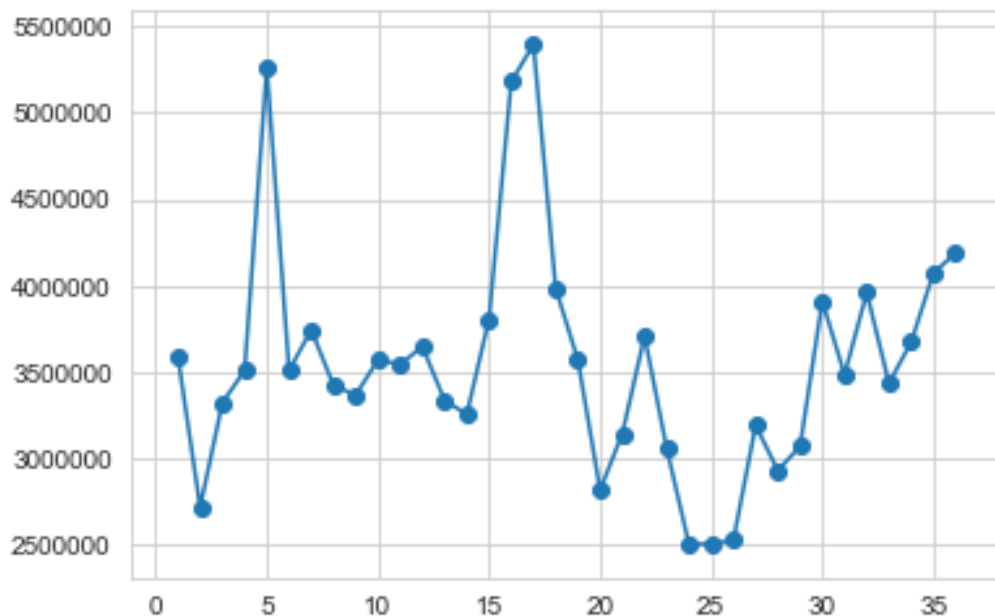
(25,) (11,)

In [384]:

```
1 Months=list(range(1,37))
2 landed = list(df_sum["Landed"])
3 dic = dict(zip(Months,landed))
4 sns.lineplot(x=dic.keys(),y=dic.values())
5 plt.scatter(x=dic.keys(),y=dic.values())
6 reg_data=pd.DataFrame(data={"Landed":dic.values(),"Year":dic.keys()})
7 reg_data.tail()
```

Out[384]:

	Landed	Year
31	3959670	32
32	3434896	33
33	3673751	34
34	4072884	35
35	4189930	36



In [385]:

```
1
2 y=reg_data['Year']
3 x=reg_data.drop(columns=['Year'])
4
5 x_train,x_test,y_train,y_test= train_test_split(x,y,test_size=0.30,
6
```

In [386]:

```
1 from sklearn.linear_model import LinearRegression
```

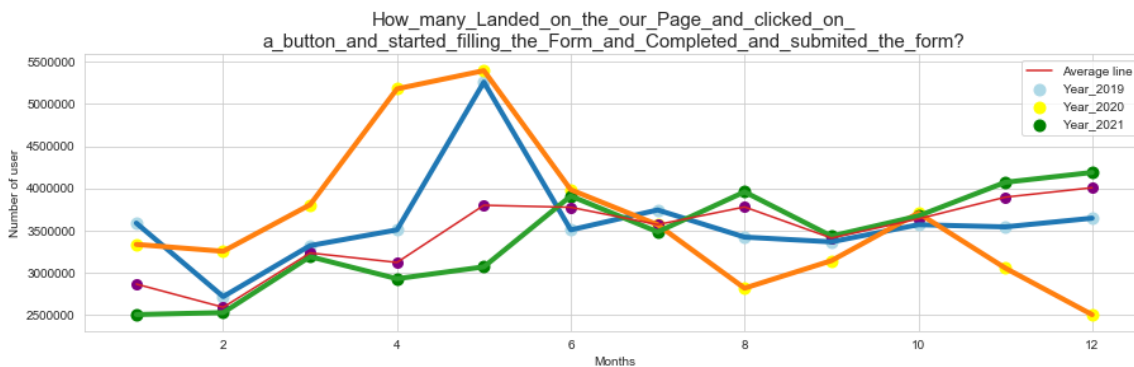
In [112]:

```
1 y_pred = reg.predict(x_test)
2 print(y_pred)
3
```

```
[3367556.33748773 3433658.67321643 3457695.88620868 33735
65.64073579
3421640.0667203 3463705.18945675 3493751.70569707 35418
26.13168158
3487742.402449 3379574.94398386 3391593.55047998]
```


In [234]:

```
1 plt.figure(figsize=(15,4))
2 sns.set_style('whitegrid')
3 plt.title("""How_many_Landed_on_the_our_Page_and_clicked_on_
4 a_button_and_started_filling_the_Form_and_Completed_and_submitted_th
5
6 sns.lineplot(x='Month',y="Landed",data=df_19,linewidth=4)
7 plt.scatter(df_19.Month,df_19.Landed,label="Year_2019",linewidth=4,
8 sns.lineplot(x='Month',y="Landed",data=df_20,linewidth=4)
9 plt.scatter(df_20.Month,df_20.Landed,label="Year_2020",linewidth=4,
10 sns.lineplot(x='Month',y="Landed",data=df_21,linewidth=4)
11 plt.scatter(df_21.Month,df_21.Landed,label="Year_2021",linewidth=4,
12 sns.lineplot(x='Month',y="Landed", data=df_22,label="Average line")
13 plt.scatter(df_22.Month,df_22.Landed,label="Year_2021",linewidth=3,
14
15
16
17 plt.xlabel("Months")
18 plt.ylabel("Number of user")
19 plt.show()
20
```

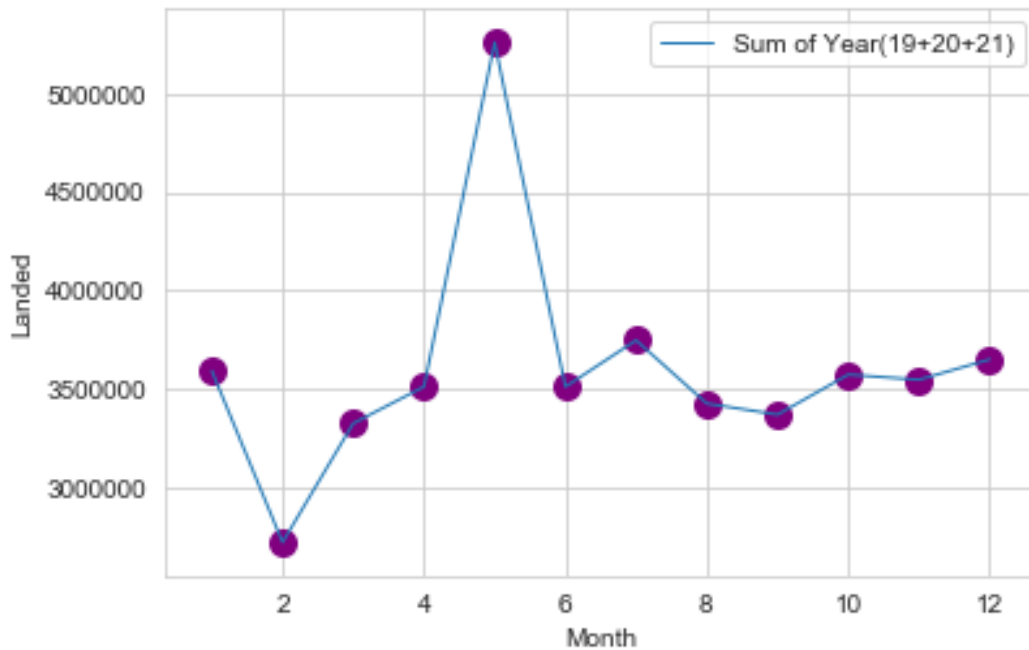


In [387]:

```
1 sum_df=pd.DataFrame()
2 sum_df["Landed"] =df19.groupby(["Month"])[ "Landed"].sum().reset_ind
3 sum_df["Month"]=df_19["Month"]
```

In [389]:

```
1 sns.lineplot(x='Month',y=sum_df["Landed"], data=sum_df,linewidth=1,  
2 plt.scatter(sum_df.Month,sum_df.Landed,linewidth=5,color="purple",)  
3 plt.show()  
4
```



In [390]:

```
1 lm=LinearRegression()  
2 lm.fit(df_22[["Month"]],df_22.Landed)
```

Out[390]:

LinearRegression()

In [391]:

```
1 x_test
2 t_22 =pd.DataFrame(data=list(range(13,25)))
3 # t_22
4 # print(type(t_22))
```

In [392]:

```
1 # y_Pred =lm.predict(t_22)
2 # y_Pred=np.array(y_Pred,dtype=np.int64)
3 # print(y_Pred)
```

In [393]:

```
1 def Pred_future(x):
2 #     y_Pred=np.array(y_Pred,dtype=np.int64)
3     v= dict(zip(t_22[0],y_Pred))
4     print("For the Next Year {} User that Land and completly fill f
5     print(v)
6     Pred_future(2022)
```

For the Next Year 2022 User that Land and completly fill
form will be probably:
{13: 4119143, 14: 4218205, 15: 4317267, 16: 4416329, 17:
4515391, 18: 4614453, 19: 4713515, 20: 4812577, 21: 49116
39, 22: 5010701, 23: 5109763, 24: 5208825}

There is lot of Analysis that i did. Now i will discuss some of the points :

- 1. I trie different method out of them i choose the Average method.
- 2. In Average method i pick value from the all the year
and sum them up with repect their corresponding year And devide the by 3
because year are given 3.
- 3. Then i trained the linear regression mode on top of it which gave me the abve
line(green).
-
-

2. Generate the overall MAPE of your prediction for the year 2021.

In [394]:

```
1 forecast=y_Pred
2 actual=list(df_21["Landed"])
3 APE = []
4 for day in range(len(actual)):
5     per_err = (actual[day] - forecast[day]) / actual[day]
6     per_err = abs(per_err)
7     APE.append(per_err)
8 MAPE = sum(APE)/len(APE)
9
10 print(f'''
11 MAPE : { round(MAPE, 2) }
12 MAPE % : { round(MAPE*100, 2) } %
13 ''')
14
```

MAPE : 0.39

MAPE % : 39.03 %

Part 5: Visualization

- Please write a code to display :
 1. A line graph for
“How_many_Landed_on_the_our_Page_and_clicked_on_a_button?” for the different “Which_Place_in_India?” over the months of the year 2019 & 2020.
 2. A line graph of the actual and projected number of
“How_many_Landed_on_the_our_Page_and_clicked_on_a_button_and_started_filling_and_Completed_and_submited_the_form?” for the months of the year 2021 & 2022

1. A line graph for

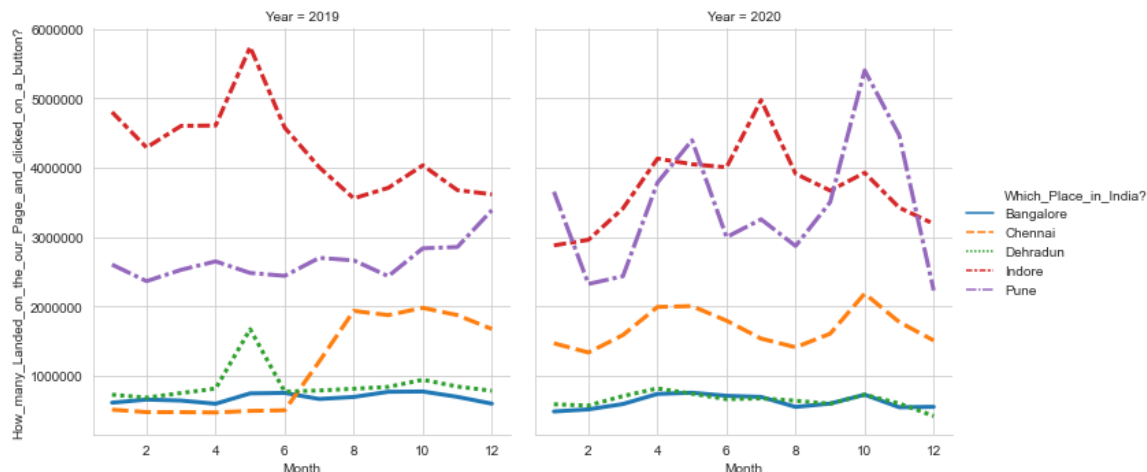
“How_many_Landed_on_the_our_Page_and_clicked_on_a_button” for the different “Which_Place_in_India?” over the months of the year 2019 & 2020.

In [395]:

```

1 q1=data[["Year","Month","How_many_Landed_on_the_our_Page_and_clicked_on_a_button"]]
2 q1=q1.groupby(["Year","Month","Which_Place_in_India?"])[["How_many_Landed_on_the_our_Page_and_clicked_on_a_button"]]
3 q1=q1.reset_index()
4 q1=q1[q1["Year"]!=2021].reset_index(drop=True)
5 sns.relplot(
6     data=q1, x="Month", y="How_many_Landed_on_the_our_Page_and_clicked_on_a_button",
7     col="Year", hue="Which_Place_in_India?",kind="line", style="Which_Place_in_India?"
8 )
9 plt.show()

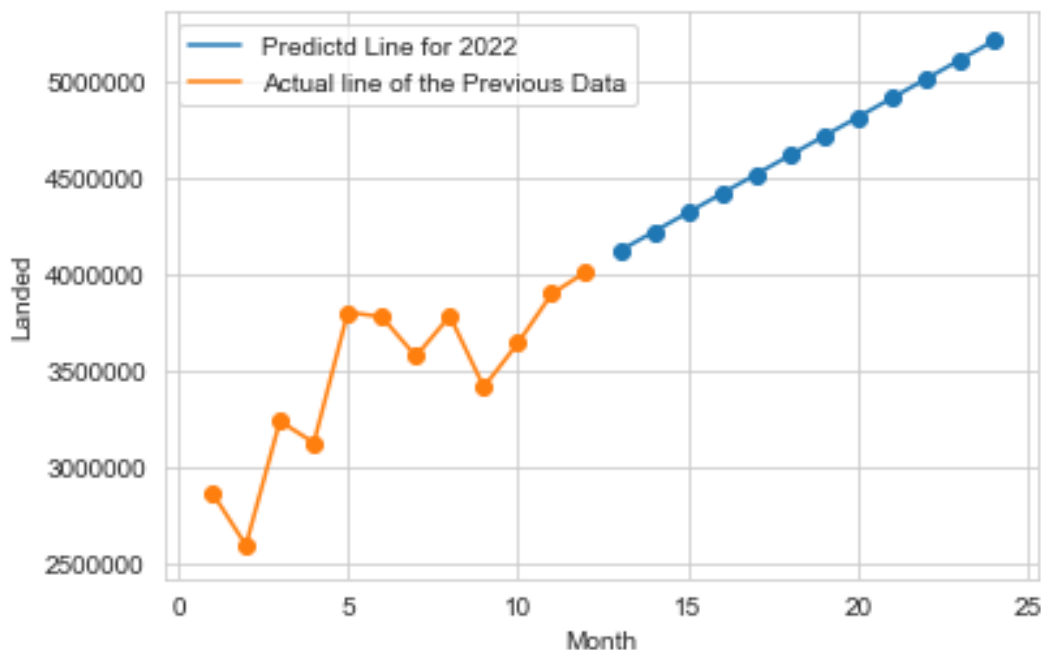
```



2. A line graph of the actual and projected number of “How_many_Landed_on_the_our_Page_and_clicked_on_a_button_and_Completed_and_submited_the_form?” for the months of the year 2021 & 2022

In [396]:

```
1 l=list(dic.keys())[:12]
2 val=list(df_22['Landed'])
3 # val.extend(y_Pred)
4 # print()
5 sns.lineplot(x=t_22[0],y=y_Pred,label="Predictd Line for 2022")
6 plt.scatter(x=t_22[0],y=y_Pred)
7
8 sns.lineplot(x=df_22["Month"],y=df_22["Landed"],label="Actual line")
9 plt.scatter(x=l,y=val)
10 plt.show()
```



Part 6: About the Previous projects

Please describe any interesting project you did in the Data Science domain in more than 250 words. Attach Github links if possible.

Description of the Problem

[Source Kaggle \(https://www.kaggle.com/c/quora-question-pairs\)](https://www.kaggle.com/c/quora-question-pairs)

Where else but Quora can a physicist help a chef with a math problem and get cooking tips in return? Quora is a place to gain and share knowledge—about anything. It's a platform to ask questions and connect with people who contribute unique insights and quality answers. This empowers people to learn from each other and to better understand the world.

Over 100 million people visit Quora every month, so it's no surprise that many people ask similarly worded questions. Multiple questions with the same intent can cause seekers to spend more time finding the best answer to their question, and make writers feel they need to answer multiple versions of the same question. Quora values canonical questions because they provide a better experience to active seekers and writers, and offer more value to both of these groups in the long term.

Currently, Quora uses a Random Forest model to identify duplicate questions. In this competition, Kagglers are challenged to tackle this natural language processing problem by applying advanced techniques to classify whether question pairs are duplicates or not. Doing so will make it easier to find high quality answers to questions resulting in an improved experience for Quora writers, seekers, and readers.

i didn't complete it yet as i run many algorithms on it and done some NLP on it using Glove it will be available soon on git hub

Part 7 : Time management

• Can you please share your thoughts, in less than 120 words, on “If you get selected, how will you manage your time for this full-time internship opportunity”

As i am learning yet it will be good for me if i got an opportunity. I will try my best to get the best of best benefit out of this intern in terms of my skill advancement and hope fully you would also get benefit from my work.