Machine Learning Term Project Report

109550186 - 李嘉玲

I. Preprocessing Procedure

• Missing value in numerical features (solved by data imputation)

```
# df1n2.columns[df1n2.isnull().any()]
         df1n2.isnull().sum()

    fixed_acidity

         volatile_acidity
citric_acid
                                     44
         residual sugar
         chlorides
free_sulfur_dioxide
total_sulfur_dioxide
         density
         sulphates
         alcohol
         dtype: int64
/ [652] df1n2 = df1n2.fillna(df1n2.mean())
/ [653] df1n2.isnull().sum()
         fixed_acidity volatile_acidity
         citric_acid
residual_sugar
         chlorides
         free_sulfur_dioxide
         total_sulfur_dioxide
         density
         sulphates
         alcohol
         class
```

The method used is filling the missing values with mean values of the non-empty value from that feature. In the code above, it can be seen that I used the panda's library *dataframe.fillna()* to replace the NaN value in the dataset.

Detecting outliers (includes data cleaning and reduction)

First, I tried to output the number of outliers to see if there is outlier in the dataset. If there is outlier, it will be dropped from the dataset and index will be reset. Again, I used the panda's library to detect and solve this issue.

• For optimized result (data transformation)

```
[663] # get optimized result
    sc = StandardScaler()

/ [663] # get optimized result
// StandardScaler()

/ [663] # get optimized result
/ [663] # get optimized result
/ StandardScaler()

/ [663] # get optimized result
/ [663] # get optimized resul
```

The standard scaler from sklearn preprocessing library is used to removes mean and scales the variables to unit variance.

• Preprocessing data for caategorical values

```
1371 # Preprocess data before training
     from keras.preprocessing.text import text to word sequence
     def preprocess data(df):
         reviews = []
         for raw in tqdm(df['Phrase']):
            # print(raw)
            text = raw.replace("'","")
            to kenized\_train\_data = text\_to\_word\_sequence(text,filters='!"#$%&()*+,-./:;<=>?@[\\]^_`{|}~\tn',split="")
             stop_words = set(stopwords.words('english'))
            removesw = [i for i in tokenized_train_data if not i in stop_words]
             rswtext = ' '.join(removesw)
            numberRemove = ''.join(num for num in rswtext if not num.isdigit())
             stemmer = PorterStemmer()
             stem_input = nltk.word_tokenize(numberRemove)
             stem_text = ' '.join([stemmer.stem(word) for word in stem_input])
            reviews.append(stem_text)
```

For categorical values, the methods used are removing special characters like parantheses and other symbols and stop words, as well as morphological and inflexional endings from the sentences given in the variable phrase of the dataframe.

II. Decision Tree Classifier Algorithm Snapshots

Decision tree classifier built using the gini and entropy measurement as shown in the functions below. Function *gini* determines the counts of each distinct label in the set. The likelihood of each label is then calculated, and one minus the sum of the squares of these probabilities is returned. If the labels are pure, then the Gini impurity will be low. While function *entropy* determines the counts of each distinct label in the set. The chance of each label is then calculated, and the result is the negative sum of the probabilities multiplied by the probabilities logarithm.

```
# functions used to construct the decision tree classifier
def gini(sequence, weights=None):
    if weights is None: # count 1 - sum of square of probabilities
        _, counts = np.unique(sequence, return_counts=True)
        p = (counts / len(sequence)) ** 2
        return 1.0 - np.sum(p)
    else:
        tot = 0
        weights = weights / weights.sum()
        for c in np.unique(sequence):
            #change prob become weighted prob
            tot = np.sum(weights[sequence == c]) ** 2
        return 1 - tot
def entropy(sequence, weights=None):
    if weights is None:
        _, counts = np.unique(sequence, return_counts=True)
        p = counts / len(sequence) #The Probability
        return -np.sum(p * np.log2(p))
        entropy = 0
        weights = weights / weights.sum()
        for c in np.unique(sequence):
            # calculate the weighted probability
            tmp = np.sum(weights[sequence == cl)
           entropy -= tmp * np.log2(tmp)
        return entropy
```

The image below shows several variables of this decision tree, which are criterion, max_depth, max_features and n_features. Criterion is used to measure the quality of a split, with default Gini impurity, but also can be entropy. Max_depth is the maximum depth the tree can have. Max_features and n_features respectively shows the most features the tree will take into account while determining the appropriate split at each node and the number of features in the input data.

```
def __init__(self, criterion='gini', max_depth=None, max_features=None):
    self.criterion = globals()[criterion]

self.max_depth = max_depth if max_depth is not None else 2 ** 100
    self.max_features = max_features
    self.n_features = None
```

Function *fit* is the one that build the decision tree. It takes input of training data and target labels as well as optional sample weight. While *get_node* function is called recursively by *fit* to build the decision tree. It takes input of input data, target labels, current node depth, and optional sample weight. Then, it will return Node, which if current depth greater than or equal to max_depth, will return a leaf node with prediction equal to most common class in data. Else, it will find best feature and threshold to split the data at this node by calling *best_split*, and it recursively calls *get_node* on the data in the left and right child nodes.

```
def fit(self, x_data, y_data, sample_weight=None):
    self.n_features = x_data.shape[1]
    self.root = self.get_node(x_data, y_data, depth=0, sample_weight=sample_weight)

def get_node(self, x, y, depth, sample_weight=None):
    weighted_counts = np.bincount(y, weights=sample_weight)
    prediction = np.argmax(weighted_counts)

    node = Node(c_value=self.criterion(y, sample_weight), prediction=prediction)
    if depth >= self.max_depth:
        return node

    node.feature_idx, node.threshold = self.best_split(x, y, sample_weight)
    if node.feature_idx is None:
        return node

left_idx = x[:, node.feature_idx] < node.threshold
    x_left, y_left = x[left_idx], y[left_idx]
    x_right, y_right = x[~left_idx], y[~left_idx]

# get child nodes recursively
    if sample_weight is not None:
        node.left = self.get_node(x_left, y_left, depth=depth + 1, sample_weight=sample_weight[left_idx])
        node.right = self.get_node(x_right, y_right, depth=depth + 1)
        node.left = self.get_node(x_left, y_left, depth=depth + 1)
        node.right = self.get_node(x_right, y_right, depth=depth + 1)
        return node</pre>
```

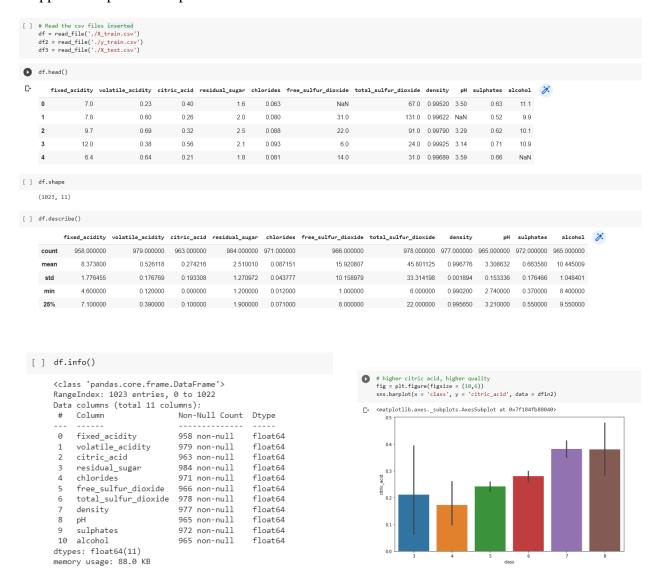
Function *best_function* is used to find best feature and threshold to split data. If max_features is set, it randomly selects a subset of the features to consider when searching for the best split. *Predict* function returns a list of predicted labels by traversing tree and returning prediction at the leaf node that the sample ends up.

```
def best_split(self, x, y, sample_weight):
    if len(y) <= 1:
    parent_c = self.criterion(y, sample_weight)
    best_infog = -2 ** 64
    best_idx, best_th = None, None
         available_features = np.random.choice(np.arange(self.n_features), size=self.max_features, replace=False)
         available_features = np.arange(self.n_features)
    for idx in available_features:
         sort_idx = np.argsort(x[:, idx])
thresholds = x[sort_idx, idx]
         labels = y[sort idx]
         for pos in range(1, len(y)):
              if thresholds[pos] == thresholds[pos - 1]:
              if sample_weight is not None:
                   sorted_sample_weight = sample_weight[sort_idx]
left_c = self.criterion(
                   labels[:pos], sorted_sample_weight[:pos])
right_c = self.criterion(
                       labels[pos:], sorted_sample_weight[pos:])
                   left_c = self.criterion(labels[:pos])
                   right_c = self.criterion(labels[pos:])
                          left c = self.criterion(labels[:pos])
                     child_c = (pos * left_c + (len(y) - pos) * right_c) / len(y)
                     if infog > best_infog:
  best_infog = infog
  best_idx = idx
  best_th = (thresholds[pos] + thresholds[pos - 1]) / 2
             return best_idx, best_th
        def predict(self, x_data):
             def util(self, x):

    cur_node = self.root

    while cur_node.left and cur_node.right:
                     if x[cur_node.feature_idx] < cur_node.threshold:
    cur_node = cur_node.left</pre>
                    else:
cur_node = cur_node.right
                 return cur node.prediction
             return np.stack([util(self, single_x) for single_x in x_data])
```

Snippets of input and output in dataset 1 and 2:



Barplot features show that some feature might not be that crucial and can have less weight compared to features with higher role in the prediction algorithm.

III. Result

Dataset 1 training dataset accuracy:

```
# predict validation set
   y pred = clf.predict(X val)
   report(y_pred, y_val)
Accuracy: 0.6
   Confusion Matrix:
   [[0 4 3 0 0]
    [ 0 62 18 1 0]
    [ 0 30 45 11 0]
    [0 3 8 16 0]
    [0 1 1 2 0]]
   Classification Report:
               precision
                         recall f1-score support
                  0.00
                        0.00
                                    0.00
            5
                  0.62
                          0.77
                                   0.69
                                              81
                                              86
                   0.60
                        0.52
                                    0.56
            6
            7
                   0.53
                           0.59
                                    0.56
                                              27
                   0.00
                           0.00
                                    0.00
                                    0.60
                                              205
      accuracy
                 0.35 0.38
                                   0.36
                                              205
     macro avg
   weighted avg
                  0.57
                            0.60
                                    0.58
                                              205
```

Dataset 2 training dataset accuracy:

```
# predict validation set
   y_pred = clf.predict(X_val_)
   report(y_pred, y_val)
Accuracy: 0.525830997196636
   Confusion Matrix:
   [[ 45 106 931
                            11]
    [ 42 357 3782 151
                           32]
    [ 20 282 12118 253
[ 26 130 4487 509
                            60]
                           116]
           28 1133 204 101]]
   Classification Report:
               precision
                         recall f1-score support
                                    0.07
                                            1132
            0
                  0.32
                          0.04
                   0.40
                           0.08
                                    0.14
                                             4364
            1
                   0.54
                            0.95
                                    0.69
                                             12733
                        0.10
                  0.44
             3
                                     0.16
                                             5268
                  0.32
                          0.07
                                             1473
            4
                                     0.11
                                     0.53
                                             24970
      accuracy
               0.40
                            0.25
                                     0.23
                                             24970
      macro avg
   weighted avg
                  0.47
                            0.53
                                    0.42
                                             24970
```

IV. Submission

- Preprocessed training dataset: X_train_final
- Preprocessed testing dataset: X_test_final
- Predicted result of testing dataset: y_test_final