ffffffi	<pre>from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay from sklearn.feature_selection import SelectKBest from sklearn.feature_selection import LabelEncoder from sklearn.feature_selection import chi2 from tility_func.utility import * from pandas_summary import DataFrameSummary from pandas_summary import bataFrameSummary from sklearn import metrics as met import joblib np.random.seed(42) np.set_printoptions(suppress=True) warnings.filterwarnings('ignore') plt.rcParams['figure.figsize'] = (12, 8)</pre>
p d	plt.rcParams['figure.figsize'] = (12, 8) plt.rcParams['font.size'] = 20 plt.rcParams['font.size'] = 20 plt.rcParams['font.size'] = (12, 8) plt.rcParams['font.size'] = 20 plt.rcParams['font.size'] = (12, 8) plt.rcParams['font.size'] = 20 plt.rcParam
Ra 2 3 4 5 6 7 8 9	class 'pandas.core.frame.DataFrame'> angeIndex: 400 entries, 0 to 399 ata columns (total 26 columns): # Column
111111111111111111111111111111111111111	11 bu 381 non-null float64 12 sc 383 non-null float64 13 sod 313 non-null float64 14 pot 312 non-null float64 15 hemo 348 non-null float64 16 pcv 330 non-null object 17 wc 295 non-null object 18 rc 270 non-null object 19 htn 398 non-null object 19 htn 398 non-null object 20 dm 398 non-null object 21 cad 398 non-null object 22 appet 399 non-null object 23 pe 399 non-null object 24 ane 399 non-null object 25 classification 400 non-null object
# d	<pre>types: float64(11), int64(1), object(14) emory usage: 81.4+ KB # rename the catagory column df.rename(columns={'classification': 'class'}, inplace=True) df.columns ndex(['id', 'age', 'bp', 'sg', 'al', 'su', 'rbc', 'pc', 'bgr', 'bu', 'sc', 'sod', 'pot', 'hemo', 'pcv', 'wc', 'rc', 'htn', 'dm', 'cad', 'apet', 'pe', 'ane', 'class'],</pre>
d Er	df.drop('id', axis=1, inplace=True) df.columns ndex(['age', 'bp', 'sg', 'al', 'su', 'rbc', 'pc', 'bc', 'bg', 'bu', 'sc', 'sod', 'pot', 'hemo', 'pcv', 'wc', 'rc', 'htn', 'dm', 'cad', 'appet', 'pe', 'ane', 'class'], dtype='object') 1 Data Preprocessing :Finding and Removing Anomalies
d C C	# Function to find the categorical and numerical features def extract_cat_num(df): cat_col=[col for col in df.columns if df[col].dtype=='0'] num_col=[col for col in df.columns if df[col].dtype!='0'] return cat_col,num_col cat_col,num_col=extract_cat_num(df) cat_col 'rbc', 'pc', 'pc', 'pc', 'pc', 'pc', 'pc',
n	'ba', 'pcv', 'wc', 'rc', 'thn', 'dm', 'cad', 'appet', 'pe', 'ane', 'class'] num_col 'age', 'bp', 'sg', 'al', 'su', 'bgr', 'bu', 'sc', 'sod', 'pot', 'hemo']
# f	dirtiness in data for col in cat_col: print('{\} has {\} values'.format(col,df[col].unique())) print("\n") bc has [nan 'normal' 'abnormal'] values c has ['normal' 'abnormal' nan] values
pc ''	cc has ['notpresent' 'present' nan] values a has ['notpresent' 'present' nan] values cv has ['44' '38' '31' '32' '35' '39' '36' '33' '29' '28' nan '16' '24' '37' '30' '34' '40' '45' '27' '48' '\t?' '52' '14' '22' '18' '42' '17' '46' '23' '19' '25' '41' '26' '15' '21' '43' '20' '\t43' '47' '9' '49' '50' '53' '51' '54'] values c has ['7800' '6000' '7500' '6700' '7300' nan '6900' '9600' '12100' '4500' '12200' '11000' '3800' '11400' '5300' '9200' '6200' '8300' '8400' '10300' '9800' '9100' '7700' '14600' '8600' '18000' '21600' '4300' '8500' '7200' '7700' '14600' '6300' '1500' '11800' '9400' '5500' '5800'
rc	'13290' '12500' '5600' '7000' '11900' '140400' '12700' '6800' '15000' '16300' '12400' '19000' '14900' '8200' '15200' '16300' '12400' '10500' '14900' '8200' '15200' '16300' '12400' '12500' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '12600' '1
dn ca	tn has ['yes' 'no' nan] values m has ['yes' 'no' ' yes' '\tno' '\tyes' nan] values ad has ['no' 'yes' '\tno' nan] values ppet has ['good' 'poor' nan] values e has ['no' 'yes' nan] values
# f	lass has ['no' 'yes' nan] values # Convert the object string features into float for col in ['pcv','wc','rc']: df[col] = df[col].str.extract('(\d+\.\d+ \d+)').astype(float) # Remove the tab in the categorical variables df['cad'].replace(to_replace=('\tno':'no'), inplace=True)
d d	<pre>df['dm'].replace(to_replace={\too':\no',\tyes':\yes',\'yes';\yes'}, inplace=True) df['class'].replace(to_replace={\too':\no',\tyes':\yes',\'yes':\yes'}, inplace=True) 2 Save the preprocessed data as.makedirs('./input/processed/', exist_ok=True) df.to_feather('./input/processed/ckd-processed') df = feather.read_dataframe('input/processed/ckd-processed') df.head()</pre>
1 2 3 4	age bp sg al su rbc pc ba bgr bu sc bu sc so hu sc bu sc so pc bu sc so pc we rc hu du sc rc hu sc rc hu sc rc hu sc
t d <c Ra Da #</c 	#Change columns of strings in dataframe to a column of categorical values train_cats(df) df.info() class 'pandas.core.frame.DataFrame'> angeIndex: 400 entries, 0 to 399 ata columns (total 25 columns): # Column Non-Null Count Dtype
2 5 6 7 8 9 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	3 al 354 non-null category 5 rbc 248 non-null category 6 pc 335 non-null category 7 pcc 396 non-null category 8 ba 396 non-null category 9 bgr 356 non-null float64 11 sc 383 non-null float64 12 sod 313 non-null float64 14 hemo 348 non-null float64 15 pcv 329 non-null float64 16 wc 294 non-null float64 17 rc 269 non-null float64 18 category
d d d	18 htn 398 non-null category 19 dm 398 non-null category 20 cad 398 non-null category 21 appet 399 non-null category 22 pe 399 non-null category 23 ane 399 non-null category 24 class 400 non-null category types: category(14), float64(11) emory usage: 41.6 KB df['class'].cat.categories ndex(['ckd', 'notckd'], dtype='object')
sg al sc oc oc oc oc oc	# collect the categorical varibales in dictionary for cols in df.columns: if df[cols].dtype.name == 'category': print(cols, dict(enumerate(df[cols].cat.categories))) g {0: 1.005, 1: 1.01, 2: 1.015, 3: 1.02, 4: 1.025} 1 {0: 0.0, 1: 1.0, 2: 2.0, 3: 3.0, 4: 4.0, 5: 5.0} u {0: 0.0, 1: 1.0, 2: 2.0, 3: 3.0, 4: 4.0, 5: 5.0} bc {0: 'abnormal', 1: 'normal'} c {0: 'abnormal', 1: 'normal'} c {0: 'notpresent', 1: 'present'} a {0: 'notpresent', 1: 'present'} tn {0: 'no', 1: 'yes'} m {0: 'no', 1: 'yes'}
ap oe ar cl	<pre>ad {0: 'no', 1: 'yes'} ppet {0: 'good', 1: 'poor'} e {0: 'no', 1: 'yes'} ne {0: 'no', 1: 'yes'} lass {0: 'ckd', 1: 'notckd'} # convert the appet and class categories for 1 for good/positive and 0 for poor/ negative to make constitent df['appet'].cat.set_categories(['poor', 'good'], ordered=True, inplace=True) df['class'].cat.set_categories(['notckd', 'ckd'], ordered=True, inplace=True) for cols in df.columns: if df[cols].dtype.name == 'category':</pre>
al surb occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occor occo occor occor occor occor occor occor occor occor occor occor occo occor occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo occo oc oc	<pre>print(cols, dict(enumerate(df[cols].cat.categories))) g {0: 1.005, 1: 1.01, 2: 1.015, 3: 1.02, 4: 1.025} l {0: 0.0, 1: 1.0, 2: 2.0, 3: 3.0, 4: 4.0, 5: 5.0} u {0: 0.0, 1: 1.0, 2: 2.0, 3: 3.0, 4: 4.0, 5: 5.0} bc {0: 'abnormal', 1: 'normal'} c {0: 'abnormal', 1: 'present'} a {0: 'notpresent', 1: 'present'} tn {0: 'no', 1: 'yes'} m {0: 'no', 1: 'yes'} ppet {0: 'noor, 1: 'yes'} ppet {0: 'poor', 1: 'yes'} e {0: 'no', 1: 'yes'} e {0: 'no', 1: 'yes'} e {0: 'no', 1: 'yes'} </pre>
# dd	lass {0: 'notckd', 1: 'ckd'} # Save the semi processed data df.to_feather('./input/processed/ckd_semi_processed') #. Creating Feature Variable and Target Variable df = feather.read_dataframe('input/processed/ckd_semi_processed') df.head()
1 2 3 4	age by sg al su rbc pc ba bg bu sc sod pc hemo pc bu sc sod pc wc rc htm dm class 48.0 80.0 1.020 1.0 0.0 NaN notpresent 121.0 36.0 1.2 NaN 11.0 48.0 70.0 1.00 1.00 NaN notpresent 121.0 36.0 1.2 NaN NaN 11.3 38.0 600.0 NaN no no no ro
al st pc ba ht dm ca pc ar	g {0: 1.005, 1: 1.01, 2: 1.015, 3: 1.02, 4: 1.025} l {0: 0.0, 1: 1.0, 2: 2.0, 3: 3.0, 4: 4.0, 5: 5.0} u {0: 0.0, 1: 1.0, 2: 2.0, 3: 3.0, 4: 4.0, 5: 5.0} bc {0: 'abnormal', 1: 'normal'} c {0: 'abnormal', 1: 'present'} a {0: 'notpresent', 1: 'present'} ftn {0: 'no', 1: 'yes'} m {0: 'no', 1: 'yes'} ppet {0: 'poor', 1: 'yes'} ppet {0: 'poor', 1: 'yes'} ne {0: 'no', 1: 'yes'} lass {0: 'no', 1: 'yes'}
f f ((# splits off the response variable, and replaced NAN by the median value of the column. features, response, nas = proc_df(df, 'class') features.shape, response.shape (400, 35), (400,)) features.info()
Rad	class 'pandas.core.frame.DataFrame'> angeIndex: 400 entries, 0 to 399 ata columns (total 35 columns): # Column Non-Null Count Dtype
111111111111111111111111111111111111111	11 sc
dt me	26 bg_na 400 non-null bool 27 bu_na 400 non-null bool 28 sc_na 400 non-null bool 29 sod_na 400 non-null bool 31 hemo_na 400 non-null bool 32 pcv_na 400 non-null bool 33 wc_na 400 non-null bool 33 wc_na 400 non-null bool 400 non-
' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' '	'bp': 80.0, 'bg': 121.0, 'bu': 42.0, 'sc': 1.3, 'sod': 138.0, 'pot': 4.4, 'hemo': 12.6499999999999, 'pcv': 40.0, 'rc': 4.8} cols = features.columns cols
3.	ndex(['age', 'bp', 'sg', 'al', 'su', 'rbc', 'pc', 'bc', 'bgr', 'bu',
× 0 1 2	**************************************
3. ×	.2 Creating Training set and Test set x_train, x_test, y_train, y_test = train_test_split(x_features_scaled, response, test_size=0.25, stratify=response, random_state=42) x_train.shape, y_train.shape, x_test.shape, y_test.shape (300, 35), (300,), (100, 35), (100,))
18 23 24	x_train.head() x_train.head()
3 3 3 3 3 3	**************************************
At	y_train.head() ttributeError
×	<pre>x_train.drop(x_train.iloc[:,24:], axis=1, inplace=True) x_test.drop(x_test.iloc[:,24:], axis=1, inplace=True) x_train.shape, y_train.shape, x_test.shape, y_test.shape x_train.head() x_test.head()</pre>
S m	** willity function to print cross k-fold cross validation scores, their mean and variance(standard deviation) ** def print_scores(model, x_train, y_train, cv, scoring): print('Cross validation scores:', cross_val_score(model, x_train.values, y_train, cv=cv, scoring=scoring, n_jobs=-1)) print('Mean_scores:', np.mean(cross_val_score(model, x_train.values, y_train, cv=cv, scoring=scoring, n_jobs=-1)) print('Variance:', np.std(cross_val_score(model, x_train.values, y_train, cv=cv, scoring=scoring, n_jobs=-1))
r p	In Random Forest Classifier Andom Forest Classifier
r s m	random_forest_preds = random_forest.predict(x_test.values) random_forest_score = met.accuracy_score(y_test, random_forest_preds) * 100 score.append(random_forest_score) model.append('Random Forest Classifier') print('Random Forest Classifier Accuracy = ', random_forest_score) cf_matrix = confusion_matrix(y_test, random_forest_preds) print('Condusion Matrix:',cf_matrix)
#	#Utility function for plotting confusion matrix def plot_cm(y_true, y_pred): cm = met.confusion_matrix(y_true, y_pred) sns.heatmap(cm, annot=True, cmap='Blues') plt.title('Confusion Matrix') plt.xlabel("y_pred") plt.ylabel("y_true") plt.show()
5	plot_cm(y_test, random_forest_preds) 5. Saving a Model for Deployment globlib.dump(scaler, './model/minmax_ckd.pkl') globlib.dump(random_forest_, './model/random_forest_ckd.pkl')
# r	<pre>joblib.dump(random_forest, './model/random_forest_ckd.pkl') # Loading the saved model for prediction random_forest_load = joblib.load('./model/random_forest_ckd.pkl') prediction = random_forest_load.predict(x_test) prediction== random_forest_preds</pre>
o x	Feature Important features_scaled.head() features=x_features_scaled.iloc[:, :24] features
n	response cat_col num_col df=df astyne('string') fillna('0')
1	<pre>df=df.astype('string').fillna('0') le=LabelEncoder() for col in cat_col: df[col]=le.fit_transform(df[col]) ind_col=[col for col in df.columns if col!='class'] dep_col='class'</pre>
i	features=df[ind_col] response=df[dep_col] imp_features=SelectKBest(score_func=chi2, k=20) imp_features=imp_features.fit(features, response) imp_features
i d f r	imp_features.scores_
ii dd dd dd ffffff	datascore=pd.DataFrame(imp_features.scores_,columns=['Score']) dfcols=pd.DataFrame(features.columns) dfcols features_rank=pd.concat([dfcols,datascore],axis=1) features_rank.columns=['features','score'] features_rank
i d d d d f f f f f f f f f f f f f f f	datascore=pd.DataFrame(imp_features.scores_,columns=['Score']) dfcols=pd.DataFrame(features.columns) dfcols features_rank=pd.concat([dfcols,datascore],axis=1) features_rank.columns=['features','score']

1. Imports modules and necessary packages