	import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns import datetime import gcsfs import gsfs import sklearn.metrics as metrics import datetime import datetime import posts import datetime import posts import datetime import datet
lç	
li	<pre>mport the best model from GCP credentials = { "type": "service_account", "project_id": "phonic-monolith-345108", "private_key_id": "09e8b01c72c4bf6f11b35998df3d41d727fb33c8", "private_key_id": "09e8b01c72c4bf6f11b35998df3d41d727fb33c8", "private_key": "BEGIN PRIVATE KEY\nMIIEvgIBADANBgkqhkiG9w0BAQEFAASCBKgwggSkAgEAAoIBAQDJyfcTdIRQuEo2\nIqWIx4U9Ari+RbF/HrZz+HBYVwL9+mpIvg0ySYlkcSdBljhx8dC4vP2YbKhbdpNp\ndZEwRy7vp35wsN2WWHx+IAeiMg0FyDSXnl3Pcl4Ic "client_email": "danielevko@phonic-monolith-345108.iam.gserviceaccount.com",</pre>
	<pre>"client_id": "112729790674474565098", "auth_uri": "https://accounts.google.com/o/oauth2/auth", "token_uri": "https://oauth2.googleapis.com/token", "auth_provider_x509_cert_url": "https://www.googleapis.com/oauth2/v1/certs", "client_x509_cert_url": "https://www.googleapis.com/robot/v1/metadata/x509/danielevko%40phonic-monolith-345108.iam.gserviceaccount.com" } gcp_json_credentials_dict = json.loads(json.dumps(credentials)) creds = service_account.Credentials_from_service_account_info(gcp_json_credentials_dict) client = storage.Client(project=gcp_json_credentials_dict['project_id'], credentials=creds)</pre>
:	<pre>bucket = client.bucket('final_project_leads') blob_model='' for blob in client.list_blobs('final_project_leads', prefix='danielev'): if 'best_model' in str(blob): print(blob) blob_model=blob</pre>
:	break cBlob: final_project_leads, danielev/best_model_random_forest, 1653483147209057> pickle_data = blob_model.download_as_string() model = pickle.loads(pickle_data)
	<pre>mport Data from GCP def import_from_gcp(file_name, bucket_name): PROJECT_NAME = 'final-project-lead-me' URL = "gs://" # Creating a pythonic file-system interface to Google Cloud Storage. fs = gcsfs.GCSFileSystem(project=PROJECT_NAME)</pre>
	<pre>for i in fs.ls(bucket_name): if file_name in i: print(i) return pd.read_csv(URL + i, storage_options={"token": credentials}, encoding="ISO-8859-1") return print("File not found") df_companies = import_from_gcp("Companies Data", "final_project_leads") df_cars = import_from_gcp("vehicles4", "final_project_leads")</pre>
f U	final_project_leads/Companies Data.csv final_project_leads/vehicles4.csv Upload a new leads file from the client df_leads = import_from_gcp("sales_leads_D0_NOT_DELETE_PART3", "final_project_leads/danielev") final_project_leads/danielev/sales_leads_D0_NOT_DELETE_PART3.csv df_leads.head()
11 22 33 44	id_lead id_lead id_lead id_lead_in_lear_in_id_lead_in_id_lear_in_id_lead_in_id_lear_in_id_lead_in_id_lear_in_id_lead_in_id_lear_in_id_lead_in_id_lear_id_lear_in_id_lear_in_id_lear_in_id_lear_in_id_lear_in_id_lear_i
C	Checking for duplication values in df_leads #checking duplicates id leads if (sum(df_leads.duplicated(subset = 'id_lead')) == 0): print("No Duplication in id_lead") else: #Number of duplications in id_lead dups = len(df_leads['id_lead'])-len(df_leads['id_lead'])-drop_duplicates()) print("{} Duplication in id_lead", format(dups)) df_leads.drop_duplicates(subset='id_lead'), keep='first', inplace=True) #checking after delete duplicates id_leads sum(df_leads.duplicated(subset = 'id_lead')) == 0 # No duplicate values
	<pre>#checking duplicates id's if (sum(df_leads.duplicated(subset = 'id')) == 0): print("No Duplication in id") else: #Number of duplications in id_lead dups = len(df_leads['id'])-len(df_leads['id'].drop_duplicates()) print("{} Duplication in id".format(dups)) df_leads.drop_duplicates(subset=['id'], keep='first', inplace=True) #checking after delete duplicates id leads sum(df_leads.duplicated(subset = 'id')) == 0</pre>
C	Description in id Calculate Age def calculate_age(born): today = datetime.date.today() return today.year - born df_leads['age'] = df_leads['year_of_birth'].apply(calculate_age)
	Data Cleaning df_companies=df_companies[df_companies['Market Cap']!='-'] df_companies=df_companies[df_companies['Market Cap'].notna()] df_companies=df_companies[df_companies['company_name'].map(lambda x: x.isascii())] df_companies.shape (947, 19)
(C	<pre>df_companies=df_companies.sort_values('profit', ascending=False) print(df_companies.shape) (947, 19) Cleaning the cars Dataset md("Cleaning cars whose year is than %i"%(df_leads['car_year'].min()))</pre>
	Cleaning cars whose year is than 1998 df_cars=df_cars[df_cars['year']>=df_leads['car_year'].min()] df_cars.shape (341327, 7)
Т	<pre>if len(df_leads[df_leads['car_year']<1997])==0: print(True) else: False True Cleaning cars that do not include the maunfactuer from df_leads</pre>
	manufacturers_to_drop=np.unique(df_cars['car'])[~np.in1d(np.unique(df_cars['car']), np.unique(df_leads['car_type']))].tolist() manufacturers_to_drop ['aston-martin', 'jaguar', 'land rover', 'tesla'] df_cars=df_cars[df_cars.car.isin(manufacturers_to_drop) == False] Cleaning cars that do not include the same model from df_leads
a	<pre>models_to_drop=np.asarray(list(set(df_cars['model'])))[~np.in1d(np.asarray(list(set(df_cars['model']))), np.unique(df_leads['car_model']))] models_to_drop array(['nan',</pre>
(df_cars=df_cars[df_cars.model.isin(models_to_drop) == False] df_cars.shape (95018, 7) df_leads: Data Transformation
Irr	this part we are going to match prices of cars in our df_leads to df_cars Create new df that is group by car_type and year df_mean_car_by_model_and_year=df_cars.groupby(['car', 'model', 'year']).mean() df_mean_car_by_model_and_year['price']=df_mean_car_by_model_and_year['price'].map(lambda x: round(x,2)) df_mean_car_by_model_and_year
C	2021 7.306675e+09 61999.00 901 rows × 2 columns calculate the average car price based on the type and the year of the car, if there is no year information, calculate the average based on the type of car temp_prices=[] for index, row in df_leads.iterrows(): #If the car_type is in df_cars it will enter
	<pre>try: #Creating a new Dataframe which grouped by the car_type new_df=df_mean_car_by_model_and_year.loc[df_mean_car_by_model_and_year.index].loc[row['car_type']].reset_index() #If the car_model is in df_cars, will select only the prices of the model if row['car_model'] in new_df['model']!=row['car_model']] #if there is info of model and year, will get the price of the model and year if row['car_year'] in new_df['year'].values: temp=new_df[(model']==row['car_model']) & (new_df['year']==row['car_year'])]['price'].values[0] temp=new_df[(new_df['model']==row['car_model']) & (new_df['year']==row['car_year'])]['price'].values[0] else: #if not, will give the minimum year price. temp_prices.append(new_df[new_df['price']==min(new_df['price'])]['price'].values[0]) else: #if there is no info of car_model, will give the car_price by year new_df=new_df.groupby(['year']).mean() if row['car_year'] in new_df.index: temp_prices.append(new_df.index) temp_prices.append(new_df.loc[row['car_year']]['price']) else: temp_prices.append(new_df.loc[min(new_df.index)]['price'])</pre>
	<pre>except: #If the car type is not located in df_cars it will calculate price by the car year new_df=df_cars.groupby(['year']).mean() if row['car_year'] in new_df.index: temp_prices.append(new_df.loc[row['car_year']]['price']) else: #if the year is not located it will calculate by the minimum year temp_prices.append(new_df.loc[min(new_df.index)]['price'])</pre>
C	df_leads['car_price']=temp_prices Convert the dates data to date_type in df_leads Convert to datetime to proper datatype df_leads['rental_period']=pd.to_datetime(df_leads['rental_period'], format='%d/%m/%Y')
0	df_leads['creation_date']=pd.to_datetime(df_leads['creation_date'], format='%d/%m/%Y') df_leads[['rental_period', 'creation_date']].head() rental_period creation_date 0 2018-06-22 2018-05-13 1 2019-12-03 2019-10-04 2 2020-07-30 2019-04-22
	3 2018-08-15 2018-02-24 4 2020-03-04 2019-12-04 Create new column based on creation_date and rental_period that called the diff of them, called 'desirable_rental_days' df_leads['desirable_rental_days']=df_leads['rental_period']-df_leads['creation_date'] df_leads['desirable_rental_days']=df_leads['desirable_rental_days'].dt.days
Т	<pre>ferify that the users filled out the rental period correctly if len(df_leads[df_leads['desirable_rental_days']<0])>=1: print(False) else: print(True) True</pre>
VE	Data Merge between df_leads and df_companies erify all the data in columns Market Cap and Profit is numeric try: df_companies['Market Cap']=pd.to_numeric(df_companies['Market Cap']) df_companies['profit']=pd.to_numeric(df_companies['profit']) print("Succeed") except: print("Error occured in converting Profit and Market Cap")
С	Succeed Create a copy of df_companies df_company_to_merge=df_companies.copy() prop unwanted columns in df_company_to_merge
M	df_company_to_merge.drop(['id_company', 'num. of employees', 'profitable', 'rank_change', 'rank', 'city', 'state', 'newcomer', 'ceo_founder', 'ceo_woman', 'prev_rank', 'CEO', 'Website', 'Ticker', 'sector'], inplace=True, axis=1) Make a copy of df_leads before merge df_leads_for_analysis=df_leads.copy() df_leads_for_analysis.sort_values(by="id_lead", inplace=True) Merge df_company_to_merge and df_leads
S	df_leads_for_analysis=df_leads.merge(df_company_to_merge, on='company_name', sort=False) fort id_lead Dataframe by id (To return to the original order) df_leads_for_analysis.sort_values(by="id_lead", inplace=True) This is the Dataset the client will get to preform sales
	<pre>df_leads_app = df_leads_for_analysis.copy() Drop unwanted columns df_leads_for_analysis.drop(['id', 'id_lead', 'first_name', 'last_name', 'email', 'year_of_birth', 'country', 'address', 'creation_date', 'rental_period', 'car_type', 'company_name', 'car_model'], axis=1, inplace=True) df_leads_for_analysis.head()</pre>
1	is_buisness gender platform department creation_time car_year desirable_rental_days revenue profit Market Cap 0 False Female Facebook Research and development 20 20 61 4749.50 40 6160.1 206.8 3542.0 4 False Male Google Research and development 15 2015 70 9360.78 66 8358.9 223.4 4294.0 18 True Male Facebook Human resources 10 2017 34 2179 450 4587.0 3434.0 1274139.0 20 False Male Website Engineering 4 2010 84 8577.74 172 9650.0 -62.4 576.9 21 False Male Google Engineering 1 2012 50 8604.47 91 8342.0 562.0 3217.2
C	Converting Categorical Features Convert the 'is_buisness' boolean column to String type df_leads_for_analysis['is_buisness'] = df_leads_for_analysis['is_buisness'].map({True: 'True', False: 'False'}) # Replace boolean to string
	Create a list of categorical columns categorical_columns = [i for i in range(len(np.array(df_leads_for_analysis.dtypes))) if np.array(df_leads_for_analysis.dtypes)[i] not in [np.dtype('float64'), np.dtype('int64')]] df_dummies=pd.get_dummies(df_leads_for_analysis[df_leads_for_analysis.columns[categorical_columns]], drop_first=True) Drop the Categorical columns from df_leads_for_analysis
C	df_leads_for_analysis.drop(df_leads_for_analysis.columns[categorical_columns], axis=1, inplace=True) Concat all the df_dummies df_leads_for_analysis = pd.concat([df_leads_for_analysis, df_dummies,], axis=1)
R	<pre>Make predictions on df_leads_for_analysis predictions = model.predict(df_leads_for_analysis) desults of predictions unique, counts = np.unique(predictions, return_counts=True) print (np.asarray((unique, counts)).T)</pre>
	[[0 716] [1 284]] df_leads_app['is_sold'] = predictions Analyizing the Model results
	<pre>sns.set_style('whitegrid') sns_plot_female_vs_male_buisness=sns.countplot(x='gender',hue='is_sold',data=df_leads_app,palette='RdBu_r') sns_plot_female_vs_male_buisness.figure.savefig("Male vs. Female is_Sold.jpeg") for p in sns_plot_female_vs_male_buisness.patches:</pre>
	400 ssold 0 1 1 1 1 1 1 1 1 1
	Female Male Graph of age sns_plot_age_vs_not_sold=df_leads_app[df_leads_app['is_sold']==0]['age'].hist(bins=30,color='red',alpha=0.5) sns_plot_age_vs_sold=df_leads_app[df_leads_app['is_sold']==1]['age'].hist(bins=30,color='green',alpha=0.5) sns_plot_age_vs_sold=df_leads_app[df_leads_app['is_sold']==1]['age'].hist(bins=30,color='green',alpha=0.5)
5 4	sns_plot_age_vs_sold=df_leads_app[df_leads_app['is_sold']==1]['age'].hist(bins=30,color='green',alpha=0.5) plt.legend(labels=["is_sold: False","is_sold: True"]) <matplotlib.legend.legend 0x1ea3259b070="" at=""> 50 40 40 40 40 40 40 40 40 40</matplotlib.legend.legend>
2	araph of creation Time & is_sold True \ False
<	sns_plot_creation_time=df_leads_app[df_leads_app['is_sold']==False]['creation_time'].hist(bins=30,color='red',alpha=0.5) sns_plot_creation_time=df_leads_app[df_leads_app['is_sold']==True]['creation_time'].hist(bins=30,color='green',alpha=0.5) plt.legend(labels=["is_sold: False","is_sold: True"]) cmatplotlib.legend.Legend at 0x1ea335ae8e0> definition of creation time & is_sold: True
4 3	is_sold: True is_sold: True is_sold: True is_sold: True
G	craph of Platform & is_sold True \ False sns.set_style('whitegrid') plt.figure(figsize = (10,8)) ax = sns.countplot(x='department', data=df_leads_app, palette='RdBu_r', hue='is_sold')
	ax.set_xticklabels(ax.get_xticklabels(),rotation = 30) for p in ax.patches: ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01)) plt.show() 160 161 162.0 162.0 163.0 164 165.0 165.0 166.0 167 168.0 16
and the second	80
Ρ	Regressed and derestropment, Human trescurses Engineering and values Controller service Administration department department
	<pre>sns.set_style('whitegrid') plt.figure(figsize = (10,8)) ax = sns.countplot(x='platform', data=df_leads_app, hue='is_sold') ax.set_xticklabels(ax.get_xticklabels(), rotation = 30) for p in ax.patches:</pre>
	187.0 175 150 133.0 131.0
	75 - 42.0
	result_data=df_leads_app.copy() result_data['Result'] = result_data['is_sold'].map({0:'No', 1:'Yes'}) # Order the cluster result_data['Result'] = result_data['is_sold'].astype('category')
	<pre>result_data.rename(columns = {'is_sold':'Total'}, inplace = True) df_groupby_segmant=result_data.groupby('Result').agg(</pre>
0	Nesult Total Septiment
U	Transform the DF to dask.dataframe and upload to the Cloud fort the file by is_sold df_leads_app.sort_values(by="is_sold", inplace=True, ascending=False) upload the file with today's date def_today_date(i): today = datetime.date.today()