

In [1]:

```
import numpy as np
import pandas as pd
from time import time
from IPython.display import display
import matplotlib.pyplot as plt
import visuals as vs
```

```
%matplotlib inline
```

```
data = pd.read_csv("census.csv")
```

```
display(data.head(n=1))
```

	age	workclass	education_level	education-num	marital-status	occupation	relationship	race	se
0	39	State-gov	Bachelors	13.0	Never-married	Adm-clerical	Not-in-family	White	Mal

In [2]:

```
n_records = data.shape[0]
```

```
n_greater_50k = data[data["income"] == ">50K"].shape[0]
```

```
n_at_most_50k = data[data["income"] == "<=50K"].shape[0]
```

```
greater_percent = n_greater_50k/n_records*100.0
```

```
print("Total number of records: {}".format(n_records))
```

```
print("Individuals making more than $50,000: {}".format(n_greater_50k))
```

```
print("Individuals making at most $50,000: {}".format(n_at_most_50k))
```

```
print("Percentage of individuals making more than $50,000: {:.2f}%".format(greater_percent))
```

```
print("Feature value for each column:\n", data.columns)
```

```
Total number of records: 45222
```

```
Individuals making more than $50,000: 11208
```

```
Individuals making at most $50,000: 34014
```

```
Percentage of individuals making more than $50,000: 24.78%
```

```
Feature value for each column:
```

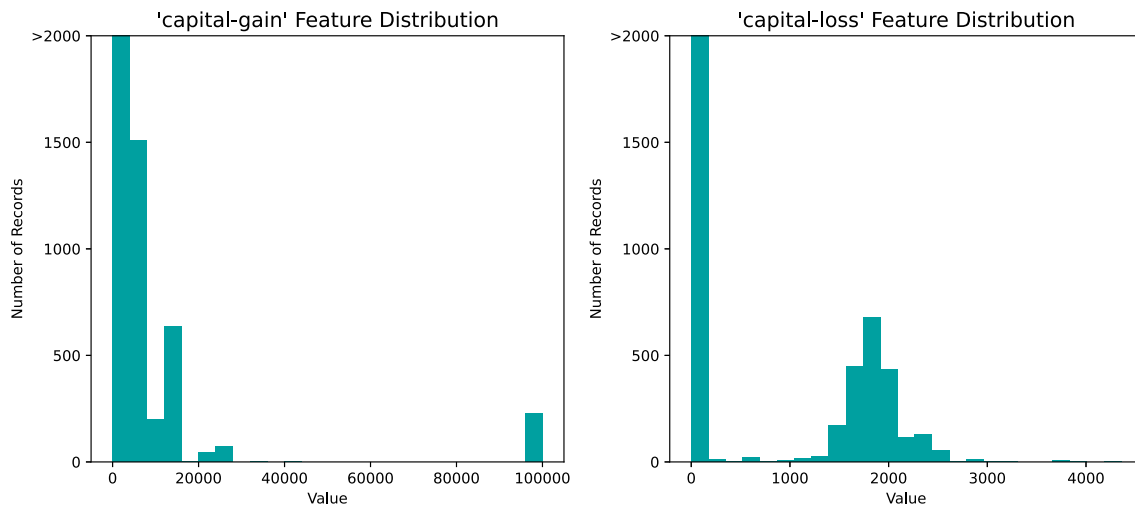
```
Index(['age', 'workclass', 'education_level', 'education-num',
       'marital-status', 'occupation', 'relationship', 'race', 'sex',
       'capital-gain', 'capital-loss', 'hours-per-week', 'native-country',
       'income'],
      dtype='object')
```

In [3]:

```
income_raw = data['income']
features_raw = data.drop('income', axis = 1)

vs.distribution(data)
```

Skewed Distributions of Continuous Census Data Features

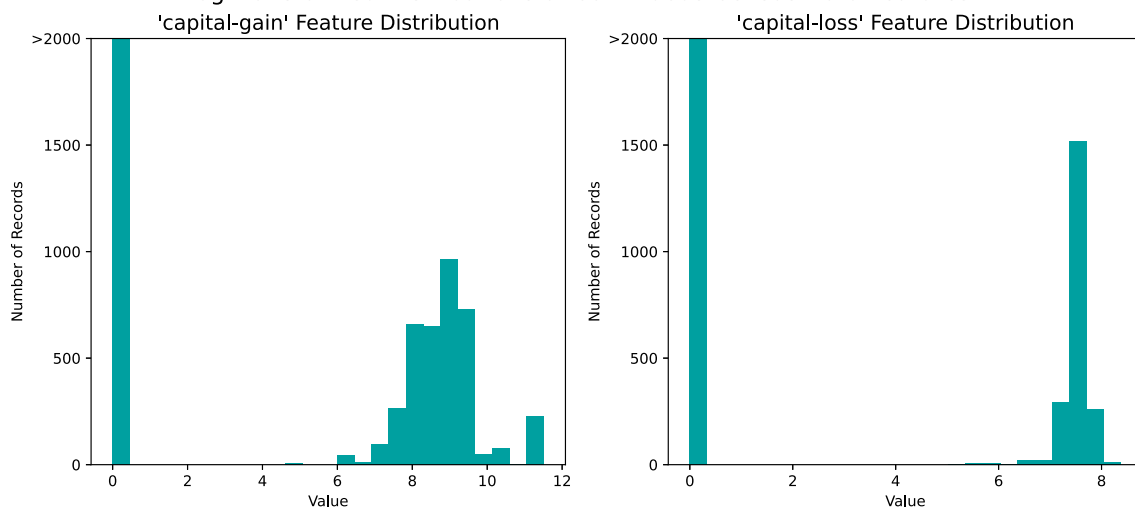


In [4]:

```
skewed = ['capital-gain', 'capital-loss']
features_log_transformed = pd.DataFrame(data = features_raw)
features_log_transformed[skewed] = features_raw[skewed].apply(lambda x: np.log(x + 1))

vs.distribution(features_log_transformed, transformed = True)
```

Log-transformed Distributions of Continuous Census Data Features



In [5]:

```
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
numerical = ['age', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']

features_log_minmax_transform = pd.DataFrame(data = features_log_transformed)
features_log_minmax_transform[numerical] = scaler.fit_transform(features_log_transformed[numerical])

display(features_log_minmax_transform.head(n = 5))
```

	age	workclass	education_level	education-num	marital-status	occupation	relationship	race
0	0.301370	State-gov	Bachelors	0.800000	Never-married	Adm-clerical	Not-in-family	White
1	0.452055	Self-emp-not-inc	Bachelors	0.800000	Married-civ-spouse	Exec-managerial	Husband	White
2	0.287671	Private	HS-grad	0.533333	Divorced	Handlers-cleaners	Not-in-family	White
3	0.493151	Private	11th	0.400000	Married-civ-spouse	Handlers-cleaners	Husband	Black
4	0.150685	Private	Bachelors	0.800000	Married-civ-spouse	Prof-specialty	Wife	Black

In [6]:

```
data2 =pd.get_dummies(features_log_minmax_transform)

income = income_raw.replace(["<=50K", ">50K"], [0, 1])

encoded = list(data2.columns)
print("{} total features ".format(len(encoded)))

print (encoded)
```

103 total features

```
['age', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week', 'workclass_Federal-gov', 'workclass_Local-gov', 'workclass_Private', 'workclass_Self-emp-inc', 'workclass_Self-emp-not-inc', 'workclass_State-gov', 'workclass_Without-pay', 'education_level_10th', 'education_level_11th', 'education_level_12th', 'education_level_1st-4th', 'education_level_5th-6th', 'education_level_7th-8th', 'education_level_9th', 'education_level_Assoc-acdm', 'education_level_Assoc-voc', 'education_level_Bachelors', 'education_level_Doctorate', 'education_level_HS-grad', 'education_level_Masters', 'education_level_Preschool', 'education_level_Prof-school', 'education_level_Some-college', 'marital-status_Divorced', 'marital-status_Married-AF-spouse', 'marital-status_Married-civ-spouse', 'marital-status_Married-spouse-absent', 'marital-status_Never-married', 'marital-status_Separated', 'marital-status_Widowed', 'occupation_Adm-clerical', 'occupation_Armed-Forces', 'occupation_Craft-repair', 'occupation_Exec-managerial', 'occupation_Farming-fishing', 'occupation_Handlers-cleaners', 'occupation_Machine-op-inspct', 'occupation_Other-service', 'occupation_Priv-house-serv', 'occupation_Prof-specialty', 'occupation_Protective-serv', 'occupation_Sales', 'occupation_Tech-support', 'occupation_Transport-moving', 'relationship_Husband', 'relationship_Not-in-family', 'relationship_Other-relative', 'relationship_Own-child', 'relationship_Unmarried', 'relationship_Wife', 'race_Amer-Indian-Eskimo', 'race_Asian-Pac-Islander', 'race_Black', 'race_Other', 'race_White', 'sex_Female', 'sex_Male', 'native-country_Cambodia', 'native-country_Canada', 'native-country_China', 'native-country_Columbia', 'native-country_Cuba', 'native-country_Dominican-Republic', 'native-country_Ecuador', 'native-country_El-Salvador', 'native-country_England', 'native-country_France', 'native-country_Germany', 'native-country_Greece', 'native-country_Guatemala', 'native-country_Haiti', 'native-country_Holand-Netherlands', 'native-country_Honduras', 'native-country_Hong', 'native-country_Hungary', 'native-country_India', 'native-country_Iran', 'native-country_Ireland', 'native-country_Italy', 'native-country_Jamaica', 'native-country_Japan', 'native-country_Laos', 'native-country_Mexico', 'native-country_Nicaragua', 'native-country_Outlying-US(Guam-USVI-etc)', 'native-country_Peru', 'native-country_Philippines', 'native-country_Poland', 'native-country_Portugal', 'native-country_Puerto-Rico', 'native-country_Scotland', 'native-country_South', 'native-country_Taiwan', 'native-country_Thailand', 'native-country_Trinidad&Tobago', 'native-country_United-States', 'native-country_Vietnam', 'native-country_Yugoslavia']
```

In [7]:

```
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test = train_test_split(data2,income,test_size = 0.2, random_state = 0)

print("Training set has {} samples".format(X_train.shape[0]))
print("Testing set has {} samples".format(X_test.shape[0]))
```

Training set has 36177 samples
Testing set has 9045 samples

In [8]:

```
print('exercise2:')
TP = np.sum(income)
FP = income.count() - TP
TN = 0
FN = 0
accuracy = float(TP)/(TP+FP)
recall = float(TP)/(TP+FN)
precision = accuracy

beta=0.5
fscore = (1 + beta ** 2)*(precision * recall)/(beta ** 2 *precision + recall)
TPR = float(TP)/(TP + FN)
FPR = float(FP)/(TN + FP)
print("Accuracy score: {:.4f}, F-score: {:.4f}".format(accuracy, fscore))
```

exercise2:

Accuracy score: 0.2478, F-score: 0.2917

In [9]:

```

from sklearn.metrics import fbeta_score, accuracy_score, confusion_matrix, roc_curve, auc
def train(learner, X_train, Y_train, X_test, Y_test):

    results = {}

    #fit/train
    # start = time()
    learner.fit(X_train, Y_train)
    # end = time()
    # results['train_time'] = end - start

    #predict
    # start = time()
    Y_pred_test = learner.predict(X_test)
    Y_pred_train = learner.predict(X_train)
    # end = time()
    # results['pred_time'] = end-start

    # results['acc_train_score'] = accuracy_score(Y_train, Y_pred_train)

    # results['acc_test_score'] = accuracy_score(Y_test, Y_pred_test)

    # results['train_f_score'] = fbeta_score(Y_train, Y_pred_train, beta=0.5)

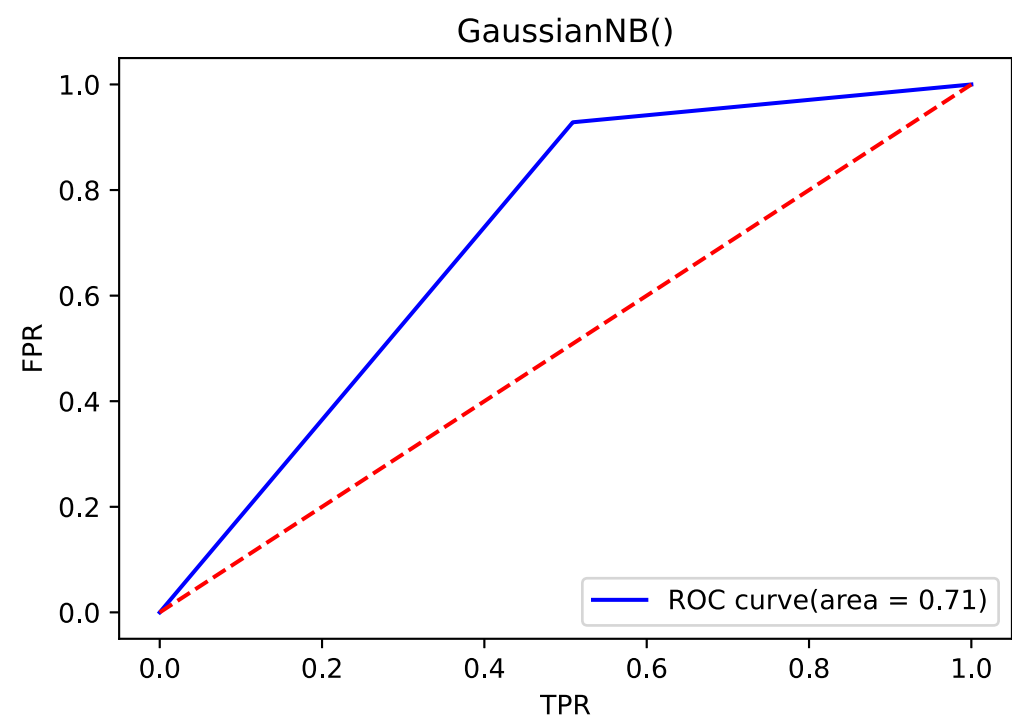
    # results['test_f_score'] = fbeta_score(Y_test, Y_pred_test, beta=0.5)

    #draw ROC
    # print(learner.name)
    fpr, tpr, threshold = roc_curve(Y_test, Y_pred_test)
    print('The function {}\'s fpr is {} and ptr is {}'.format(learner, fpr, tpr))
    roc_auc = auc(fpr, tpr)
    plt.figure()
    plt.title(learner)
    plt.plot(fpr, tpr, 'b', label = 'ROC curve(area = %.2f)' % roc_auc)
    plt.legend(loc = "lower right")
    plt.plot([0, 1], [0, 1], 'r--', label = 'random')
    plt.xlabel('TPR')
    plt.ylabel('FPR')
    plt.show()

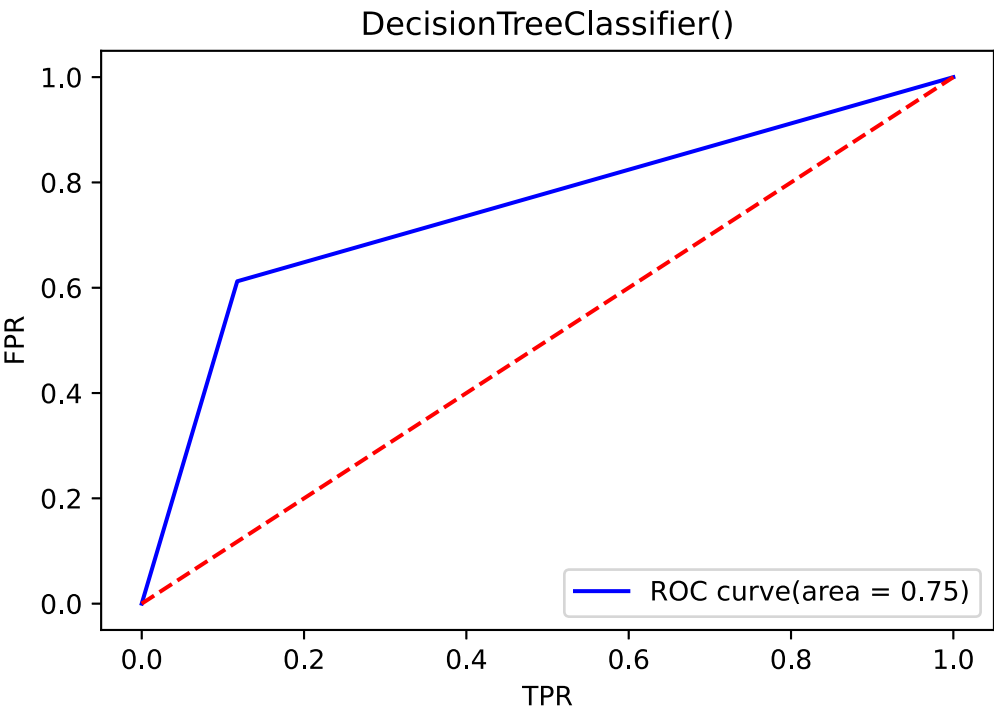
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
learner = [GaussianNB(), DecisionTreeClassifier(), BaggingClassifier(), AdaBoostClassifier(), Random
ForestClassifier(), KNeighborsClassifier(), SVC(), LogisticRegression()]
for i in learner:
    train(i, X_train, Y_train, X_test, Y_test)

```

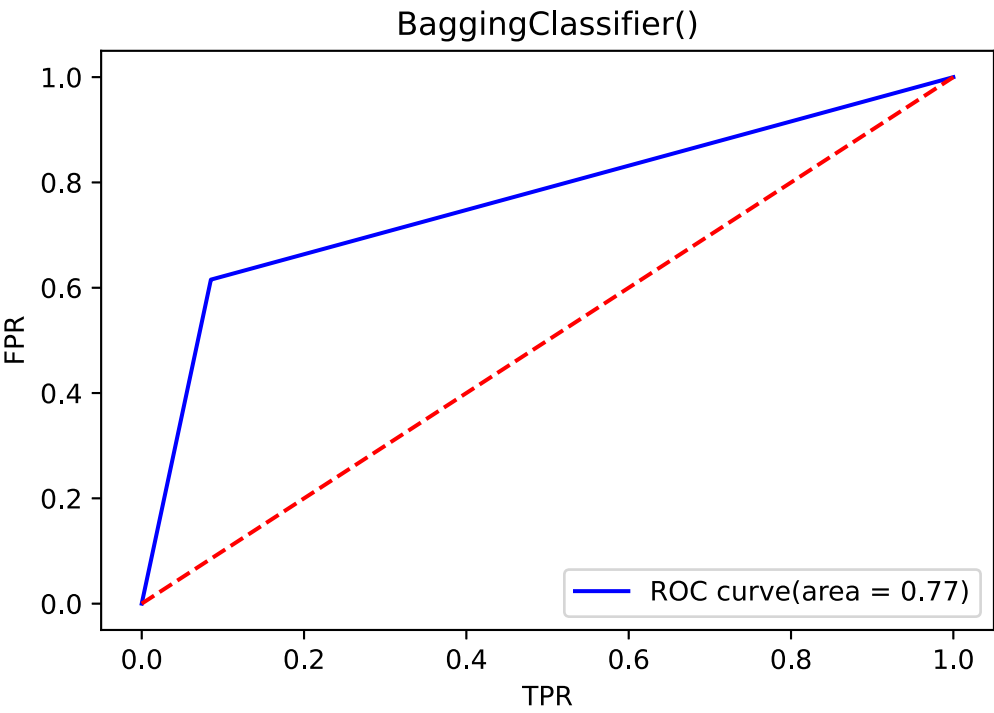
The function GaussianNB()'s fpr is [0. 0.50891813 1.] and ptr is [0. 0.92834467 1.]



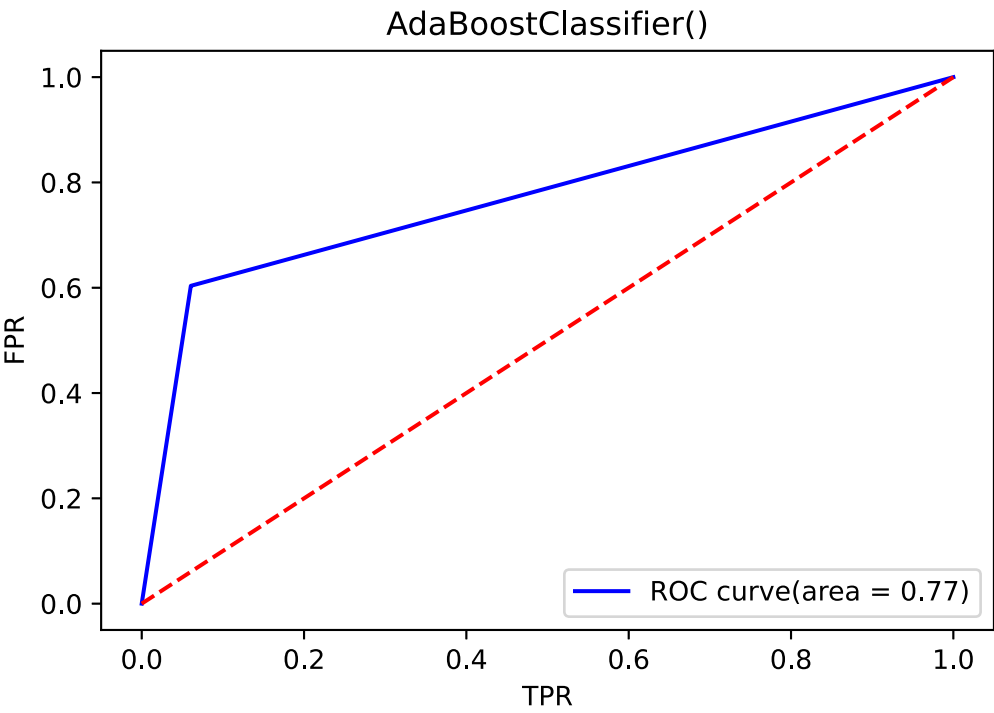
The function DecisionTreeClassifier()'s fpr is [0. 0.11769006 1.] and ptr is [0. 0.6122449 1.]



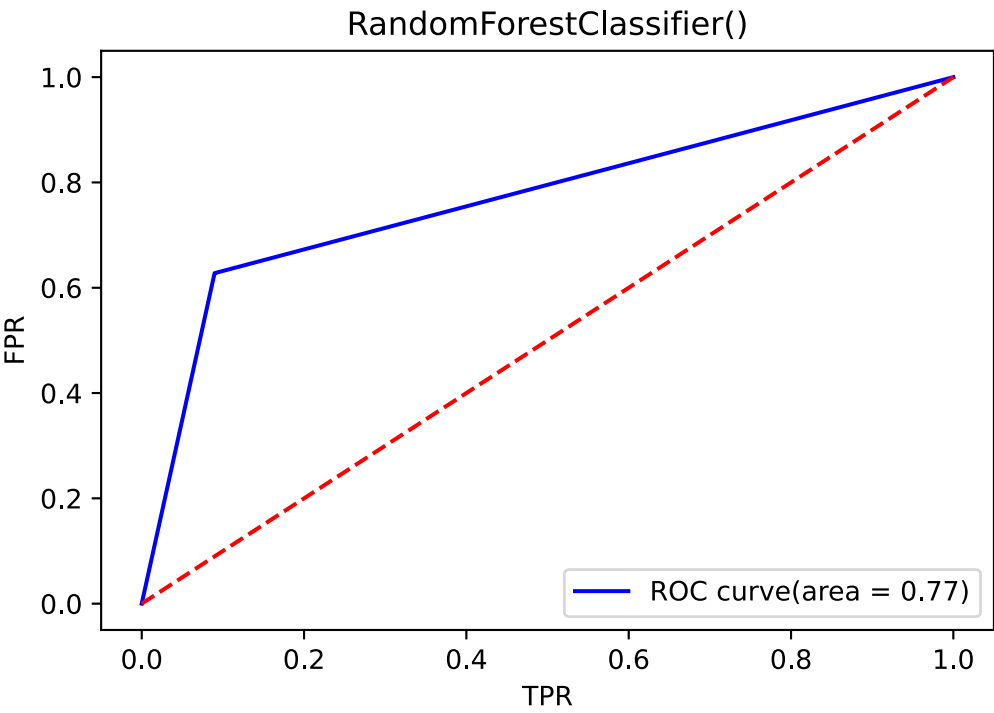
The function BaggingClassifier()'s fpr is [0.08523392 1.] and p
tr is[0.6154195 1.]



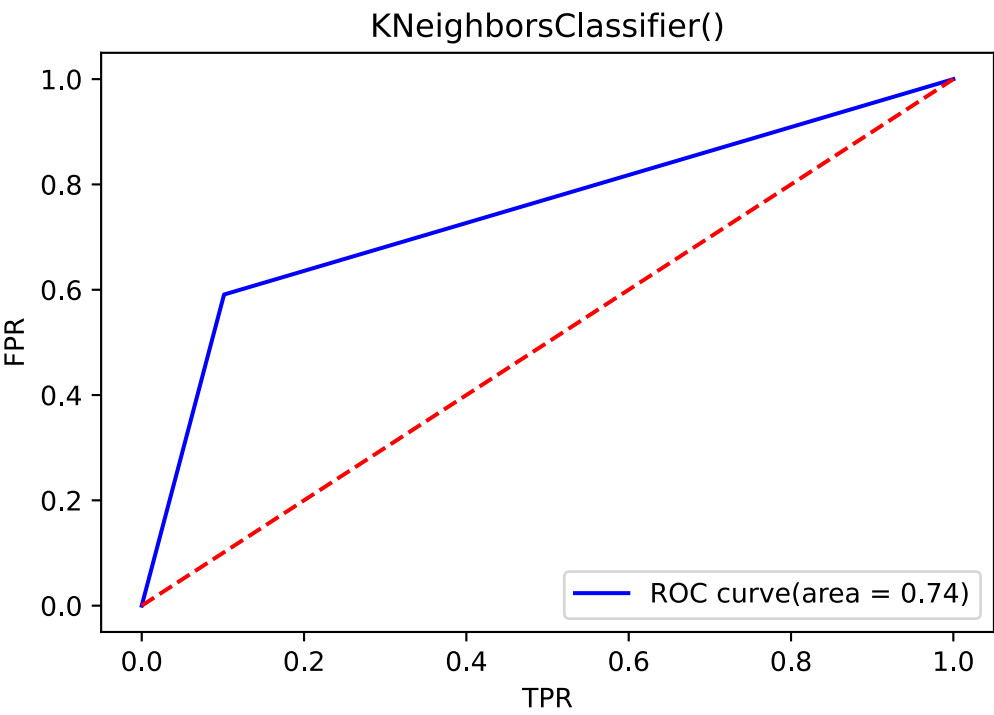
The function AdaBoostClassifier()'s fpr is [0.06052632 1.] and ptr is[0.60362812 1.]



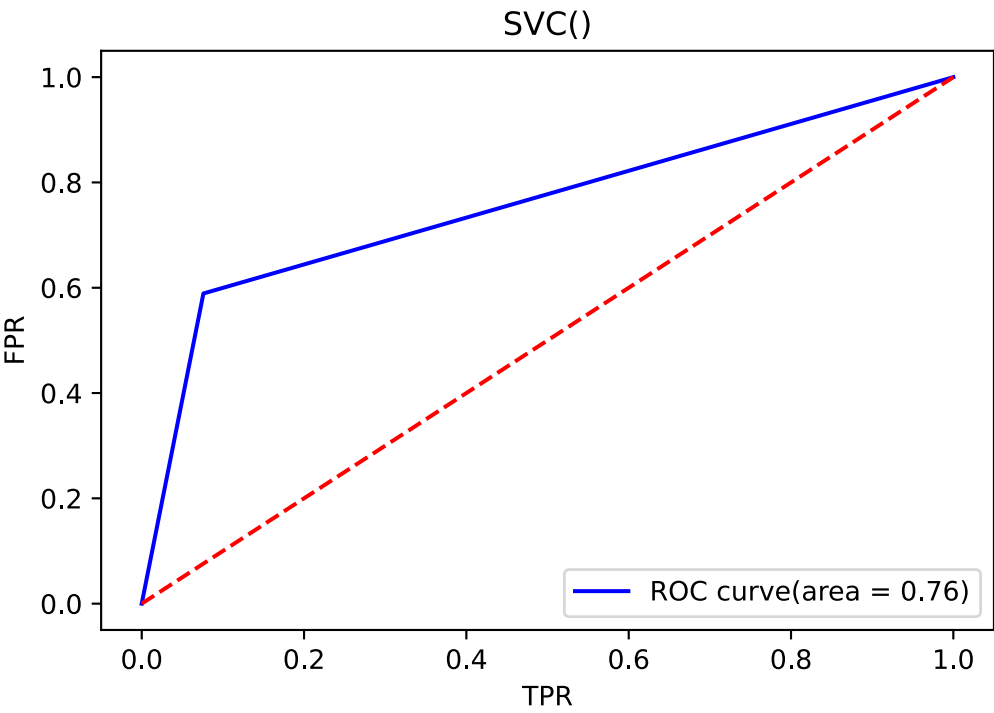
The function RandomForestClassifier()'s fpr is [0.08976608 1.0] and ptr is [0.6276644 1.0]



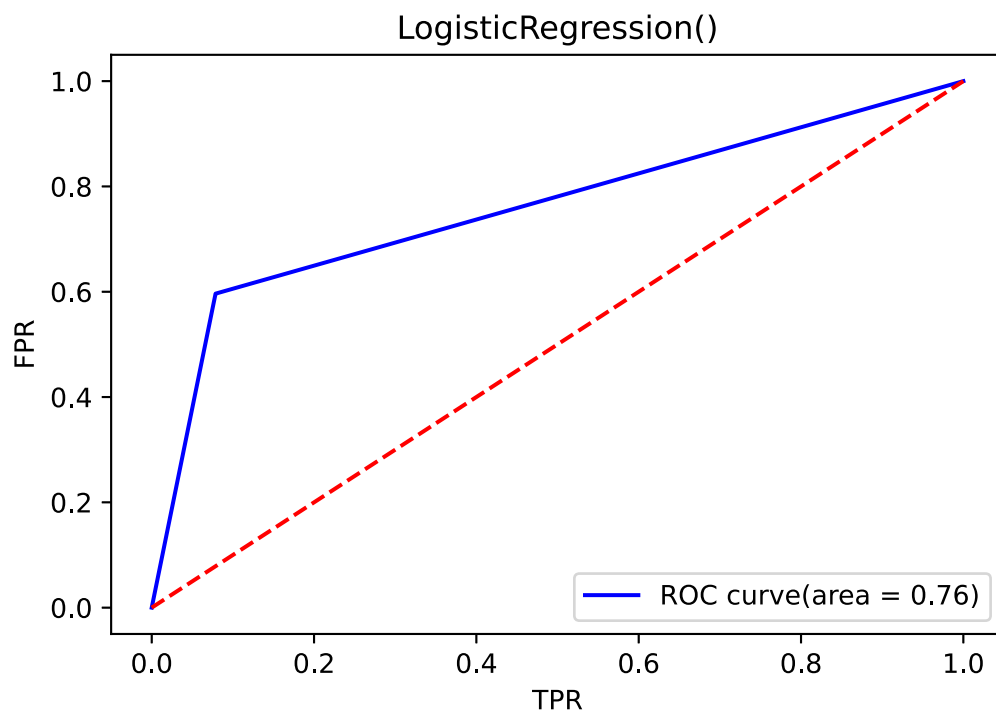
The function KNeighborsClassifier()'s fpr is [0. 0.10146199 1.] and ptr is [0. 0.59092971 1.]



The function SVC()'s fpr is [0.07602339 1.0] and ptr is [0.58911565 1.0]



The function `LogisticRegression()`'s `fpr` is `[0.07865497 1.]` and `ptr` is `[0.59637188 1.]`



In [16]:

```
RFC = RandomForestClassifier()  
RFC.fit(X_train, Y_train)
```

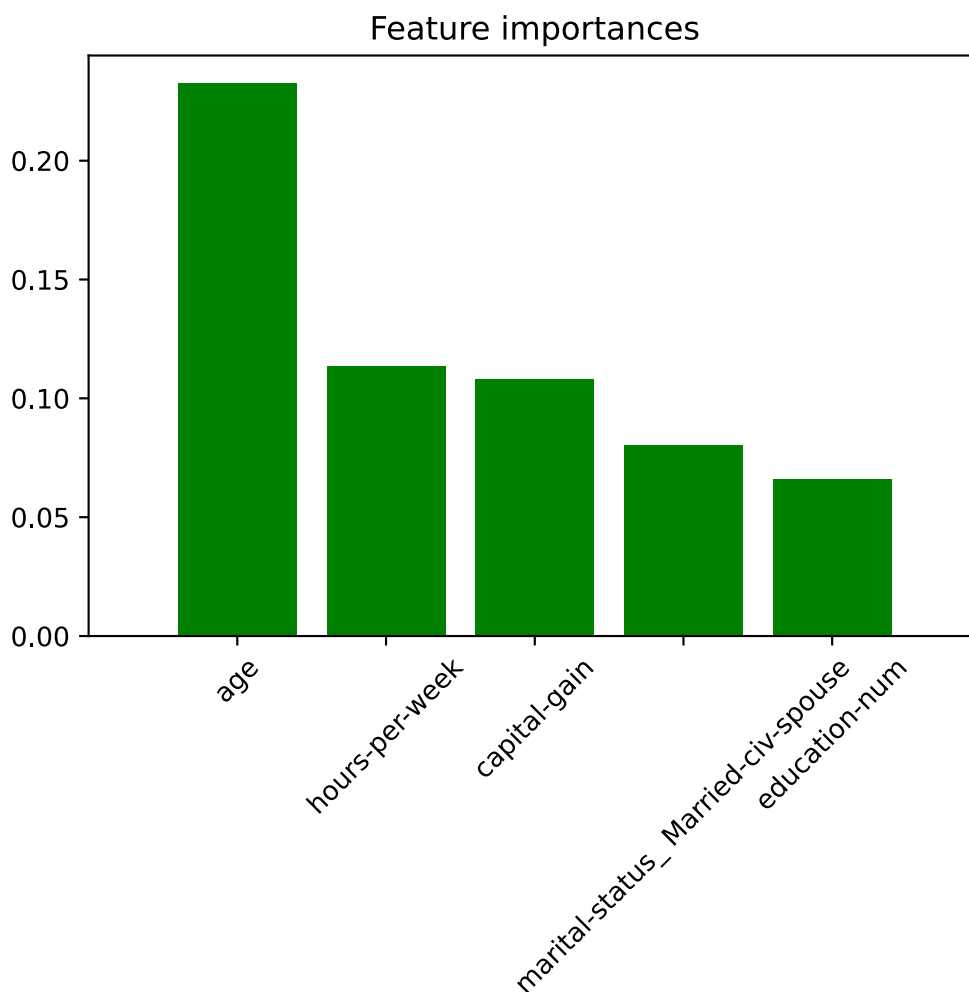
Out[16]:

```
RandomForestClassifier()
```

In [24]:

```
Y_pred_test = RFC.predict(X_test)
Y_pred_train = RFC.predict(X_train)
score = RFC.score(X_test, Y_test)
features = list(X_test.columns)
importances = RFC.feature_importances_
indices = np.argsort(importances)[::-1]
# print top 5 important features
num_features = 5
# num_features = len(importances)

plt.figure()
plt.title("Feature importances")
plt.bar(range(num_features), importances[indices[0:num_features]], color="g", align="center")
plt.xticks(range(num_features), [features[i] for i in indices], rotation='45')
plt.xlim([-1, num_features])
plt.show()
```



Exercise 3:

通过上面方法的 ROC 图我们观察到 Ensemble Methods 有最好的分类结果。

- (1) 我们选取了 Decision Trees、AdaBoost、RandomForest 方法并画出了各个方法排名前五的特征，如上图所示。
- (2) 土地覆被测绘是地球观测卫星传感器的主要应用之一，它利用遥感和地理空间数据来识别位于目标区域表面的材料和物体。通常，目标材料的类别包括道路，建筑物，河流，湖泊和植被。基于人工神经网络的一些不同的集成学习方法，带 Boosting 的决策树，随机森林和多分类器系统的自动设计，被提出来有效地识别土地覆盖物。
- (3) 随机森林可以处理很高维度的数据，并且不用降维，无需做特征选择，且可以判断特征的重要程度，训练速度较快
- (4) 随机森林算法在某些噪音较大的分类或回归上会过拟合，对于不同取值属性的数据，取值划分较多的属性会对随机森林产生更大的影响。
- (5) 在这个数据集中，数据的维度比较高，所以使用随机森林算法可以比较快速的完成拟合，在随机森林算法处理完可以从中得出重要特征，更利于分析这个问题得出收入高低的主要影响因素。