

In [210]: model = LogisticRegression(penalty='l2')

• **#### DecisionTreeClassifier**

parameter:

1.criterion:

2.splitter:

3.max_features:

4.max_depth:

5.min_impurity_split:

Either "gini" or "entropy" can be used, the former representing the Gini coefficient and the latter representing the information gain. Generally speaking, it is enough to use the default Gini coefficient "gini".

Either "best" or "random" can be used. The former finds the optimal division point among all the division points of the feature. The latter is to randomly find the locally optimal dividing point among the partial dividing points.

Many types of values can be used, the default is "None", which means that all features are considered when dividing; if it is "log2", it means that at most log2N features are considered when dividing;

The maximum depth of the decision tree can not be entered by default. If it is not entered, the decision tree will not limit the depth of the subtree when building the subtree. Generally speaking, this value can be ignored when there are few data or features.

This value limits the growth of the decision tree. If the impurity of a node (Gini coefficient, information gain, mean square error, absolute difference) is less than this threshold, the node will no longer generate child nodes. is the leaf node.

In [211]: # model = tree.DecisionTreeClassifier()

• **#### GradientBoostingClassifier: Gradient Boosting is a Boosting method. Its main idea is that each time a model is built, the gradient descent direction of the model loss function is established before.**

In [212]: # model = GradientBoostingClassifier(n_estimators=200)

• **#### svm:**

The SVM classifier is a binary or discriminative model that distinguishes between two types of data

In [213]: # model = SVC(kernel='rbf', probability=True)

• **#### RandomForestClassifier**

The goal of the ensemble algorithm is to consider the modeling results of multiple evaluators and aggregate them to obtain a comprehensive result, in order to obtain better regression or classification performance than a single model.

In [214]: # model = RandomForestClassifier(n_estimators=1000,max_depth=10)

In [215]: model.fit(X_train, y_train)
prediction = model.predict(X_test)
report = classification_report(y_test, prediction, output_dict=True)

D:\Anaconda\lib\site-packages\sklearn\linear_model_logistic.py:818: ConvergenceWarning: lbfgs failed to converge (status=1):
 https://scikit-learn.org/stable/modules/preprocessing.html
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
 https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
 https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,

Evaluate

In [216]: df_evaluate = pd.DataFrame([["5330254",report["macro avg"],["precision"],report["macro avg"],["recall"],report["accuracy"]],
 columns=["zid","average_precision","average_recall","accuracy"]])
df_evaluate

Out[216]:

	zid	average_precision	average_recall	accuracy
0	5330254	0.676872	0.516665	0.83375

prediction output

In [217]: df_predict = pd.concat([pd.Series(X_test.index.values),pd.Series(prediction)],axis=1)
df_predict.columns = ['SK_ID_CURR','predicted_target']
df_predict

Out[217]:

	SK_ID_CURR	predicted_target
0	0	0
1	1	0
2	2	0
3	3	0
4	4	0
...
11995	11995	0
11996	11996	0
11997	11997	0
11998	11998	0
11999	11999	0

12000 rows × 2 columns

In [218]: df_predict.predicted_target.unique()

Out[218]: array([0, 1], dtype=int64)

Regression

In [219]: from sklearn import tree
from sklearn import neighbors
from sklearn.ensemble import RandomForestRegressor,AdaBoostRegressor,GradientBoostingRegressor,BaggingRegressor
from sklearn.linear_model import Ridge,LinearRegression
from sklearn.neighbors import KNeighborsRegressor
from sklearn.linear_model import RidgeCV
from sklearn.metrics import mean_squared_error, accuracy_score
from scipy.stats import pearsonr

In [220]: # Training set
X_train = trainData_regression_correlationFilter
y_train = targetColSeries_regression
y_trainMean = y_train.mean()
y_trainStd = y_train.std()
y_train = y_train.apply(lambda x: (x-y_trainMean)/y_trainStd)

test set
X_test = testData_regression_data
y_test = testData_regression_target
y_testMean = y_test.mean()
y_testStd = y_test.std()
y_test = y_test.apply(lambda x: (x-y_trainMean)/y_trainStd)

Ensure that the dimensions of the training set and the test set are the same
features = X_train.columns.values.tolist()
features = list(set(features).intersection(set(X_test.columns.values.tolist())))
X_train = X_train[features]
X_test = X_test[features]

LinearRegression

In [221]: # model = LinearRegression()

DecisionTreeRegressor

In [222]: # model = tree.DecisionTreeRegressor()

KNeighborsRegressor

In [223]: # model = KNeighborsRegressor()

RandomForestRegressor

In [224]: # model = RandomForestRegressor(n_estimators=20)

AdaBoostRegressor

In [225]: # model = AdaBoostRegressor(n_estimators=50)

GradientBoostingRegressor

In [226]: # model = GradientBoostingRegressor(n_estimators=100)

BaggingRegressor

In [227]: # model = BaggingRegressor()

Ridge Choose Ridge Regression Hyperparameter α

In [228]: ridgecv = RidgeCV(alphas=[0.01, 0.1, 0.5, 1, 3, 5, 7, 10, 20, 100])
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[228]: 5.0

In [229]: model = Ridge(alpha=ridgecv.alpha_)

In [230]: model.fit(X_train, y_train)
prediction = model.predict(X_test)

In [231]: df_evaluate = pd.DataFrame([["5330254",mean_squared_error(y_test, prediction),pearsonr(y_test, prediction)[0]]],
 columns=["zid","MSE","correlation"])
df_evaluate

Out[231]:

	zid	MSE	correlation
0	5330254	0.068077	0.508406

In [232]: df_predict = pd.concat([pd.Series(X_test.index.values),pd.Series(prediction.tolist())].apply(lambda x: x
y_testStd*(y_testMean)),axis=1)
df_predict.columns = ['SK_ID_CURR','predicted_income']
df_predict

Out[232]:

	SK_ID_CURR	predicted_income
0	0	168714.488444
1	1	143296.779749
2	2	173886.171656
3	3	197307.293532
4	4	153267.476605
...
11995	11995	156266.118464
11996	11996	171268.291516
11997	11997	184468.110381
11998	11998	168076.120025
11999	11999	167650.478022

12000 rows × 2 columns

In []: