

### Categorical Visualization

```
[1]: import pandas as pd

df = pd.read_csv('train.csv', sep='\\t')

from os import path
import matplotlib.pyplot as plt

col_to_transform = list(set(df.columns)-set(df_get_numeric_data().columns))
new_col=[]
for col in col_to_transform:
    col = '_' +
    new_col.append(col)

good = df[df['Label']==1]
bad = df[df['Label']==0]

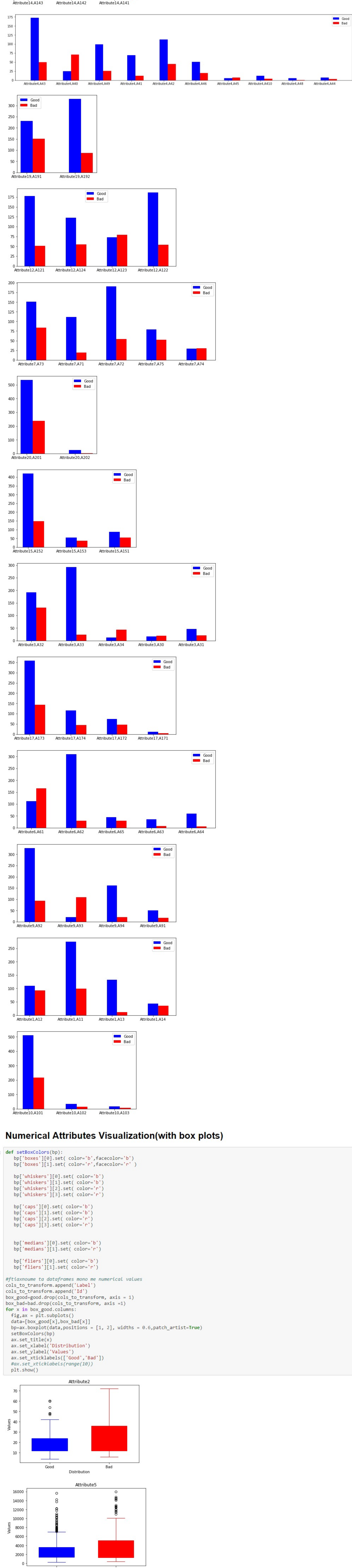
col_to_drop = df_get_numeric_data().columns
good_cat = good.drop(col_to_drop, axis = 1)
good_dict = good_cat.to_dict(orient='records')
bad_cat = bad.drop(col_to_drop, axis = 1)
bad_dict = bad_cat.to_dict(orient='records')

from sklearn.feature_extraction import DictVectorizer as DV

vectorizer = DV(separator='_', sparse=False)
goodvec = vectorizer.fit_transform(good_dict)
badvec = vectorizer.transform(bad_dict)
goodd = pd.DataFrame(goodvec,columns=vectorizer.get_feature_names())
badd = pd.DataFrame(badvec,columns=vectorizer.get_feature_names())

for single_col in new_col:
    filter_col=[col for col in list(vectorizer.get_feature_names()) if col.startswith(single_col)]
    goodd[filter_col]
    badd[filter_col]

from collections import Counter
cnt = Counter()
for index,row in goodd.iterrows():
    for i in filter_col:
        if(row[i]==1):
            cnt[i]+=1
for index,row in badd.iterrows():
    for i in filter_col:
        if(row[i]==1):
            cnt[i]-=1
import numpy as np
l = np.arange(len(cnt.keys()))
width=0.5
fig = plt.figure(figsize=(2 * len(cnt.keys()), 4))
ax = fig.add_subplot(111)
bars1 = ax.bar(l, cnt.values(), width,linewidth=2 * len(cnt.keys()) * width,align='center',color='b')
bars2 = ax.bar(l+width,cnt.values(),width,linewidth=2 * len(cnt.keys()) * width,align='center',color='r')
plt.xticks(cnt.keys())
ax.legend(['bars1[]', bars2[]], ('Good', 'Bad'))
plt.show()
```



Good                  Bad  
Distribution

---

*Attribute8*

The figure displays five box plots, each representing a different attribute. The y-axis for all plots is labeled 'Values'. The x-axis is labeled 'Distribution' with two categories: 'Good' and 'Bad'.

- Distribution:** The y-axis ranges from 1.0 to 4.0. The 'Good' category (blue box) has a median around 2.5, with whiskers extending from approximately 1.2 to 3.8. The 'Bad' category (red box) has a median around 2.5, with whiskers extending from approximately 1.2 to 3.8.
- Attribute11:** The y-axis ranges from 1.0 to 4.0. The 'Good' category (blue box) has a median around 2.5, with whiskers extending from approximately 1.2 to 3.8. The 'Bad' category (red box) has a median around 2.5, with whiskers extending from approximately 1.2 to 3.8.
- Attribute13:** The y-axis ranges from 20 to 70. The 'Good' category (blue box) has a median around 35, with whiskers extending from approximately 25 to 65. The 'Bad' category (red box) has a median around 35, with whiskers extending from approximately 25 to 65. Outliers are present for both categories at approximately 68.
- Attribute16:** The y-axis ranges from 1.0 to 4.0. The 'Good' category (blue box) has a median around 2.0, with whiskers extending from approximately 1.2 to 3.2. The 'Bad' category (red box) has a median around 2.0, with whiskers extending from approximately 1.2 to 3.2. Outliers are present for both categories at approximately 3.8.
- Attribute18:** The y-axis ranges from 1.0 to 2.0. The 'Good' category (blue box) has a median around 1.0, with whiskers extending from approximately 0.8 to 1.0. The 'Bad' category (red box) has a median around 1.0, with whiskers extending from approximately 0.8 to 1.0. Outliers are present for both categories at approximately 2.0.

**Classification : SVM , RandomForest , GaussianNE**

```
In [5]: import pandas as pd
import numpy as np

dataframe = pd.read_csv('train.tsv', sep='\t')
Y = dataframe['Label']

dataframe = dataframe.drop('Label',axis = 1)
dataframe = dataframe.drop('Y',axis = 1)

dataX = pd.get_dummies(dataframe)
X = dataX.as_matrix()

In [4]: from scipy import interp
from sklearn import svm
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import Kfold
from sklearn.model_selection import ShuffleSplit
from sklearn.metrics import accuracy_score

# SVM
k = 10
kf = ShuffleSplit(n_splits=k, test_size=0.1, random_state=0)
acc = 0

for train_index, test_index in kf.split(dataX):
    parameters = {'C':[1, 10]}
    svm = svm.LinearSVC()
    clf = GridSearchCV(svm, parameters)
    clf = clf.fit(X[train_index], Y[train_index])
    yPred = clf.predict(X[test_index])

    # Calculating accuracy
    acc += accuracy_score(yPred, Y[test_index])

ea = np.zeros((1,3)) #storing mean accuracy for EvaluationMetric.csv
en[0,0] = acc/k
print ("Average Accuracy of SVC_linear : " + str(acc/k))

Average Accuracy of SVC_linear : 0.6675

In [5]: from sklearn.ensemble import RandomForestClassifier

kf = ShuffleSplit(n_splits=k, test_size=0.1, random_state=0)
acc = 0

for train_index, test_index in kf.split(dataX):
    clf = RandomForestClassifier(n_estimators = 20, n_jobs = -1)
    clf = clf.fit(X[train_index], Y[train_index])
    yPred = clf.predict(X[test_index])

    # Calculating accuracy
    acc += accuracy_score(yPred, Y[test_index])

Average Accuracy of RandomForest : 0.76375

In [6]: from sklearn.naive_bayes import GaussianNB

kf = ShuffleSplit(n_splits=k, test_size=0.1, random_state=0)
acc = 0

for train_index, test_index in kf.split(X):
    clf = GaussianNB()

    clf = clf.fit(X[train_index], Y[train_index])
    yPred = clf.predict(X[test_index])

    # Calculating metrics
    acc += accuracy_score(yPred, Y[test_index])

print(Average Accuracy GaussianNB: " + str(acc/k))
en[0,2] = acc/k

Average Accuracy GaussianNB: 0.705

Output_Files

In [7]: #evaluation metric
columns = ["svm","RandomForest","GaussianNB"]
rows = ["Accuracy"]
mycsv = pd.DataFrame(columns=columns, index=rows, dtype=object)
mycsv.to_csv("EvaluationMetric_10fold.csv")
print ("Created EvaluationMetric_10fold.csv file!")

Created EvaluationMetric_10fold.csv file!

In [8]: test = pd.read_csv('test.tsv', sep='\t')

clf = RandomForestClassifier(n_estimators = 20, n_jobs = -1)
clf = clf.fit(X, Y)
yPred = clf.predict(X)

testX = test.drop('id',axis = 1)
testX = pd.get_dummies(testX)

predicted = clf.predict(testX)

output = np.zeros((len(testPred), 3), dtype=float)
```

```
for i,j in zip(pr
    if i == 1:
        output[i]
```

```

        output[j][i] = "Good"
    elif i == 1:
        output[j][i] = "Bad"
    output[j][0] = test.iloc[j]["ID"]
    output[j][1] = "Bad"

my_csv = pd.DataFrame(data=output, columns=["Client_ID", "Predicted_Label"])
my_csv.to_csv("testset_categories.csv", sep="\\", index=False, header=False)
print("Created testset_categories.csv file")

Created testset_categories.csv file!

```

## Information Gain & Accuracy

```

In [9]: import math
def nt_entropy(data):
    gdict={}
    vdata=[Label']

```

```
for i in range(1, n):
    if i==1:
        goods+=1
    p_good= goods/data.shape[0]
```

```

i=(p_good * math.log(p_good))
i=(p_bad * math.log(p_bad))
s_fin = s_fin2
return s_fin

def info_sum_attribute(arg):
    temp_list=list(of.columns)
    nonused_cols=[] for v in temp_list if not v.startswith('Label')]
    nonused_cols.remove(arg)
    newdf=df.drop(nonused_cols,axis=1)
    sum=0
    myvals=newdf[arg].value_counts()
    mycount=0
    for i in myvals.index:
        of=temp-newdf[newdf[arg]==i]
    import math
    def my_entropy(data):
        goods=[Label']
        v=data[Label']
        for i in V:
            if i==1:
                goods+=1
            if data.shape[0]==0:
                return 0
            else:
                p_good=goods/data.shape[0]
                p_bad=1-p_good
                i=(p_good * math.log(p_good))
                i=(p_bad * math.log(p_bad))
                s_fin = s_fin2
                return s_fin
    def info_sum_attribute(arg):
        temp_list=list(of.columns)
        nonused_cols=[] for v in temp_list if not v.startswith('Label')]
        nonused_cols.remove(arg)
        newdf=df.drop(nonused_cols,axis=1)
        sum=0
        myvals=newdf[arg].value_counts()
        mycount=0
        for i in myvals.index:
            of=temp-newdf[newdf[arg]==i]
            temp_my_entropy(of_temp)
            amount_to_sum=myvals[mycount]/newdf.shape[0] * temp2
            sum+= amount_to_sum
            mycount+=1
        return sum
    def info_num_sum(arg):
        temp_list=list(of.columns)
        nonused_cols=[] for v in temp_list if not v.startswith('Label')]
        nonused_cols.remove(arg)
        newdf=df.drop(nonused_cols,axis=1)
        sum=0
        (counters,bins)=mp.histogram(newdf[arg],5)
        print(counters,bins)
        mycount=0
        j=0
        for i in counters:
            if j<4:
                of=temp-newdf[(newdf[arg]>=bins[j]) & (newdf[arg]<bins[j+1])]
            else:
                of=temp-newdf[(newdf[arg]>=bins[j]) & (newdf[arg]<=bins[j+1])]
            j+=1
            temp2_my_entropy(of_temp)
            amount_to_sum=i/newdf.shape[0] * temp2
            sum+= amount_to_sum
        return sum
    info_list=[]
    print('Entropy : ', str(my_entropy(df)))
    cols_to_drop=[] for v in cols_to_drop if not v.startswith('Label')]
    cols_to_drop.append(cols_to_drop,axis=1)
    losdata=catdata.drop(Label',axis=1)
    for i in losdata.columns:
        infogain=my_entropy(df)-info_sum_attribute(i)
        info_list.append((infogain,i))
    cols_to_transform=[] for i in cols_to_transform if not i.startswith('Label')]
    numerals=catdata.drop(cols_to_transform,axis=1)
    losdata=numerals.data.drop(Label',axis=1)
    print('Counters & bins')
    for i in losdata.columns:
        infogain=num_entropy(df)-info_num_sum(i)
        info_list.append((infogain,i))
    Entropy = 0.8975737373656
    Counters & bins
    [[56 39 79 44 42] [ 4. 17.6 31.2 44.8 58.4 72.]
    [547 165 22 240] [ 250. 3389. 6528. 9667. 12886. 15945.]
    [131 184 6 120 188] [ 1. 1.6 2.2 2.8 3.4 4.]
    [107 242 6 122 193] [ 1. 1.6 2.2 2.8 3.4 4.]
    [130 275 140 54 122] [ 19. 36.2 41.4 52.6 63.8 75.]
    [513 261 0 22 41] [ 1. 1.6 2.2 2.8 3.4 4.]
    [682 0 0 0 118] [ 1. 1.2 1.4 1.6 1.8 2.]

In [18]: info_list.sort(key=lambda x: x[0])
import operator
acc=[]
info_list.sort(key=lambda x: x[0])
df = pd.read_csv('train.csv', sep='t')
t1s=[Label', 'S']
df=df.drop(t1s,axis=1)

for col in info_list:
    if len(df.columns)>2:
        break
    df = df.drop(col[1],axis=1)
    datax = pd.get_dummies(df)
    X = datax.as_matrix()

    k = 10
    k = 5
    df = ShuffleSplit(n_splits=k, test_size=0.1, random_state=0)
    acc = 0

    for train_index, test_index in kf.split(datax):
        X_train=X[train_index]
        y_train=y[test_index]
        y_pred = clf.predict(X_train)
        acc += accuracy_score(y_pred, y[test_index])

mean_acc = acc/k
accs.append(mean_acc)

info_df=pd.DataFrame(info_train)
info_df.of=df['featureSelection=InfoGainInfoTrain']
# Plot confusions
from numpy import array
a = np.zeros((len(info_list),),dtype=object)
count=0
for col in info_list:
    [a[count] = col[1]]
    count+= 1

```

[illegible]

```
data = pd.get_dummies(tr)
X = data.as_matrix(X)
clf = RandomForestClassifier(n_estimators = 20 ,n_jobs = -1)
clf = clf.fit(X, Y)
yPred = clf.predict(X)

testX = test.drop('Id',axis = 1)
testX = pd.get_dummies(testX)
```