

Theoretical homework #1

Kamil Salakhiev

April 12, 2016

1 BEST MODEL DESCRIPTION

1.1 DATA PREPARATION

Data was written using pandas library:

```
train = pd.read_csv("train.csv", sep=',', header = 0)
```

To extract features and categories appropriate methods `get_features()` and `get_categories` were prepared:

```
X = get_features(train)
y = get_categories(train)
```

```
def get_features(df, encoder=None):
    cat=['workclass', 'education', 'marital-status', 'occupation', \
         'relationship', 'race', 'sex', 'native-country']
    num=['age', 'capital-gain', 'capital-loss', 'hours-per-week']
    if (len(cat) == 0 or len(num) == 0):
        return 'con'
    if encoder is None:
        Xseries_categorical = [get_one_hot_encoding(df[x]) \
                               for x in cat]
    else:
        Xseries_categorical = [get_one_hot_encoding_with \
                               (df[x], encoder) for x in cat]
    Xseries_numerical = [df[x] for x in num]
```

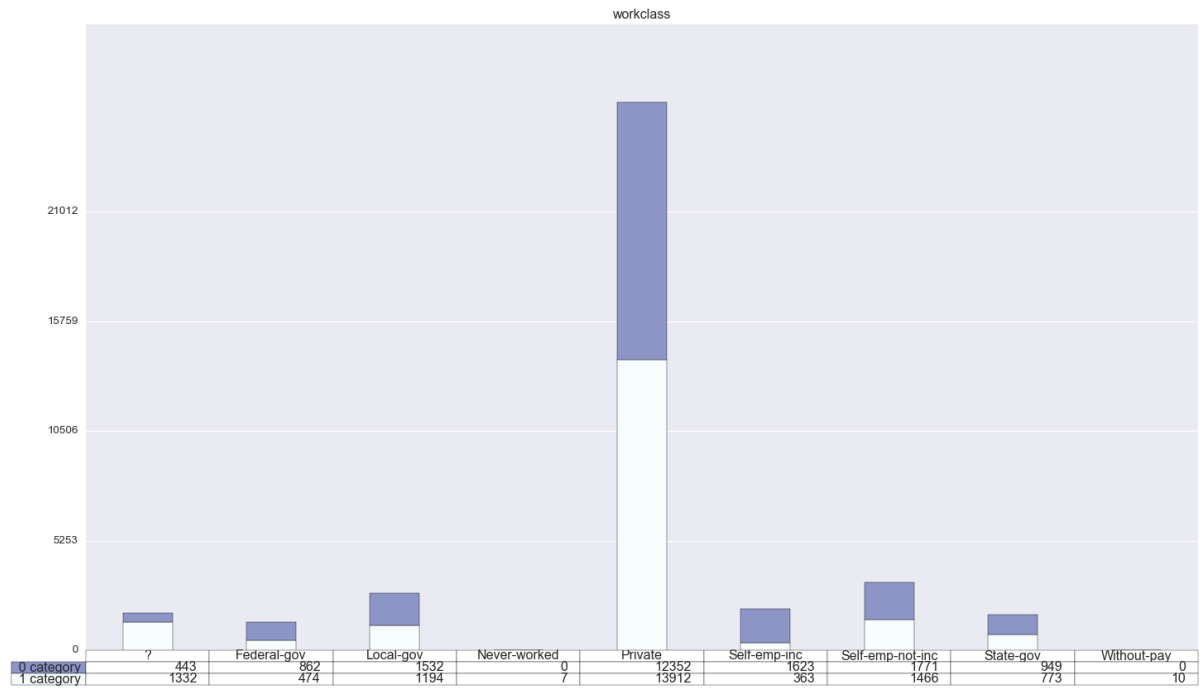
```
        X = pd.concat(Xseries_categorical+Xseries_numerical , axis=1)
        return X.values
def get_categories(df):
    return df['Category'].values
```

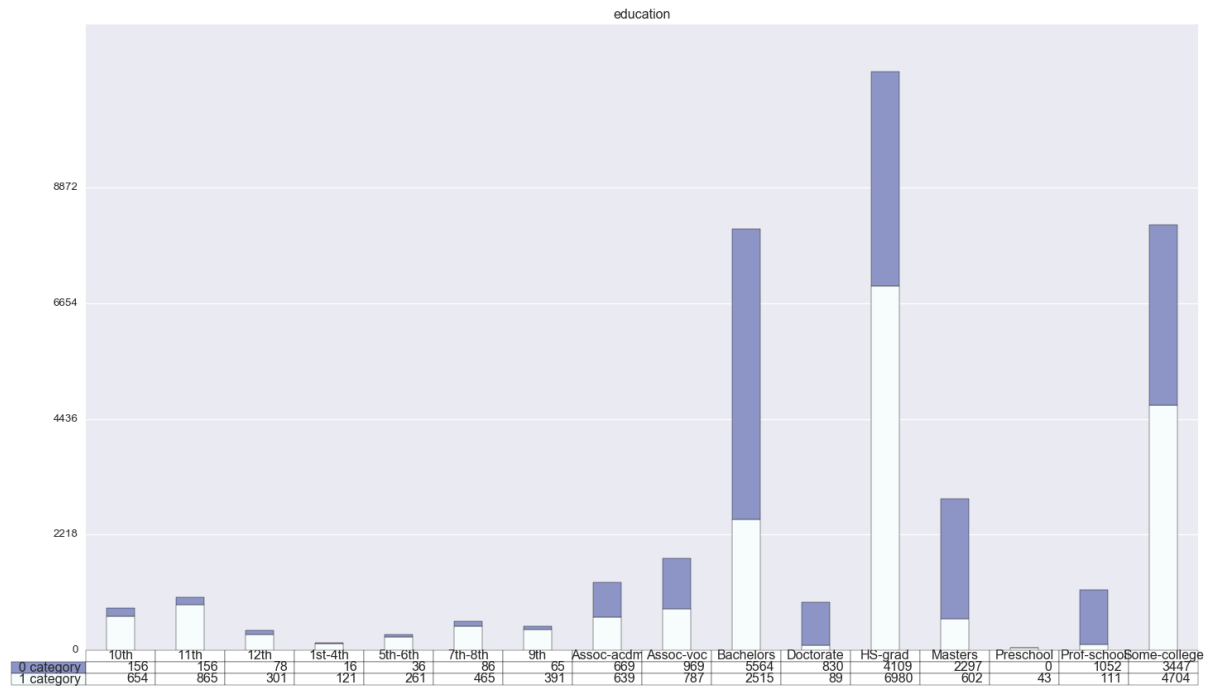
To select features, lists cat and num were varied in order obtain the best set. However, eventually, sklearn library were charged for feature selection. Recursive feature elimination were used for this goal, which recursively considered smaller and smaller sets of features. The quality of features' sets were checked by AdaBoost estimator with Decision tree with maximum depth equal one. Before AdaBoost estimator were chosen, plenty of other classifiers were considered(see later in report). But AdaBoost with recursive feature elimination gave the best performance. To end up with number of estimators = 900, and SAMME.R algorithm GridSearchCV were used. Calculated cross validation score was 0.870

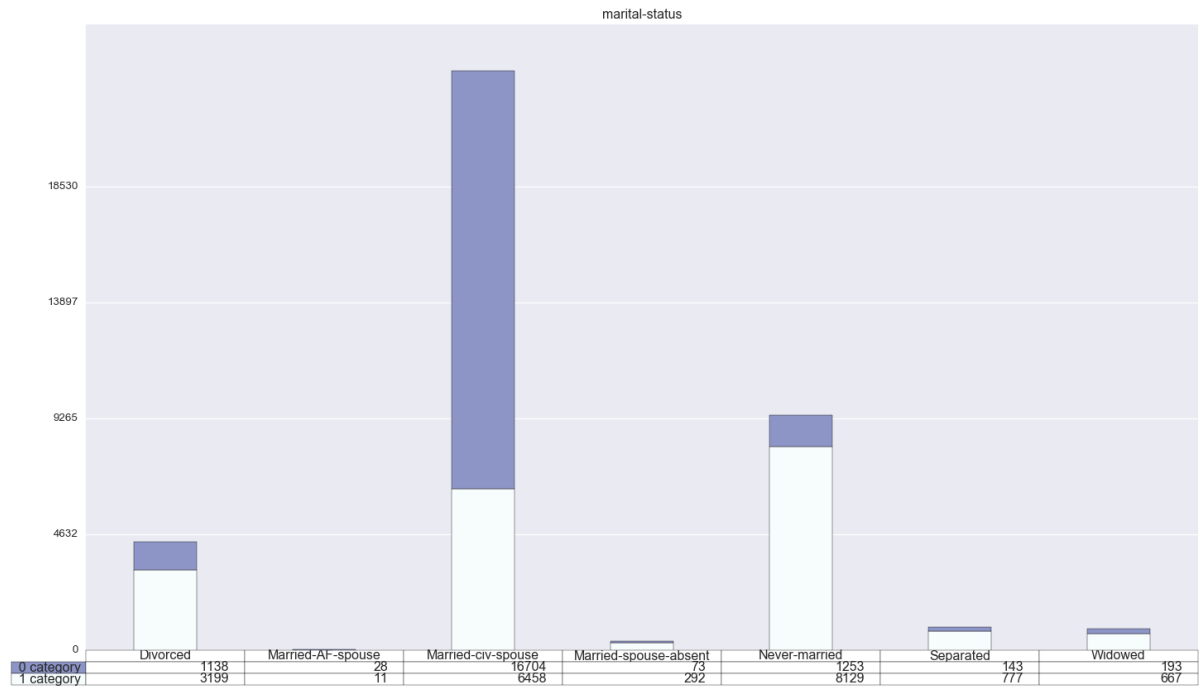
```
ada=AdaBoostClassifier(
    base_estimator=DecisionTreeClassifier(max_depth=1),
    algorithm='SAMME.R',n_estimators=900)
rfecv = RFECV(estimator=ada, step=1, scoring='accuracy',verbose=1)
```

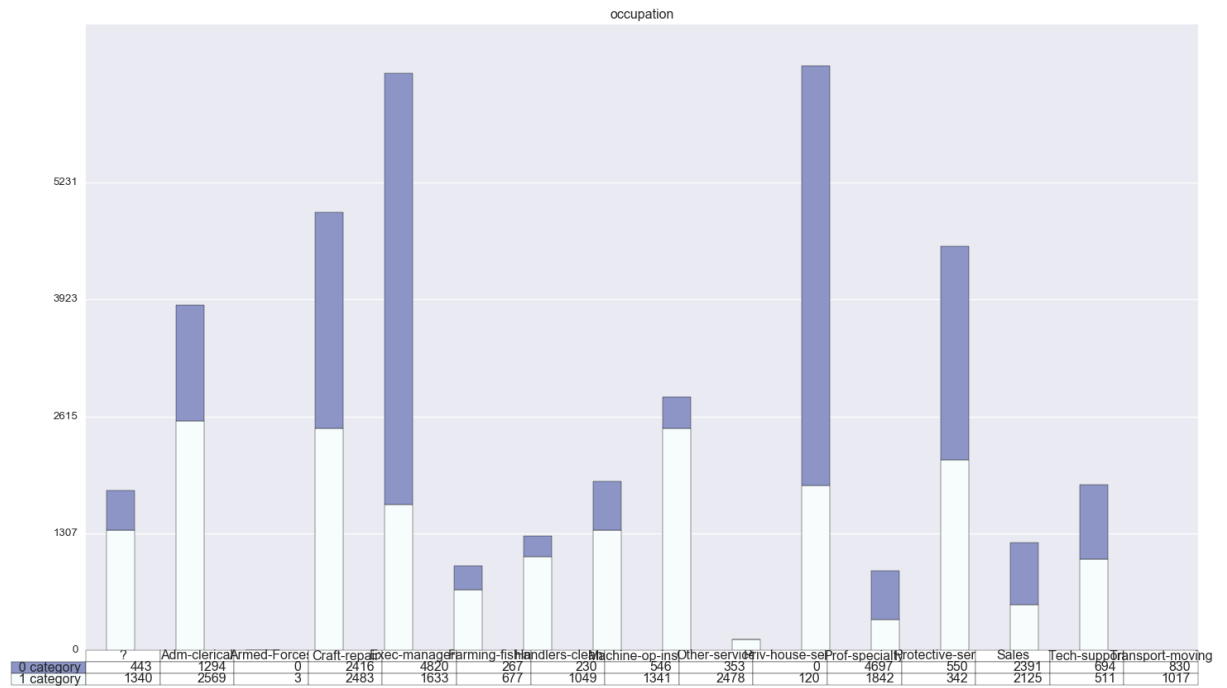
2 GRAPHS

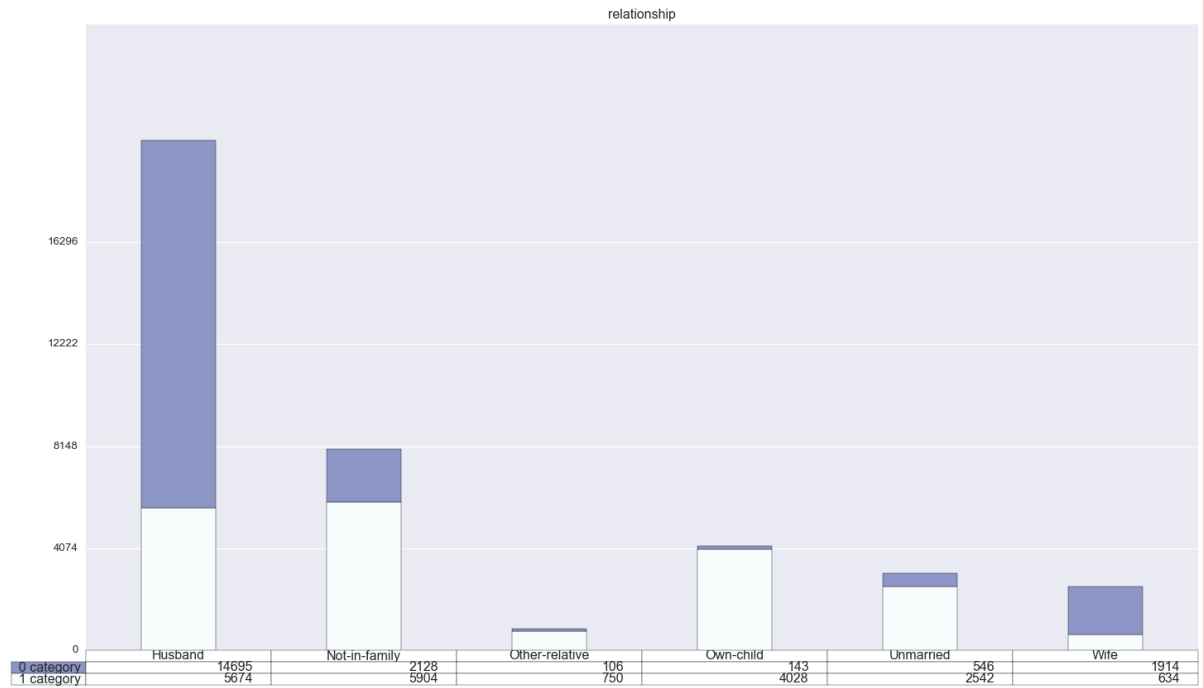
Before I ended up with AdaBoost classifier, many other classifiers were used. To select appropriate features, dataset were visualized. All used graphics are placed below.

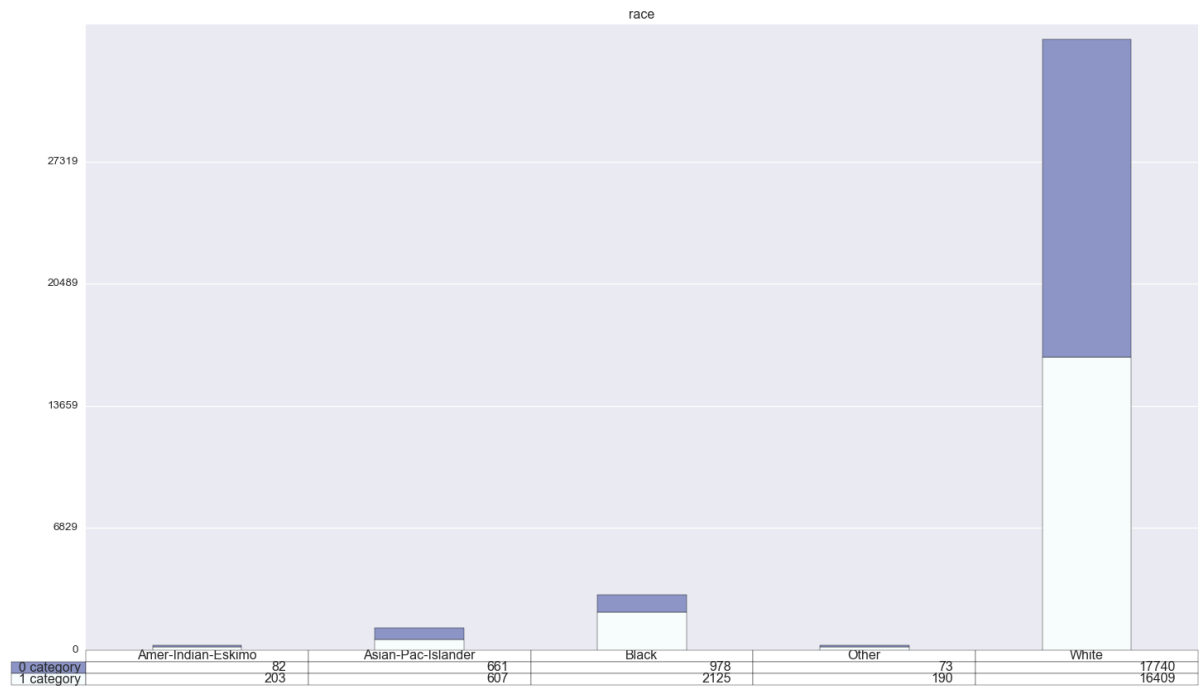


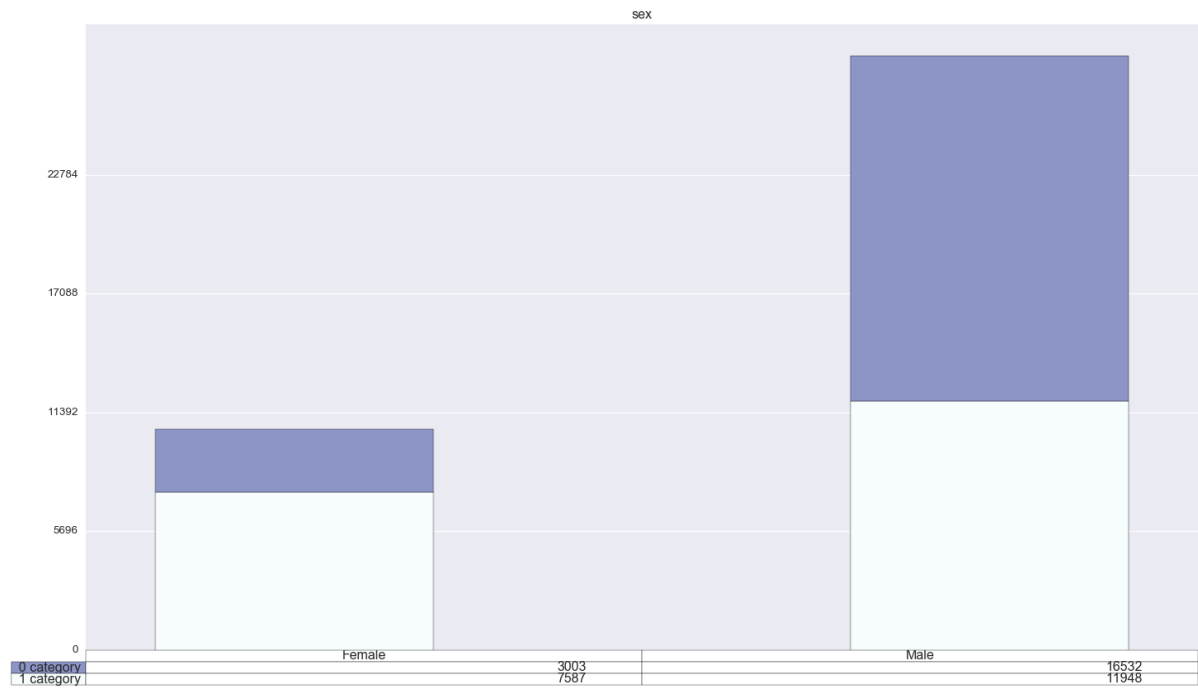














Based on graphics, I decided to remain the following set of features: ['sex', 'education', 'native-country', 'marital-status', 'relationship', 'occupation', 'capital-gain', 'capital-loss']. Other features showed less certainty. For example workclass feature shows, that based on this feature it is hard to decide to which category person belongs, while marital status demonstrates good reliability. For example if we have 'never-married' person we can see that this person most likely belongs to the first category.

3 OTHER MODELS

The following models were obtained with gridsearch (aforementioned features were supplied to the models):

- SVM with RBF kernel and gamma=0.1 with PCA dimensionality reduction to the 40 components. Obtained accuracy:0.862
- Decision tree with PCA dimensionality reduction to 2 components. Obtained accuracy: 0.862
- Random forest classifier with PCA dimensionality reduction to 3 components with max depth 11 and with 15 estimators. Obtained accuracy: 0.863

- KNN classifier with PCA dimensionality reduction to 20 components with 13 nearest neighbors and power distance parameter $p = 1$. Obtained accuracy: 0.860
- Bagging with KNN with 5 nearest neighbors and $p=1$. Obtained accuracy: 0.863

After deriving several strong classifiers, they were put into 1 voting classifier. As result I got improving to 0.866 accuracy.

However, again, AdaBoost beat that result with accuracy on cross validation 0.870. Latter model were used for make submission on kaggle with accuracy 0.86809.

Classifier	Accuracy on CV
SVM	0.862
Decision tree	0.862
Random forest	0.863
KNN	0.860
Bagging on KNN	0.863
Voting Classifier	0.866
AdaBoost	0.870

Table 1: Classifiers' accuracy