<pre># need import import import # need from sk from sk</pre>	<pre>to library the package pandas as pd numpy as np xgboost as xgb to import the necessary dataset klearn.model_selection import train_test_split klearn.metrics import balanced_accuracy_score, roc_auc_score, make_scorer</pre>
from sk from sk from sk from sk #load t	<pre>klearn.model_selection import GridSearchCV klearn.metrics import confusion_matrix klearn.metrics import plot_confusion_matrix klearn.multioutput import MultiOutputRegressor the data to see the dataframe looks like = pd.read_csv(r'german.csv') man.head()</pre>
<pre># becau # use i german. b= germ print(b)</pre>	use the purpose question is a interview question, so we did not it as a prediction in the model, we will remove it .drop(['Purpose'],
<pre># usefu german. c= germ print(c # check</pre>	sue numebr id is different for each individual, thus it is not that ul to use for classification .drop(['number'], axis=1, inplace= True) man.head() c) king for missing data and identify to see what to do about it man.dtypes
<pre># sense e= germ print(e</pre>	to find the unique intepretation of the varibale and see if they make e to exists man['Sex'].unique() e) man['Housing'].unique()
g= germ print(g # when # need # this german.	<pre>man['Saving accounts'].unique() g) you do this , it will allow you see how the catergories are, and we to replace all the white spaces in all the columns with underscores. is because we can draw a nicer picture of XGBoost trees later .replace(' ', '_', regex= True, inplace = True)</pre>
<pre>h= germ print(h i= germ print(i ## spli # the f</pre>	man.head() n) man['Saving accounts'].unique()
#2. # we wi #classi # in th #return #this p # thus,	the cloumn of data that we want to predict ill use capital X to define the columns of data that we use to make ifications and small y to define the thing we want to predict his case we want to predict whether some one is high risk or low risk to not the credit that someone take in the bank according to the attribute that person have in the factors we want to predict the Default value here
A= X.he print(A y= germ B= y.he print(B ## Form	man['Default'].copy() ead()
<pre># categ C= X.dt print (</pre>	<pre>types (C) get_dummies(X, columns=['Sex']).head()</pre>
E= X_en	eed to varify that y only contain 0 and 1 unique function
print (## Buil # we ha # builc # we di	(F) Id a basic Preliinary XGboost Model aev successfully split the model into traiing and testing sets, we will d a model now. iscover that this data is imbalended by dividing the numebr of people consider not return, where y=1, divide by the number of people in the
<pre>print(j # we s #money # using # who c</pre>	see only 30% of people consider high risk that they did not return the y. thus, when we soikt the data into training and testing, we will slipt g stratificaion in order to maintain the same percentage of people do not return the borrow money in both the traiing set and testing set.
<pre>k= sum(print(k</pre>	(y_test)/len(y_test)
<pre># and y # retur # we ca clf_xgb</pre>	<pre>k and 1 =30%, this shows that stratify works as expected that both y_test y_train have the same number of high risk people profile that does not rm money. an then build the preliminary model b= xgb.XGBClassifier(objective = 'binary:logistic', missing = 1,</pre>
# we wi # the t	y_train, verbose=True, early_stopping_rounds=10, eval_metric='aucpr', eval_set=[(X_test, y_test)]) ee the best iteration is 22 is 0.63548 ill then see how the performance on the testing dataset by running testing dateset down the model and drawing the confusion matrix
<pre>print(t temp = print(t</pre>	<pre>f_xgb.predict(X_train) t1.shape) clf_xgb.predict(X_test) temp.shape) onfusion_matrix(clf_xgb,</pre>
#has lo #who ar	y_test, values_format='d', display_labels=["low risk","high risk"]) an see in this confusion matrix, we see that of 175 epople that ower risk, 149 (85.14%) is correctly classified. and of the 75 people re high risk, 42 (56%) is correctly classified. so this means this xgboost l was not that awesome.this may becasue the imbalanced dataset in the
<pre># begin # bank, # try t ## then</pre>	ning. but becasue high risk customer profile cost a lot more money to the , we want to capture more of people who are gih risk. thus we want to to improve the prediction using Cross Vlidation to optimize the parameters. In we want use AUC for evaluation, because grivdDearchCV takes long time the optimize parameters in a short period of time.
'le 'ga 're 'sc	<pre>grid= { ax_depth':[3,4,5], earning_rate':[0.1,0.01,0.05], amma':[0,0.25,1.0], eg_lambda':[0,1.0,10.0], cale_pos_weight':[1,3,5]</pre> = xgb.XGBRegressor(seed = 20)
#clf.fi #print(<pre>GridSearchCV(estimator=xgbr,</pre>
# # # # # # #	<pre>gridsearch(v) estimator=xgb.XGBClassifier(), param_grid={ 'max_depth':[3,4,5], 'learning_rate':[0.1,0.01,0.05], 'gamma':[0,0.25,1.0], 'reg_lambda':[0,1.0,10.0], 'scale_pos_weight':[1,3,5], }, cv=3, scoring='neg_mean_squared_error', verbose=0, n_jobs=-1)</pre>
<pre>#grid_r #best_p optimal est</pre>	<pre>result = MultiOutputRegressor(gsc).fit(X_train, y_train) params = grid_result.estimators_[0].best_params_ l_params = GridSearchCV(timator=xgb.XGBClassifier(objective='binary:logistic',</pre>
scc ver n_j cv=) #from s	<pre>colsample_bytree=0.5), ram_grid=param_grid, pring='roc_auc', rbose=0, jobs=10,</pre>
<pre>#sorted # from optimal #print(</pre>	<pre>d(sklearn.metrics.SCORERS.keys() sklearn.model_selection import metrics.roc_auc_score l_params.fit(X_train,</pre>
<pre>#we can # learn #scale_ # if th #optima</pre>	optimal_params.best_params_) in see the optimal output said by round one is gamma= 0.25, ining rate= 0.1, maxdepth=3,reg_lammbda=10.0, _pos_weight=3 the calue is a middle vaue then we try to keep that, that is already a al in the case, then if it is a lower boudn or upper bound of that candidate all test out for more values that based on that
## RROLL param_g 'ma 'le 'ga 're	UND 2
par sco ver	<pre>l_params = GridSearchCV(timator=xgb.XGBClassifier(objective='binary:logistic',</pre>
n_j cv=)	jobs=10,
# these #gamma' # 'scal	optimal_params.best_params_) e are the round 2 resuts we have ': 0.25, 'learning_rate': 0.1, 'max_depth': 2, 'reg_lambda': 20, le_pos_weight': 3
	we want to build our final XGBoost model and evlaute the results o= xgb.XGBClassifier(seed=21,
_	<pre>subsample=0.9,</pre>
	ave the following results: GBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
	<pre>validate_parameters=1, verbosity=None) hen plot another matrix and see how it goes, onfusion_matrix(clf_xgb,</pre>
<pre># out o #but fo #is muc #expens # i thi</pre>	values_format='d',
<pre># risk # pubni # peope #if we</pre>	poeple, because they are the potential cost more money and resourse to ish and prevent sth like that happening, whereas for the wrongly identify el we can give some for example couple and they will not be that offened identify them wrongly compared to leting go dangerous people out. D= xgb.XGBClassifier(seed=21, objective= 'binary:logistic', gamma= 0.25,
_	<pre>learn_rate=0.1, max_depth=2, reg_lambda=20, scale_pos_weight=3, subsample=0.9, colsample_bytree=0.5, n_estimators=1)</pre>
# outpu # XGBCI # # # # # # #	lassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
for imp	<pre>validate_parameters=1, verbosity=None) lf_xgb.get_booster() portance_type in ('weight', 'gain', 'cover', 'total_gain', 'total_cover'): ('%s:'% importance_type, bst.get_score(importance_type=importance_type)) arams= {'shape':'box', 'skyle': 'filled, rounded',</pre>
_	<pre>'fillcolor': '#78cbe'} arams={'shape':'box',</pre>
graph_d	<pre>data= xgb.to_graphviz(clf_xgb,num_tree=0,size="10,10",</pre>
#you wa numbe 0 1 2	level give you a choice of yes and no, it helps you determine which way and to go, then the leaf indicate the probability for each action you take. er Age Sex Job Housing Saving accounts Checking account \ 0 67 male 2 own NaN little 1 22 female 2 own little moderate 2 49 male 1 own little NaN
4	3 45 male 2 free little little 4 53 male 2 free little little it amount Duration Purpose Default 1169 6 radio/TV 0 5951 48 radio/TV 1 2096 12 education 0 7882 42 furniture/equipment 0 4870 24 car 1 er Age Sex Job Housing Saving accounts Checking account \
0 1 2 3 4	0 67 male 2 own NaN little 1 22 female 2 own little moderate 2 49 male 1 own little NaN 3 45 male 2 free little little 4 53 male 2 free little little it amount Duration Default 1169 6 0 5951 48 1 2096 12 0
Age 0 67 1 22 2 49 3 45 4 53	7882 42 0 4870 24 1 Sex Job Housing Saving accounts Checking account Credit amount \ male 2 own NaN little 1169 female 2 own little moderate 5951 male 1 own little NaN 2096 male 2 free little little 7882 male 2 free little little 4870
Durat Durat Durat Page Sex Job Housing	tion Default 6 0 48 1 12 0 42 0 24 1 int64 object int64 object
Saving a Checking Credit a Duration Default dtype: c ['male' ['own' ' [nan 'li	accounts object g account object amount int64 n int64 object object 'female'] 'free' 'rent'] ittle' 'quite rich' 'rich' 'moderate']
Age 0 67 1 22 2 49 3 45 4 53 Durat 0	Sex Job Housing Saving accounts Checking account Credit amount \ male 2 own NaN little 1169 female 2 own little moderate 5951 male 1 own little NaN 2096 male 2 free little little 7882 male 2 free little little 4870 tion Default 6 0 48 1
1 2 3 4 [nan 'li Age 0 67 1 22 2 49 3 45	48 1 12 0 42 0 24 1 ittle' 'quite_rich' 'rich' 'moderate'] Sex Job Housing Saving accounts Checking account Credit amount \ male 2 own NaN little 1169 female 2 own little moderate 5951 male 1 own little NaN 2096 male 2 free little little 7882
3 45 4 53 Durat 0 1 2 3 4 0 0	male2freelittlelittle7882male2freelittlelittle4870
1 1 2 0 3 0 4 1 Name: De Age Sex Job Housing	accounts object
Checking Credit a Duration dtype: Age 0 67 1 22 2 49 3 45	accounts object g account object amount int64 n int64 object Job Housing Saving accounts Checking account Credit amount Duration \ 2 own NaN little 1169 6 2 own little moderate 5951 48 1 own little NaN 2096 12 2 free little little 7882 42
4 53 Sex_f 0 1 2 3 4 Age	2 free little little 4870 24 female Sex_male 0
0 67 1 22 2 49 3 45 4 53 Housi 0 1	2 1169 6 0 1 0 2 5951 48 1 0 0 1 2096 12 0 1 0 2 7882 42 0 1 1 2 4870 24 0 1 1 ing_own Housing_rent Saving accounts_little \ 1 0 0 1 1 0 1
3	0 0 1 0 0 1 0 0 1 0 0 1 0 0 0 1 0
Check 0 1 2 3 4 [0 1] 0.3	king account_little
0.3 [0] [1] [2] [3] [4] [5] [6]	<pre>validation_0-aucpr:0.57051 validation_0-aucpr:0.57840 validation_0-aucpr:0.58978 validation_0-aucpr:0.54545 validation_0-aucpr:0.58128 validation_0-aucpr:0.55663 validation_0-aucpr:0.56723 validation_0-aucpr:0.56452</pre>
[9] [10] [11] [12] (750,) (250,) C:\Users e use of	<pre>validation_0-aucpr:0.58035 validation_0-aucpr:0.57406 validation_0-aucpr:0.57924 validation_0-aucpr:0.57970 validation_0-aucpr:0.58447 s\yorkuniversity\AppData\Roaming\Python\Python39\site-packages\xgboost\sklearn.py:1224: UserWarningf label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove to the following: 1) Pass option use label encoder=False when constructing XGBClassifier object;</pre>
2) Encoor warning. C:\Users e use of warning. 2) Encoor warning. {'gamma'}	de your labels (y) as integers starting with 0, i.e. 0, 1, 2,, [num_class - 1]. ngs.warn(label_encoder_deprecation_msg, UserWarning) s\yorkuniversity\AppData\Roaming\Python\Python39\site-packages\xgboost\sklearn.py:1224: UserWarnin f label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove , do the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; de your labels (y) as integers starting with 0, i.e. 0, 1, 2,, [num_class - 1]. ngs.warn(label_encoder_deprecation_msg, UserWarning) ': 0.25, 'learning_rate': 0.1, 'max_depth': 3, 'reg_lambda': 10.0, 'scale_pos_weight': 3} ': 0.25, 'learning_rate': 0.1, 'max_depth': 2, 'reg_lambda': 20, 'scale_pos_weight': 3}
Paramete This o	12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:576: ers: { "learn_rate" } might not be used. could be a false alarm, with some parameters getting used by language bindings but being mistakenly passed down to XGBoost core, or some parameter actually being used etting flagged wrongly here. Please open an issue if you find any such cases. validation_0-aucpr:0.51926 validation_0-aucpr:0.53845
[2] [3] [4] [5] [6] [7] [8] [9] [10]	<pre>validation_0-aucpr:0.53845 validation_0-aucpr:0.65706 validation_0-aucpr:0.69142 validation_0-aucpr:0.64115 validation_0-aucpr:0.63039 validation_0-aucpr:0.63246 validation_0-aucpr:0.65629 validation_0-aucpr:0.66024 validation_0-aucpr:0.66951 validation_0-aucpr:0.66155</pre>
[11] [12] [13] [14] [15] [16] [17] [18] [19]	<pre>validation_0-aucpr:0.66155 validation_0-aucpr:0.66247 validation_0-aucpr:0.65216 validation_0-aucpr:0.66287 validation_0-aucpr:0.63634 validation_0-aucpr:0.64834 validation_0-aucpr:0.65134 validation_0-aucpr:0.65217 validation_0-aucpr:0.65029</pre>
C:\Users e use of warning, 2) Encod warnir [01:09:1 Paramete This of then k	s\yorkuniversity\AppData\Roaming\Python\Python39\site-packages\xgboost\sklearn.py:1224: UserWarning f label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove to the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; de your labels (y) as integers starting with 0, i.e. 0, 1, 2,, [num_class - 1]. ngs.warn(label_encoder_deprecation_msg, UserWarning) 12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:576: ers: { "learn_rate" } might not be used. could be a false alarm, with some parameters getting used by language bindings but being mistakenly passed down to XGBoost core, or some parameter actually being used
then k but ge [01:09:1 in XGBooror' to weight: gain: {'e': 57.2	being mistakenly passed down to XGBoost core, or some parameter actually being used etting flagged wrongly here. Please open an issue if you find any such cases. 12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:1115: Star ost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior. {'Duration': 1.0, 'Checking account_little': 1.0, 'Checking account_moderate': 1.0} 'Duration': 11.417625427246094, 'Checking account_little': 53.7900505065918, 'Checking account_mode 282169342041016}
cover: { total_ga derate': total_co C:\Users e use of warning, 2) Encoo	{'Duration': 90.5, 'Checking account_little': 267.0, 'Checking account_moderate': 176.5} ain: {'Duration': 11.417625427246094, 'Checking account_little': 53.7900505065918, 'Checking account: 57.282169342041016} over: {'Duration': 90.5, 'Checking account_little': 267.0, 'Checking account_moderate': 176.5} s\yorkuniversity\AppData\Roaming\Python\Python39\site-packages\xgboost\sklearn.py:1224: UserWarning f label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove to the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; de your labels (y) as integers starting with 0, i.e. 0, 1, 2,, [num_class - 1].
CalledPr C:\Users 399 400 > 401 402	Could not open "xgboost_tree_risk_customer.pdf" for writing: Permission denied rocessError
403 ~\anacor 280 281> 282 283 284 ~\anacor	<pre>#interpreting the tree nda3\lib\site-packages\graphviz\files.py in view(self, filename, directory, cleanup, quiet, quiet_v</pre>
241 242 > 243 244 245 ~\anacor 223 224	<pre>formatter, quiet, quiet_view)</pre>
224> 225 226 227 ~\anacor *kwargs) 183 184> 185	<pre>run(cmd, capture_output=True, cwd=cwd, check=True, quiet=quiet) return rendered nda3\lib\site-packages\graphviz\backend.py in run(cmd, input, capture_output, check, encoding, quie) if check and proc.returncode: raise CalledProcessError(proc.returncode, cmd,</pre>
186 187 CalledPr ero exit denied	output=out, stderr=err) rocessError: Command '['dot.bat', '-Kdot', '-Tpdf', '-O', 'xgboost_tree_risk_customer']' returned ret status 1. [stderr: b'Error: Could not open "xgboost_tree_risk_customer.pdf" for writing: Permiss \r\n'] -140
low risk lapel lapel lapel high risk	- 100 - 80 - 60
low risk	low risk high risk Predicted label - 40 - 100 - 90 - 80
rue label	- 70 - 60 - 50 - 40
high risk	- 30 - 20

Predicted label

In []:

In []:

In [3]: # -*- coding: utf-8 -*-

Spyder Editor

This is a temporary script file.