

# Деректерді өңдеу мамандарына арналған практикалық статистика (Python)

## 6-тарау. Статистикалық Машиналық Оқыту

(c) 2019 Peter C. Bruce, Andrew Bruce, Peter Gedeck

Import required Python packages.

In [1]:

```
import math
import os
import random
from pathlib import Path
from collections import defaultdict
from itertools import product

import pandas as pd
import numpy as np

from sklearn import preprocessing
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

from xgboost import XGBClassifier

from dmba import plotDecisionTree, textDecisionTree

import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib.patches import Ellipse
```

```
%matplotlib inline
```

```
/opt/conda/lib/python3.9/site-packages/xgboost/compat.py:36: FutureWarning: pandas.Int64Index is deprecated and will be removed from pandas in a future version. Use pandas.Index with the appropriate dtype instead.
```

```
from pandas import MultiIndex, Int64Index
no display found. Using non-interactive Agg backend
```

In [2]:

```
try:
    import common
    DATA = common.dataDirectory()
except ImportError:
    DATA = Path().resolve() / 'data'
```

Define paths to data sets. If you don't keep your data in the same directory as the code, adapt the path names.

In [3]:

```
LOAN200_CSV = DATA / 'loan200.csv'
```

```
LOAN3000_CSV = DATA / 'loan3000.csv'  
LOAN_DATA_CSV = DATA / 'loan_data.csv.gz'
```

Set this if the notebook crashes in the XGBoost part.

In [4]:

```
os.environ['KMP_DUPLICATE_LIB_OK'] = 'TRUE'
```

## K-Nearest Neighbors

### A Small Example: Predicting Loan Default

#### К-жақын көршілер

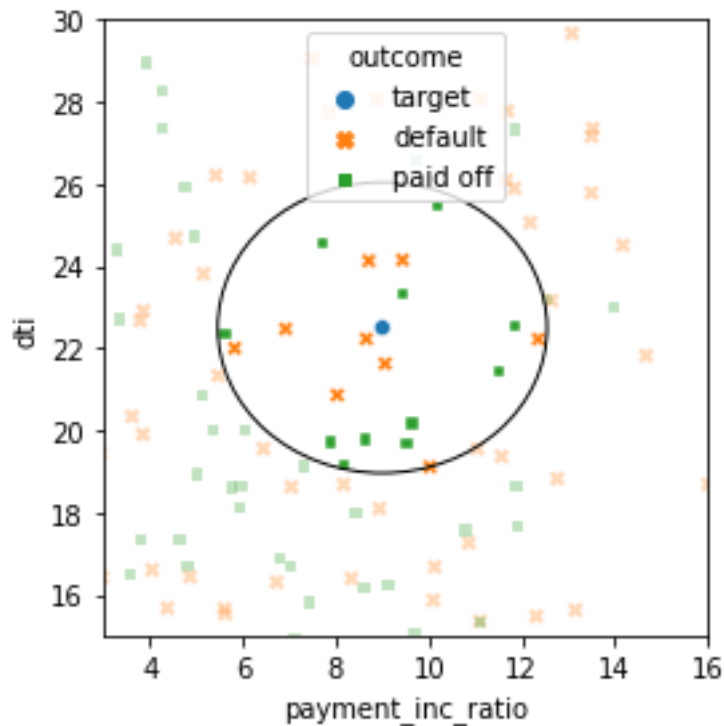
#### Шағын мысал: несие бойынша дефолтты болжау

In [5]:

```
loan200 = pd.read_csv(LOAN200_CSV)  
  
predictors = ['payment_inc_ratio', 'dti']  
outcome = 'outcome'  
  
newloan = loan200.loc[0:0, predictors]  
X = loan200.loc[1:, predictors]  
y = loan200.loc[1:, outcome]  
  
knn = KNeighborsClassifier(n_neighbors=20)  
knn.fit(X, y)  
knn.predict(newloan)  
print(knn.predict_proba(newloan))  
[[0.45 0.55]]
```

In [6]:

```
nbrs = knn.kneighbors(newloan)  
maxDistance = np.max(nbrs[0][0])  
  
fig, ax = plt.subplots(figsize=(4, 4))  
sns.scatterplot(x='payment_inc_ratio', y='dti', style='outcome',  
                hue='outcome', data=loan200, alpha=0.3, ax=ax)  
sns.scatterplot(x='payment_inc_ratio', y='dti', style='outcome',  
                hue='outcome',  
                data=pd.concat([loan200.loc[0:0, :], loan200.loc[nbrs[1][0] +  
1, :]]),  
                ax=ax, legend=False)  
ellipse = Ellipse(xy=newloan.values[0],  
                  width=2 * maxDistance, height=2 * maxDistance,  
                  edgecolor='black', fc=None, lw=1)  
ax.add_patch(ellipse)  
ax.set_xlim(3, 16)  
ax.set_ylim(15, 30)  
  
plt.tight_layout()  
plt.show()
```



## Standardization (Normalization, Z-Scores)

## Стандарттау (қалыпқа келтіру, Z-ұпайлар)

In [7]:

```
loan_data = pd.read_csv(LOAN_DATA_CSV)
loan_data = loan_data.drop(columns=['Unnamed: 0', 'status'])
loan_data['outcome'] = pd.Categorical(loan_data['outcome'],
                                     categories=['paid off', 'default'],
                                     ordered=True)
```

```
predictors = ['payment_inc_ratio', 'dti', 'revol_bal', 'revol_util']
outcome = 'outcome'
```

```
newloan = loan_data.loc[0:0, predictors]
print(newloan)
X = loan_data.loc[1:, predictors]
y = loan_data.loc[1:, outcome]
```

```
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X, y)
```

```
nbrs = knn.kneighbors(newloan)
print(X.iloc[nbrs[1][0], :])
```

	payment_inc_ratio	dti	revol_bal	revol_util
0	2.3932	1.0	1687	9.4
35536	1.47212	1.46	1686	10.0
33651	3.38178	6.37	1688	8.4
25863	2.36303	1.39	1691	3.5
42953	1.28160	7.14	1684	3.9
43599	4.12244	8.98	1684	7.2

In [8]:

```
newloan = loan_data.loc[0:0, predictors]
X = loan_data.loc[1:, predictors]
y = loan_data.loc[1:, outcome]

scaler = preprocessing.StandardScaler()
scaler.fit(X * 1.0)

X_std = scaler.transform(X * 1.0)
newloan_std = scaler.transform(newloan * 1.0)

knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_std, y)

nbrs = knn.kneighbors(newloan_std)
print(X.iloc[nbrs[1][0], :])
```

	payment_inc_ratio	dti	revol_bal	revol_util
2080	2.61091	1.03	1218	9.7
1438	2.34343	0.51	278	9.9
30215	2.71200	1.34	1075	8.5
28542	2.39760	0.74	2917	7.4
44737	2.34309	1.37	488	7.2

## KNN as a Feature Engine

### Функция қозғалтқышы ретінде белгілі

In [9]:

```
loan_data = pd.read_csv(LOAN_DATA_CSV)
loan_data = loan_data.drop(columns=['Unnamed: 0', 'status'])
loan_data['outcome'] = pd.Categorical(loan_data['outcome'],
                                     categories=['paid off', 'default'],
                                     ordered=True)

predictors = ['dti', 'revol_bal', 'revol_util', 'open_acc',
              'delinq_2yrs_zero', 'pub_rec_zero']
outcome = 'outcome'

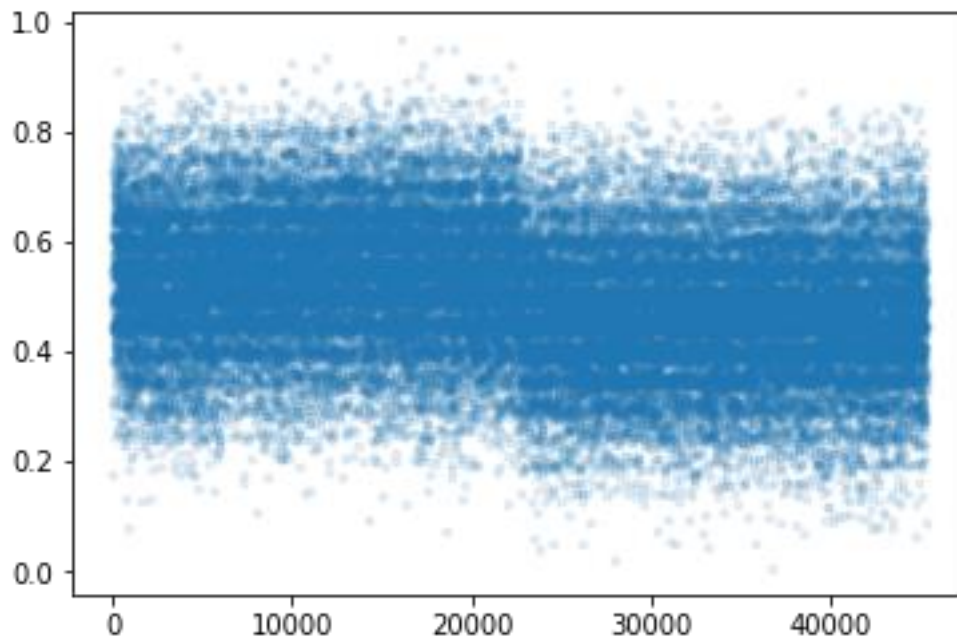
X = loan_data[predictors]
y = loan_data[outcome]

knn = KNeighborsClassifier(n_neighbors=20)
knn.fit(X, y)
plt.scatter(range(len(X)), [bs + random.gauss(0, 0.015) for bs in
                             knn.predict_proba(X[:,0]),
                             alpha=0.1, marker='.'])
knn.predict_proba(X[:, 0])

loan_data['borrower_score'] = knn.predict_proba(X[:, 0])
print(loan_data['borrower_score'].describe())
```

count	45342.000000
mean	0.501098
std	0.128736
min	0.000000
25%	0.400000
50%	0.500000

```
75%          0.600000
max          0.950000
Name: borrower_score, dtype: float64
```



## Tree Models

### A Simple Example

#### Үш модель

#### Қарапайым мысал

The package *scikit-learn* has the class `DecisionTreeClassifier` to build a decision tree model. The function `plotDecisionTree` from the *dmba* package can be used to visualize the tree.

In [10]:

```
loan3000 = pd.read_csv(LOAN3000_CSV)

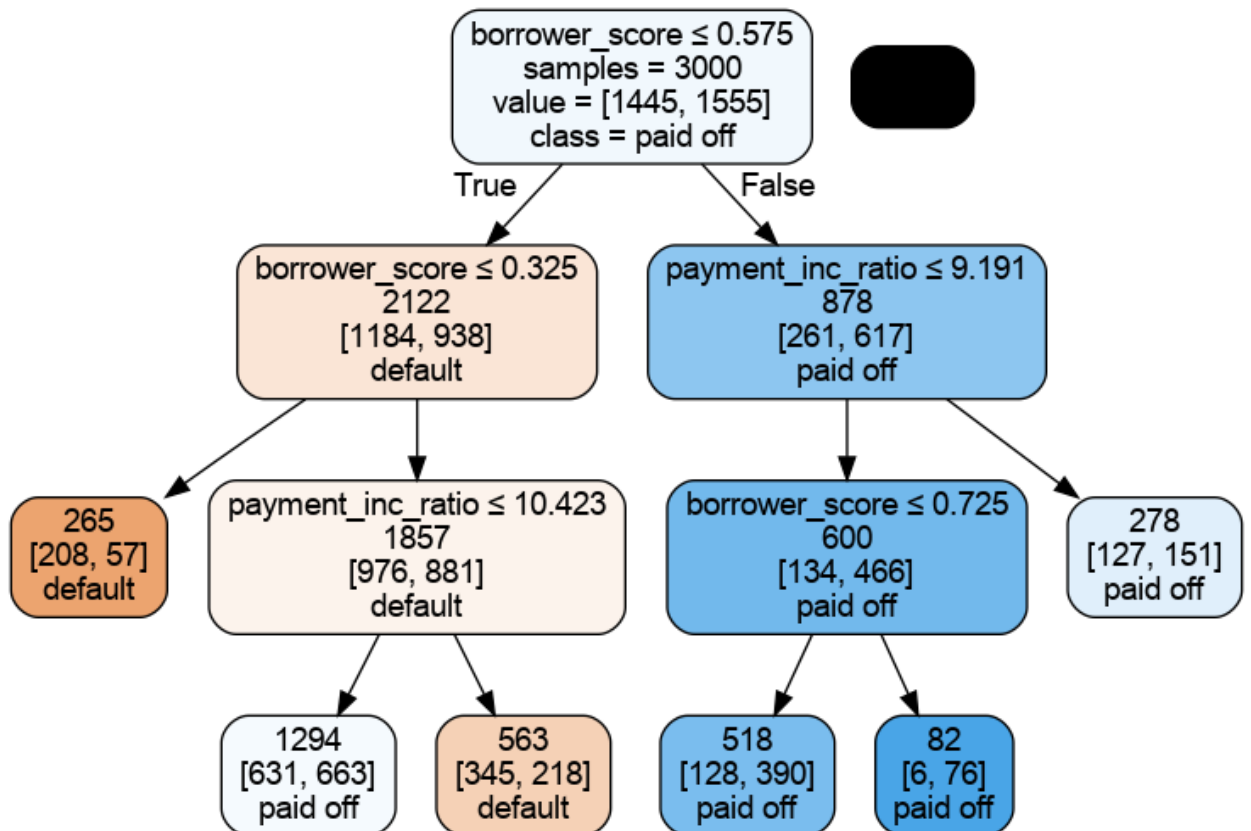
predictors = ['borrower_score', 'payment_inc_ratio']
outcome = 'outcome'

X = loan3000[predictors]
y = loan3000[outcome]

loan_tree = DecisionTreeClassifier(random_state=1, criterion='entropy',
                                   min_impurity_decrease=0.003)

loan_tree.fit(X, y)
plotDecisionTree(loan_tree, feature_names=predictors,
                 class_names=loan_tree.classes_)
```

Out[10]:



In [11]:

```

print(textDecisionTree(loan_tree))

node=0 test node: go to node 1 if 0 <= 0.5750000178813934 else to node 6
node=1 test node: go to node 2 if 0 <= 0.32500000298023224 else to node 3
node=2 leaf node: [[0.785, 0.215]]
node=3 test node: go to node 4 if 1 <= 10.42264986038208 else to node 5
node=4 leaf node: [[0.488, 0.512]]
node=5 leaf node: [[0.613, 0.387]]
node=6 test node: go to node 7 if 1 <= 9.19082498550415 else to node 10
node=7 test node: go to node 8 if 0 <= 0.7249999940395355 else to node 9
node=8 leaf node: [[0.247, 0.753]]
node=9 leaf node: [[0.073, 0.927]]
node=10 leaf node: [[0.457, 0.543]]
  
```

## The Recursive Partitioning Algorithm

## Рекурсивті бөлу алгоритмі

In [12]:

```

fig, ax = plt.subplots(figsize=(6, 4))

loan3000.loc[loan3000.outcome=='paid off'].plot(
    x='borrower_score', y='payment_inc_ratio', style='.',
    markerfacecolor='none', markeredgecolor='C1', ax=ax)
loan3000.loc[loan3000.outcome=='default'].plot(
    x='borrower_score', y='payment_inc_ratio', style='o',
    markerfacecolor='none', markeredgecolor='C0', ax=ax)
ax.legend(['paid off', 'default']);
ax.set_xlim(0, 1)
ax.set_ylim(0, 25)
ax.set_xlabel('borrower_score')
  
```

```

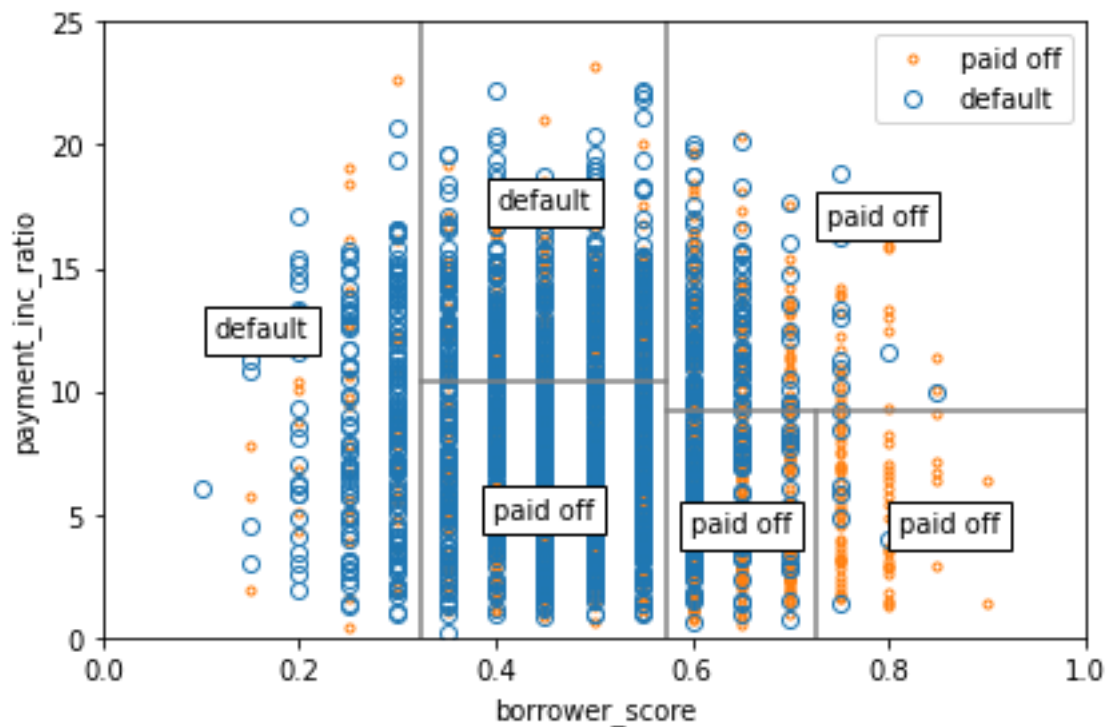
ax.set_ylabel('payment_inc_ratio')

x0 = 0.575
x1a = 0.325; y1b = 9.191
y2a = 10.423; x2b = 0.725
ax.plot((x0, x0), (0, 25), color='grey')
ax.plot((x1a, x1a), (0, 25), color='grey')
ax.plot((x0, 1), (y1b, y1b), color='grey')
ax.plot((x1a, x0), (y2a, y2a), color='grey')
ax.plot((x2b, x2b), (0, y1b), color='grey')

labels = [('default', (x1a / 2, 25 / 2)),
          ('default', ((x0 + x1a) / 2, (25 + y2a) / 2)),
          ('paid off', ((x0 + x1a) / 2, y2a / 2)),
          ('paid off', ((1 + x0) / 2, (y1b + 25) / 2)),
          ('paid off', ((1 + x2b) / 2, (y1b + 0) / 2)),
          ('paid off', ((x0 + x2b) / 2, (y1b + 0) / 2)),
          ]
for label, (x, y) in labels:
    ax.text(x, y, label, bbox={'facecolor':'white'},
            verticalalignment='center', horizontalalignment='center')

plt.tight_layout()
plt.show()

```



## Measuring Homogeneity or Impurity

## Біртектілікті немесе қоспаны өлшеу

```

def entropyFunction(x):
    if x == 0: return 0
    return -x * math.log(x, 2) - (1 - x) * math.log(1 - x, 2)

```

In [13]:

```

def giniFunction(x):
    return x * (1 - x)

x = np.linspace(0, 0.5, 50)
impure = pd.DataFrame({
    'x': x,
    'Accuracy': 2 * x,
    'Gini': [giniFunction(xi) / giniFunction(.5) for xi in x],
    'Entropy': [entropyFunction(xi) for xi in x],
})

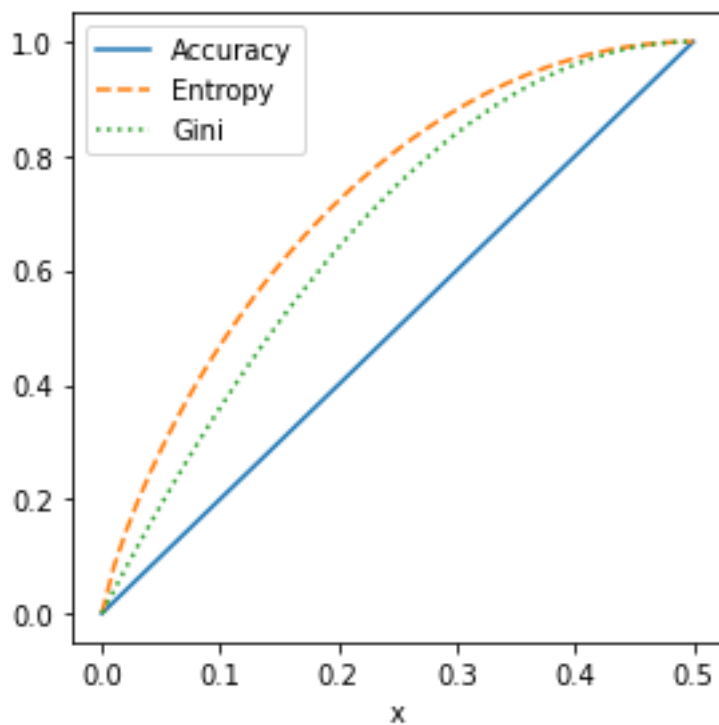
fig, ax = plt.subplots(figsize=(4, 4))

impure.plot(x='x', y='Accuracy', ax=ax, linestyle='solid')
impure.plot(x='x', y='Entropy', ax=ax, linestyle='--')
impure.plot(x='x', y='Gini', ax=ax, linestyle=':')

plt.tight_layout()
plt.show()

```

In [14]:



# Bagging and the Random Forest

## Random Forest

## Bagging және кездейсоқ орман

## Кездейсоқ орман

In [15]:

```

predictors = ['borrower_score', 'payment_inc_ratio']

```



```

outcome = 'outcome'

X = loan3000[predictors]
y = loan3000[outcome]

rf = RandomForestClassifier(n_estimators=500, random_state=1,
                           oob_score=True)

rf.fit(X, y)
print(rf.oob_decision_function_)
[[0.18131868 0.81868132]
 [0.26704545 0.73295455]
 [0.93333333 0.06666667]
 ...
 [1.          0.          ]
 [0.73157895 0.26842105]
 [0.68085106 0.31914894]]

```

In [16]:

```

n_estimator = list(range(20, 510, 5))
oobScores = []
for n in n_estimator:
    rf = RandomForestClassifier(n_estimators=n,
                              criterion='entropy', max_depth=5,
                              random_state=1, oob_score=True)

    rf.fit(X, y)
    oobScores.append(rf.oob_score_)

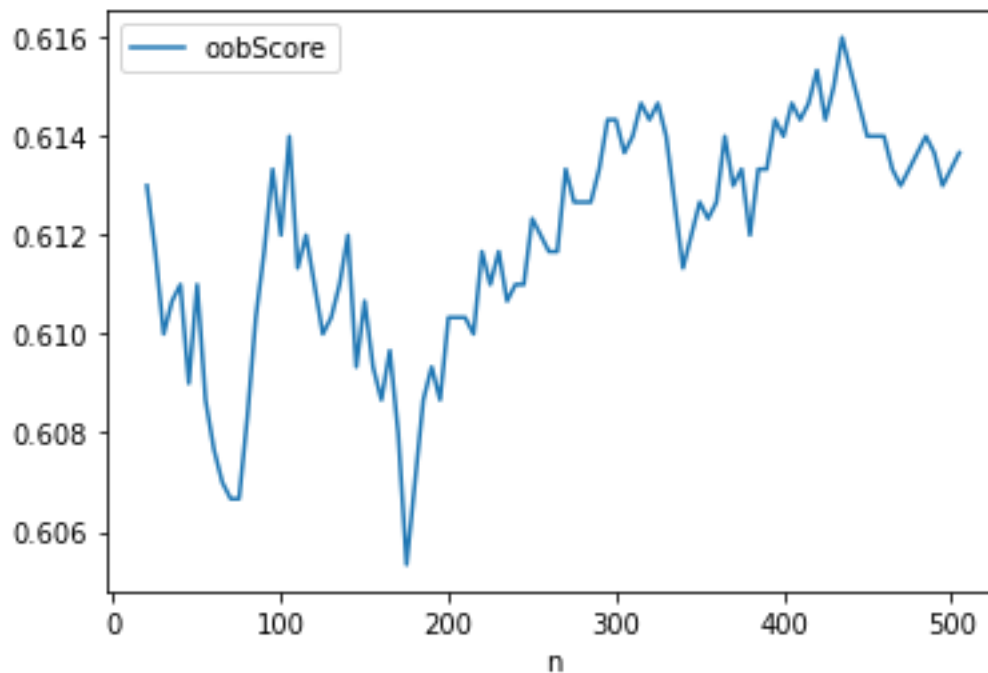
```

In [17]:

```

pd.DataFrame({
    'n': n_estimator,
    'oobScore': oobScores
}).plot(x='n', y='oobScore')
plt.show()

```



In [18]:

```

predictions = X.copy()
predictions['prediction'] = rf.predict(X)
predictions.head()

```

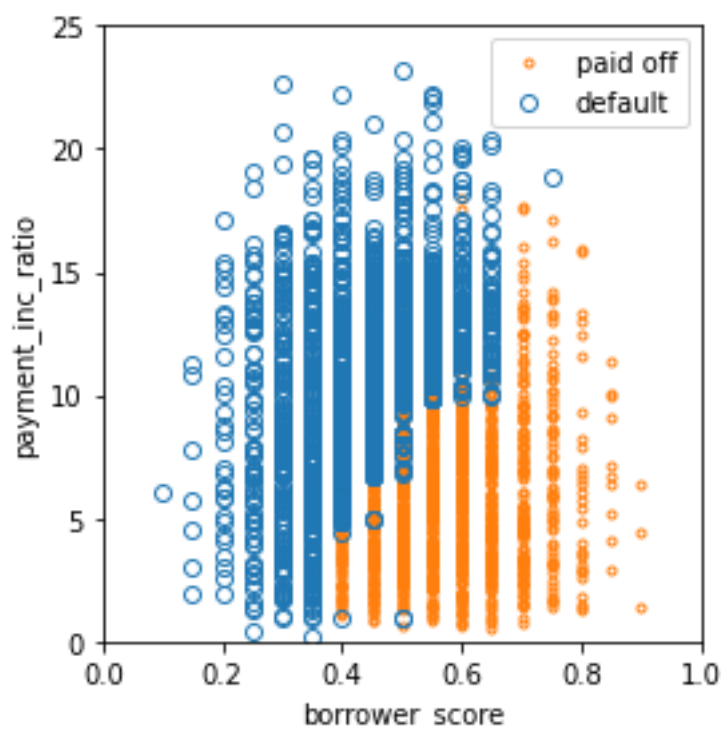
```

fig, ax = plt.subplots(figsize=(4, 4))

predictions.loc[predictions.prediction=='paid off'].plot(
    x='borrower_score', y='payment_inc_ratio', style='.',
    markerfacecolor='none', markeredgecolor='C1', ax=ax)
predictions.loc[predictions.prediction=='default'].plot(
    x='borrower_score', y='payment_inc_ratio', style='o',
    markerfacecolor='none', markeredgecolor='C0', ax=ax)
ax.legend(['paid off', 'default']);
ax.set_xlim(0, 1)
ax.set_ylim(0, 25)
ax.set_xlabel('borrower_score')
ax.set_ylabel('payment_inc_ratio')

plt.tight_layout()
plt.show()

```



## Variable importance

### Айнымалы маңыздылығы

This is different to R. The accuracy decrease is not available out of the box.

In [19]:

```

predictors = ['loan_amnt', 'term', 'annual_inc', 'dti',
              'payment_inc_ratio', 'revol_bal', 'revol_util',
              'purpose', 'delinq_2yrs_zero', 'pub_rec_zero',
              'open_acc', 'grade', 'emp_length', 'purpose_',
              'home_', 'emp_len_', 'borrower_score']
outcome = 'outcome'

X = pd.get_dummies(loan_data[predictors], drop_first=True)
y = loan_data[outcome]

```

```

rf_all = RandomForestClassifier(n_estimators=500, random_state=1)
rf_all.fit(X, y)

rf_all_entropy = RandomForestClassifier(n_estimators=500, random_state=1,
                                       criterion='entropy')

print(rf_all_entropy.fit(X, y))

RandomForestClassifier(criterion='entropy', n_estimators=500, random_state=1)
                                                                    In [20]:

rf = RandomForestClassifier(n_estimators=500)
scores = defaultdict(list)

# crossvalidate the scores on a number of different random splits of the data
for _ in range(3):
    train_X, valid_X, train_y, valid_y = train_test_split(X, y,
                                                         test_size=0.3)

    rf.fit(train_X, train_y)
    acc = metrics.accuracy_score(valid_y, rf.predict(valid_X))
    for column in X.columns:
        X_t = valid_X.copy()
        X_t[column] = np.random.permutation(X_t[column].values)
        shuff_acc = metrics.accuracy_score(valid_y, rf.predict(X_t))
        scores[column].append((acc-shuff_acc)/acc)
print('Features sorted by their score:')
print(sorted([(round(np.mean(score), 4), feat) for
              feat, score in scores.items()], reverse=True))

```

```

Features sorted by their score:
[(0.0739, 'borrower_score'), (0.0352, 'grade'), (0.0286, 'term_60 months'), (
0.0137, 'annual_inc'), (0.009, 'payment_inc_ratio'), (0.0021, 'purpose_small
_business'), (0.002, 'dti'), (0.0016, 'purpose_small_business'), (0.0011, 'ho
me_RENT'), (0.0011, 'delinq_2yrs_zero'), (0.001, 'revol_bal'), (0.0004, 'pub
_rec_zero'), (0.0003, 'home_OWN'), (0.0003, 'emp_len_ > 1 Year'), (0.0002,
'purpose_other'), (0.0002, 'emp_length'), (0.0001, 'purpose_vacation'), (0.0
001, 'purpose_house'), (0.0001, 'purpose_home_improvement'), (-0.0, 'purpose
moving'), (-0.0001, 'purpose_medical'), (-0.0001, 'purpose_medical'), (-0.00
02, 'purpose_wedding'), (-0.0004, 'purpose_major_purchase'), (-0.0004, 'open
_acc'), (-0.0005, 'purpose_major_purchase'), (-0.0006, 'purpose_home_improve
ment'), (-0.0007, 'purpose_other'), (-0.0007, 'purpose_debt_consolidation'),
(-0.0008, 'revol_util'), (-0.0013, 'purpose_credit_card'), (-0.0019, 'loan_am
nt'), (-0.0022, 'purpose_debt_consolidation')]

```

```

                                                                    In [21]:

importances = rf_all.feature_importances_

df = pd.DataFrame({
    'feature': X.columns,
    'Accuracy decrease': [np.mean(scores[column]) for column in
                        X.columns],
    'Gini decrease': rf_all.feature_importances_,
    'Entropy decrease': rf_all_entropy.feature_importances_,
})
df = df.sort_values('Accuracy decrease')

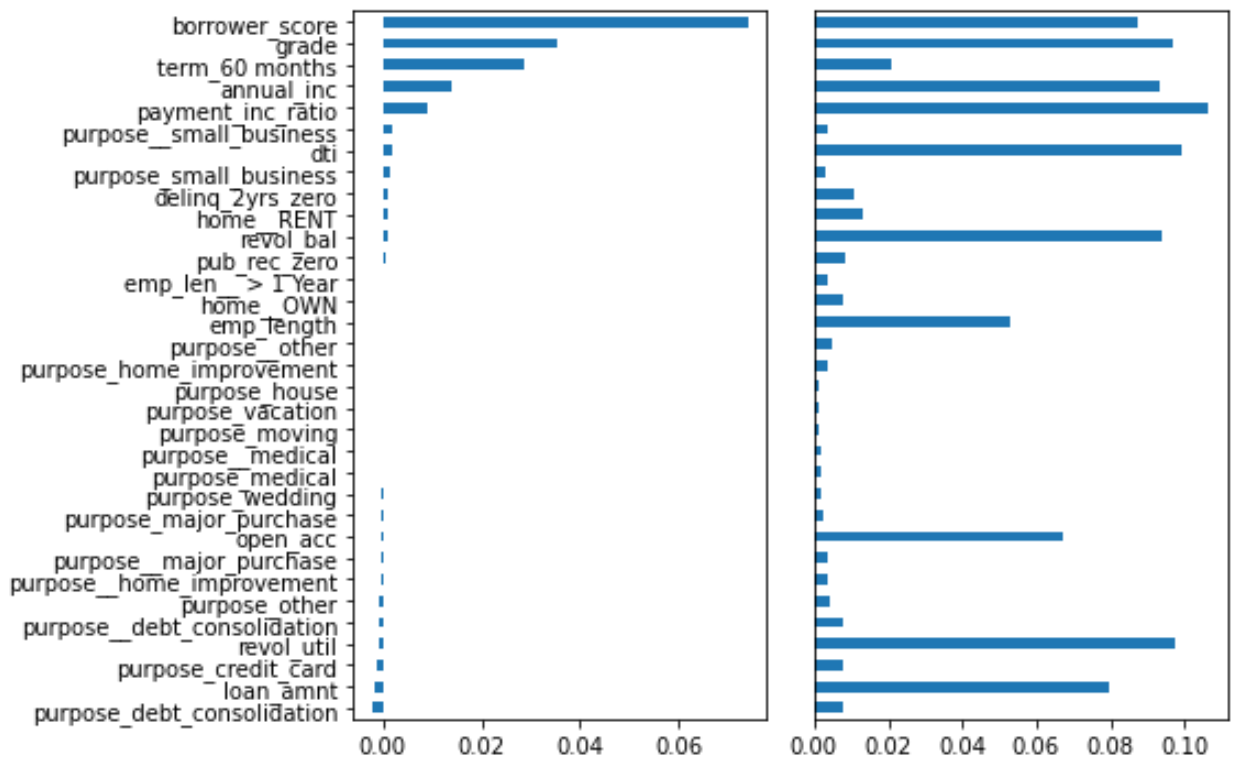
fig, axes = plt.subplots(ncols=2, figsize=(8, 5))
ax = df.plot(kind='barh', x='feature', y='Accuracy decrease',
             legend=False, ax=axes[0])
ax.set_ylabel('')

ax = df.plot(kind='barh', x='feature', y='Gini decrease',
             legend=False, ax=axes[1])

```

```
ax.set_ylabel('')
ax.get_yaxis().set_visible(False)

plt.tight_layout()
plt.show()
```



# Boosting

## XGBoost

In [22]:

```
predictors = ['borrower_score', 'payment_inc_ratio']
outcome = 'outcome'

X = loan3000[predictors]
y = pd.Series([1 if o == 'default' else 0 for o in loan3000[outcome]])

xgb = XGBClassifier(objective='binary:logistic', subsample=.63,
                    use_label_encoder=False, eval_metric='error')
print(xgb.fit(X, y))

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=1, enable_categorical=False,
              eval_metric='error', gamma=0, gpu_id=-1, importance_type=None,
              interaction_constraints='', learning_rate=0.300000012,
              max_delta_step=0, max_depth=6, min_child_weight=1, missing=nan,
              monotone_constraints='()', n_estimators=100, n_jobs=8,
              num_parallel_tree=1, predictor='auto', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=0.63,
              tree_method='exact', use_label_encoder=False,
              validate_parameters=1, verbosity=None)
```

```
/opt/conda/lib/python3.9/site-packages/xgboost/data.py:262: FutureWarning: pandas.Int64Index is deprecated and will be removed from pandas in a future version. Use pandas.Index with the appropriate dtype instead.
```

```
elif isinstance(data.columns, (pd.Int64Index, pd.RangeIndex)):
```

In [23]:

```
xgb_df = X.copy()
xgb_df['prediction'] = ['default' if p == 1 else 'paid off' for p in
xgb.predict(X)]
xgb_df['prob_default'] = xgb.predict_proba(X)[:, 0]
print(xgb_df.head())
```

	borrower_score	payment_inc_ratio	prediction	prob_default
0	0.40	5.11135	paid off	0.828856
1	0.40	5.43165	default	0.260156
2	0.70	9.23003	default	0.320805
3	0.40	2.33482	paid off	0.678005
4	0.45	12.10320	default	0.140204

```
/opt/conda/lib/python3.9/site-packages/xgboost/data.py:262: FutureWarning: pandas.Int64Index is deprecated and will be removed from pandas in a future version. Use pandas.Index with the appropriate dtype instead.
```

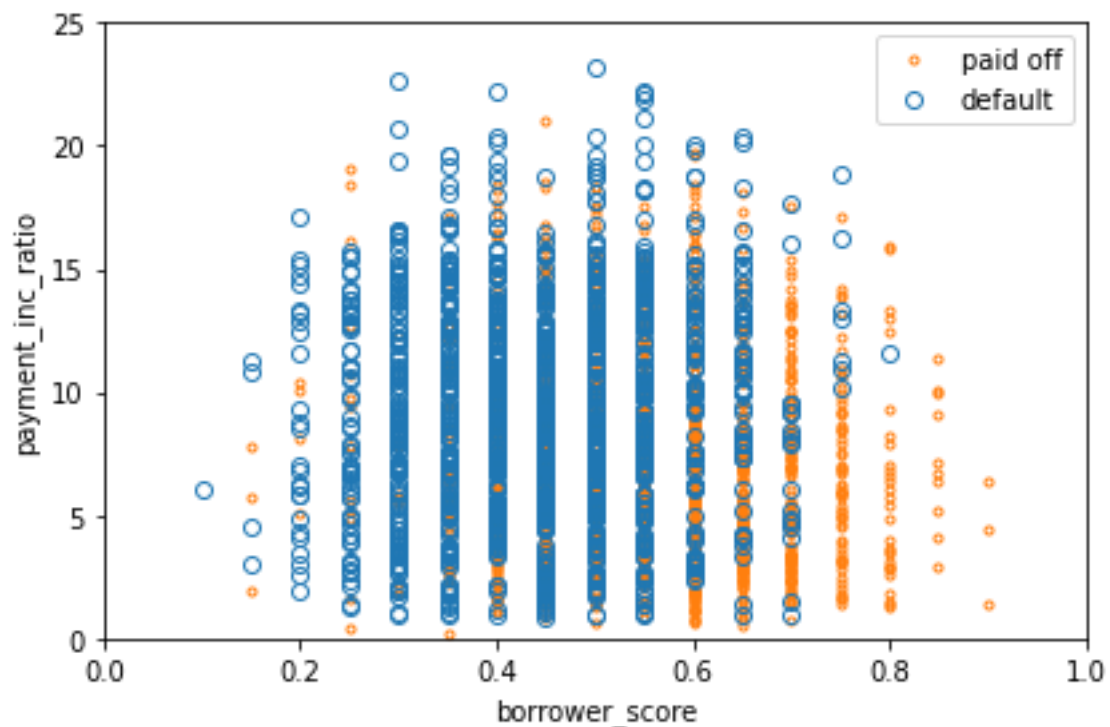
```
elif isinstance(data.columns, (pd.Int64Index, pd.RangeIndex)):
```

In [24]:

```
fig, ax = plt.subplots(figsize=(6, 4))
```

```
xgb_df.loc[xgb_df.prediction=='paid off'].plot(
    x='borrower_score', y='payment_inc_ratio', style='.',
    markerfacecolor='none', markeredgecolor='C1', ax=ax)
xgb_df.loc[xgb_df.prediction=='default'].plot(
    x='borrower_score', y='payment_inc_ratio', style='o',
    markerfacecolor='none', markeredgecolor='C0', ax=ax)
ax.legend(['paid off', 'default']);
ax.set_xlim(0, 1)
ax.set_ylim(0, 25)
ax.set_xlabel('borrower_score')
ax.set_ylabel('payment_inc_ratio')
```

```
plt.tight_layout()
plt.show()
```



## Regularization: Avoiding Overfitting

### Реттеу: қайта оқытудан аулақ болу

In [25]:

```

predictors = ['loan_amnt', 'term', 'annual_inc', 'dti',
              'payment_inc_ratio', 'revol_bal', 'revol_util',
              'purpose', 'delinq_2yrs_zero', 'pub_rec_zero',
              'open_acc', 'grade', 'emp_length', 'purpose_',
              'home_', 'emp_len_', 'borrower_score']
outcome = 'outcome'

X = pd.get_dummies(loan_data[predictors], drop_first=True)
y = pd.Series([1 if o == 'default' else 0 for o in loan_data[outcome]])

train_X, valid_X, train_y, valid_y = train_test_split(X, y, test_size=10000)

xgb_default = XGBClassifier(objective='binary:logistic', n_estimators=250,
                           max_depth=6,
                           reg_lambda=0, learning_rate=0.3, subsample=1,
                           use_label_encoder=False, eval_metric='error')
xgb_default.fit(train_X, train_y)

xgb_penalty = XGBClassifier(objective='binary:logistic', n_estimators=250,
                           max_depth=6,
                           reg_lambda=1000, learning_rate=0.1,
                           subsample=0.63,
                           use_label_encoder=False, eval_metric='error')
print(xgb_penalty.fit(train_X, train_y))

/opt/conda/lib/python3.9/site-packages/xgboost/data.py:262: FutureWarning: pa
ndas.Int64Index is deprecated and will be removed from pandas in a future ver
sion. Use pandas.Index with the appropriate dtype instead.
  elif isinstance(data.columns, (pd.Int64Index, pd.RangeIndex)):

```

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=1, enable_categorical=False,
              eval_metric='error', gamma=0, gpu_id=-1, importance_type=None,
              interaction_constraints='', learning_rate=0.1, max_delta_step=0,
              max_depth=6, min_child_weight=1, missing=nan,
              monotone_constraints='()', n_estimators=250, n_jobs=8,
              num_parallel_tree=1, predictor='auto', random_state=0,
              reg_alpha=0, reg_lambda=1000, scale_pos_weight=1, subsample=0.6,
              tree_method='exact', use_label_encoder=False,
              validate_parameters=1, verbosity=None)
```

In [26]:

```
pred_default = xgb_default.predict_proba(train_X)[:, 1]
error_default = abs(train_y - pred_default) > 0.5
print('default (train): ', np.mean(error_default))

pred_default = xgb_default.predict_proba(valid_X)[:, 1]
error_default = abs(valid_y - pred_default) > 0.5
print('default: ', np.mean(error_default))

pred_penalty = xgb_penalty.predict_proba(valid_X)[:, 1]
error_penalty = abs(valid_y - pred_penalty) > 0.5
print('penalty: ', np.mean(error_penalty))

default (train): 0.12678965536755135
default: 0.3553
penalty: 0.3311
```

In [27]:

```
results = []
for ntree_limit in range(1, 250):
    iteration_range = [1, ntree_limit + 1]
    train_default = xgb_default.predict_proba(train_X,
                                              iteration_range=iteration_range)[:, 1]
    train_penalty = xgb_penalty.predict_proba(train_X,
                                              iteration_range=iteration_range)[:, 1]
    pred_default = xgb_default.predict_proba(valid_X,
                                             iteration_range=iteration_range)[:, 1]
    pred_penalty = xgb_penalty.predict_proba(valid_X,
                                             iteration_range=iteration_range)[:, 1]
    results.append({
        'iterations': ntree_limit,
        'default train': np.mean(abs(train_y - train_default) > 0.5),
        'penalty train': np.mean(abs(train_y - train_penalty) > 0.5),
        'default test': np.mean(abs(valid_y - pred_default) > 0.5),
        'penalty test': np.mean(abs(valid_y - pred_penalty) > 0.5),
    })
```

```
results = pd.DataFrame(results)
print(results.head())
```

	iterations	default train	penalty train	default test	penalty test
0	1	0.342680	0.337021	0.3600	0.3532
1	2	0.331419	0.334559	0.3530	0.3483
2	3	0.321883	0.333795	0.3407	0.3511
3	4	0.317300	0.334446	0.3412	0.3508
4	5	0.315772	0.338662	0.3404	0.3510

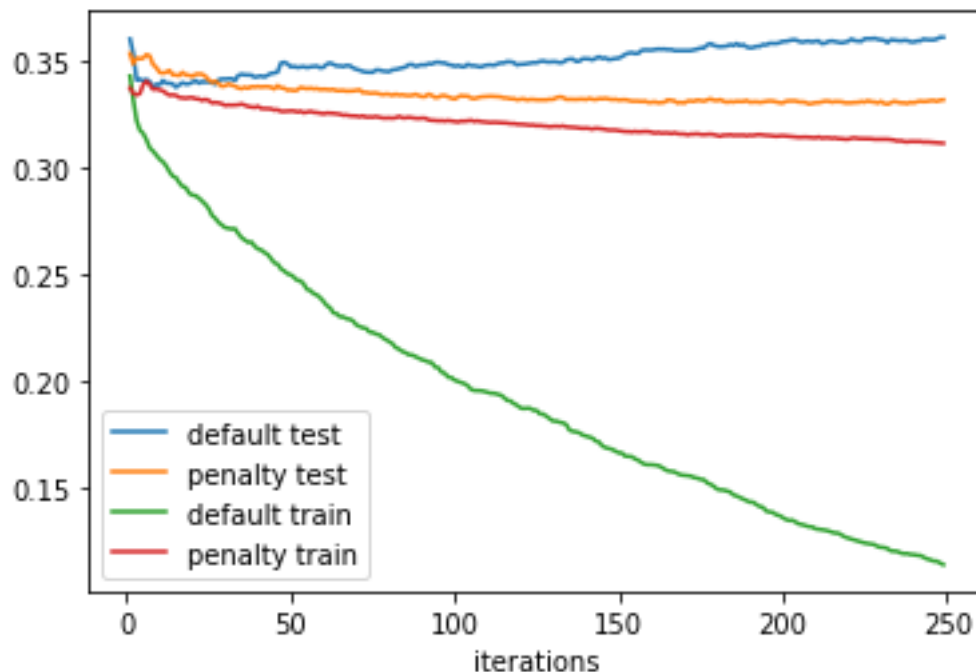
In [28]:

```
ax = results.plot(x='iterations', y='default test')
```

```

results.plot(x='iterations', y='penalty test', ax=ax)
results.plot(x='iterations', y='default train', ax=ax)
results.plot(x='iterations', y='penalty train', ax=ax)
plt.show()

```



## Hyperparameters and Cross-Validation

## Гиперпараметрлер және Кросс-тексеру

In [29]:

```

idx = np.random.choice(range(5), size=len(X), replace=True)
error = []
for eta, max_depth in product([0.1, 0.5, 0.9], [3, 6, 9]):
    xgb = XGBClassifier(objective='binary:logistic', n_estimators=250,
                        max_depth=max_depth, learning_rate=eta,
                        use_label_encoder=False, eval_metric='error')
    cv_error = []
    for k in range(5):
        fold_idx = idx == k
        train_X = X.loc[~fold_idx]; train_y = y[~fold_idx]
        valid_X = X.loc[fold_idx]; valid_y = y[fold_idx]

        xgb.fit(train_X, train_y)
        pred = xgb.predict_proba(valid_X)[:, 1]
        cv_error.append(np.mean(abs(valid_y - pred) > 0.5))
    error.append({
        'eta': eta,
        'max_depth': max_depth,
        'avg_error': np.mean(cv_error)
    })
print(error[-1])
errors = pd.DataFrame(error)
print(errors)

```



```
/opt/conda/lib/python3.9/site-packages/xgboost/data.py:262: FutureWarning: pa
ndas.Int64Index is deprecated and will be removed from pandas in a future ver
sion. Use pandas.Index with the appropriate dtype instead.
```

```
elif isinstance(data.columns, (pd.Int64Index, pd.RangeIndex)):
{'eta': 0.1, 'max_depth': 3, 'avg_error': 0.32895002759265257}
{'eta': 0.1, 'max_depth': 6, 'avg_error': 0.3360983351965724}
{'eta': 0.1, 'max_depth': 9, 'avg_error': 0.3461055248368196}
{'eta': 0.5, 'max_depth': 3, 'avg_error': 0.3425077028439957}
{'eta': 0.5, 'max_depth': 6, 'avg_error': 0.3694194098086156}
{'eta': 0.5, 'max_depth': 9, 'avg_error': 0.3764740475715814}
{'eta': 0.9, 'max_depth': 3, 'avg_error': 0.3519677896500327}
{'eta': 0.9, 'max_depth': 6, 'avg_error': 0.3871414745778243}
{'eta': 0.9, 'max_depth': 9, 'avg_error': 0.38706390437335264}
eta  max_depth  avg_error
0  0.1          3    0.328950
1  0.1          6    0.336098
2  0.1          9    0.346106
3  0.5          3    0.342508
4  0.5          6    0.369419
5  0.5          9    0.376474
6  0.9          3    0.351968
7  0.9          6    0.387141
8  0.9          9    0.387064
```

In [30]:

```
print(errors.pivot_table(index='eta', columns='max_depth',
values='avg_error') * 100)
```

```
max_depth      3          6          9
eta
0.1      32.895003  33.609834  34.610552
0.5      34.250770  36.941941  37.647405
0.9      35.196779  38.714147  38.706390
```