

МИНОБРНАУКИ РОССИИ
САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ
ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ
«ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)
Кафедра МОЭВМ

ОТЧЕТ
по лабораторной работе №7
по дисциплине «Машинное обучение»
Тема: Классификация (Байесовские методы, деревья)

Студент гр. 6307

Золотухин М. А.

Преподаватель

Жангиров Т. Р.

Санкт-Петербург

2020

Байесовские методы

Проведем классификацию наблюдений наивным байесовским методом.

```
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
y_pred = gnb.fit(X_train, y_train).predict(X_test)
print((y_test != y_pred).sum()) #количество наблюдений, который были неправильно
определены
```

4

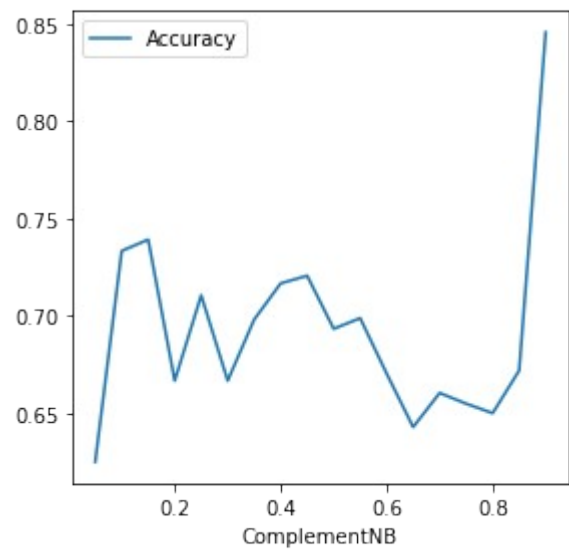
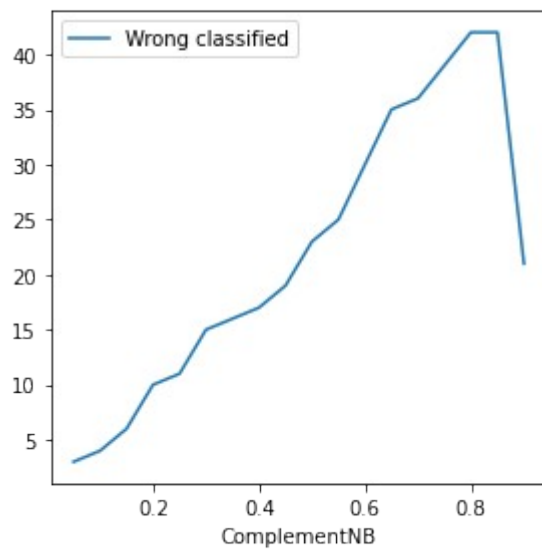
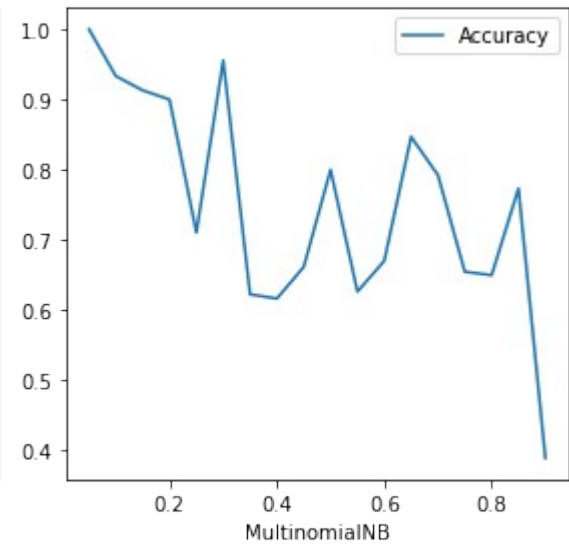
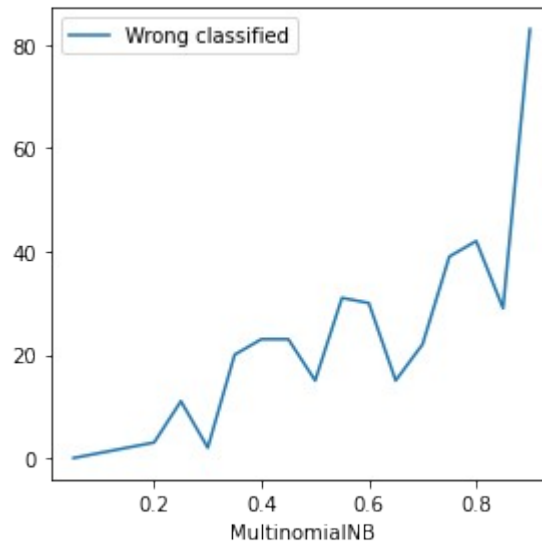
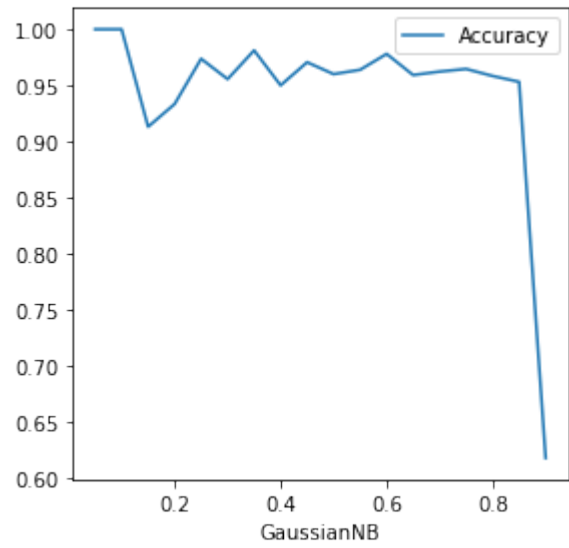
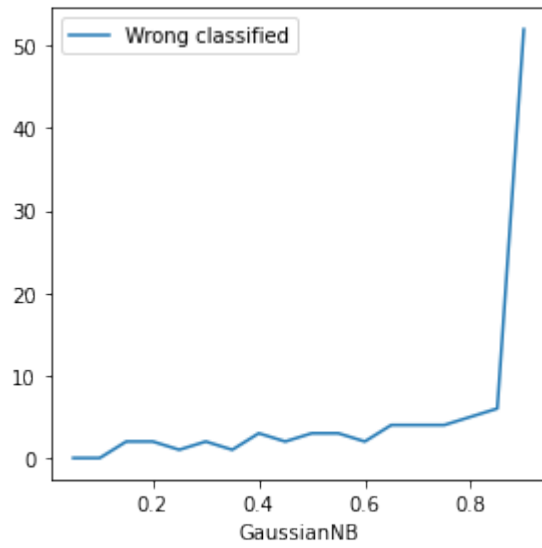
Атрибуты данного классификатора:

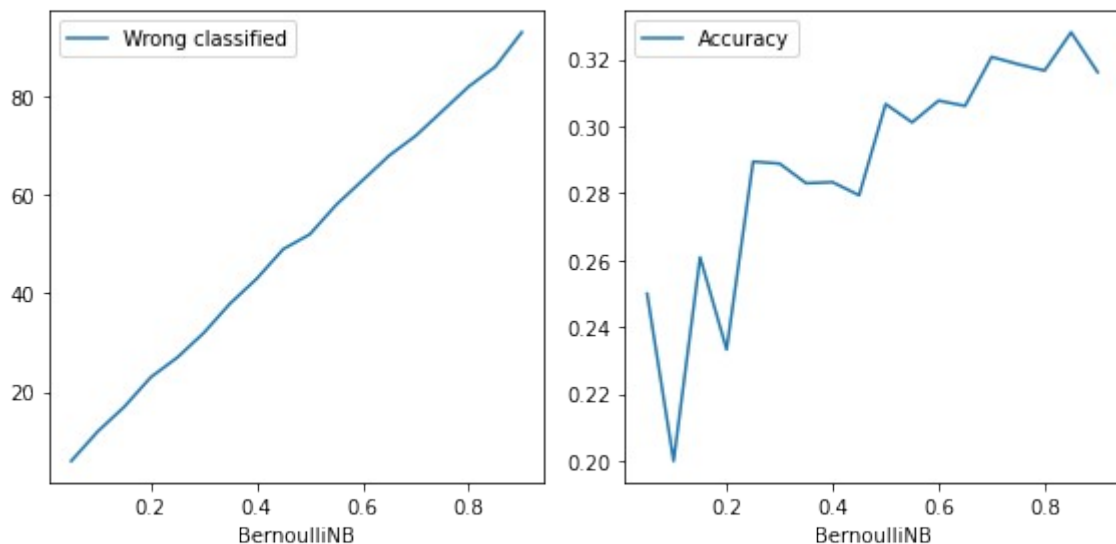
- *classcount* - количество тренировочных сэмплов, которые наблюдаются в каждом классе.
- *classprior*- вероятность каждого из классов.
- *classes_* - метки классов, известные классификатору.
- *epsilon_* - абсолютное добавочное значение для дисперсий.
- *sigma_* - дисперсия каждой фичи в каждом классе.
- *theta_* - среднее каждой фичи каждого класса.

Точность классификации:

```
print(gnb.score(X_test, y_test))
0.9466666666666667
```

Построим график зависимости неправильно классифицированных наблюдений и точности классификации от размера тестовой выборки. Размер тестовой выборки изменяйте от 0.05 до 0.95 с шагом 0.05. Параметр *random_state* сделаем равным номеру своей зачетной книжки. Проведем также классификацию используя [MultinomialNB](#), [ComplementNB](#), [BernoulliNB](#).





Gaussian — стандартный алгоритм, подразумевает, что фичи распределяются нормально.

Multinomial — алгоритм подразумевает, что данные подчиняются мультиномиальному распределению.

Complement — алгоритм, который используется для не сбалансированных датасетов.

Bernoulli — алгоритм подразумевает, что данные подчиняются распределению Бернулли.

Классифицирующие деревья

Проведем классификацию при помощи деревьев на тех же данных.

```
from sklearn import tree

clf = tree.DecisionTreeClassifier()

y_pred = clf.fit(X_train, y_train).predict(X_test)
print((y_test != y_pred).sum())
4
```

Используя функцию score() выведите точность классификации

```
print("Score: ", clf.score(X_test, y_test))
Score: 0.9466666666666667
```

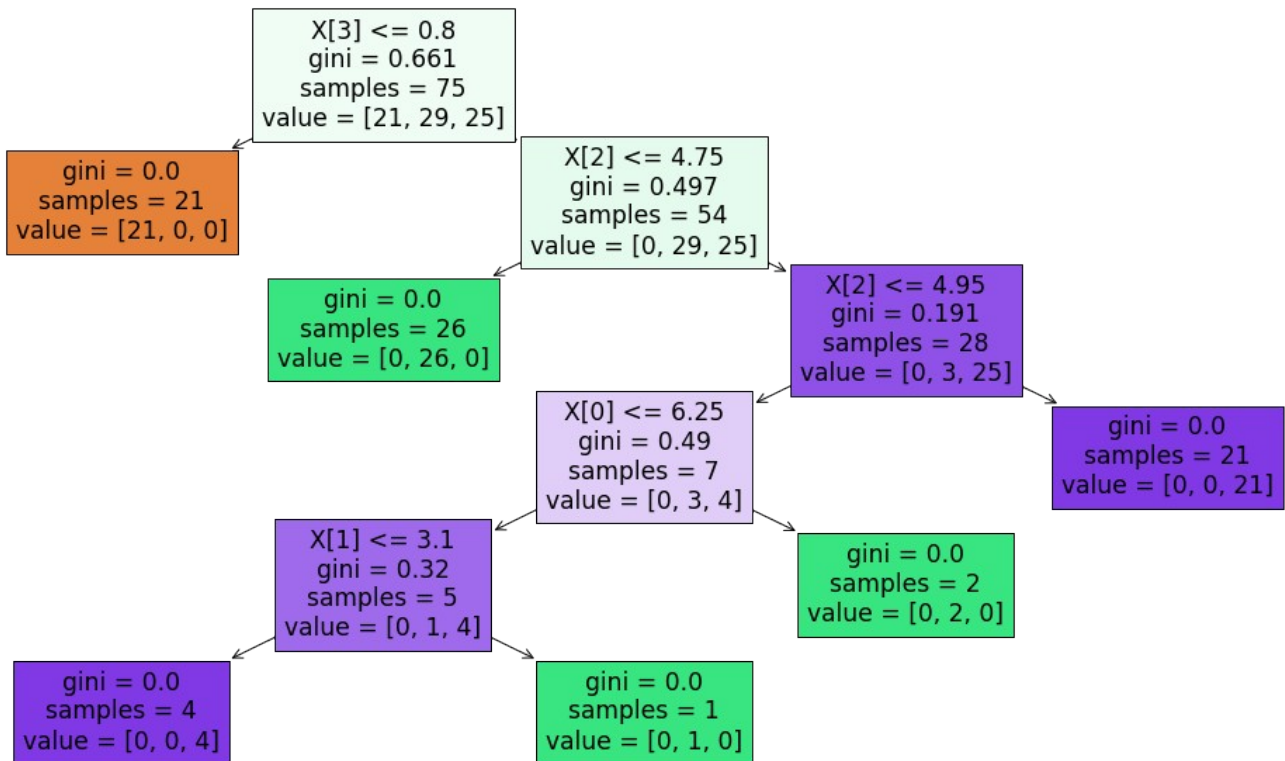
Выведите характеристики дерева, количество листьев и глубину, используя функции `get_n_leaves` и `get_depth`

```
print('Num of leaves: ', clf.get_n_leaves())
print('Depth: ', clf.get_depth())
Num of leaves: 6
Depth: 5
```

Выведем изображение полученного дерева

```
import matplotlib.pyplot as plt

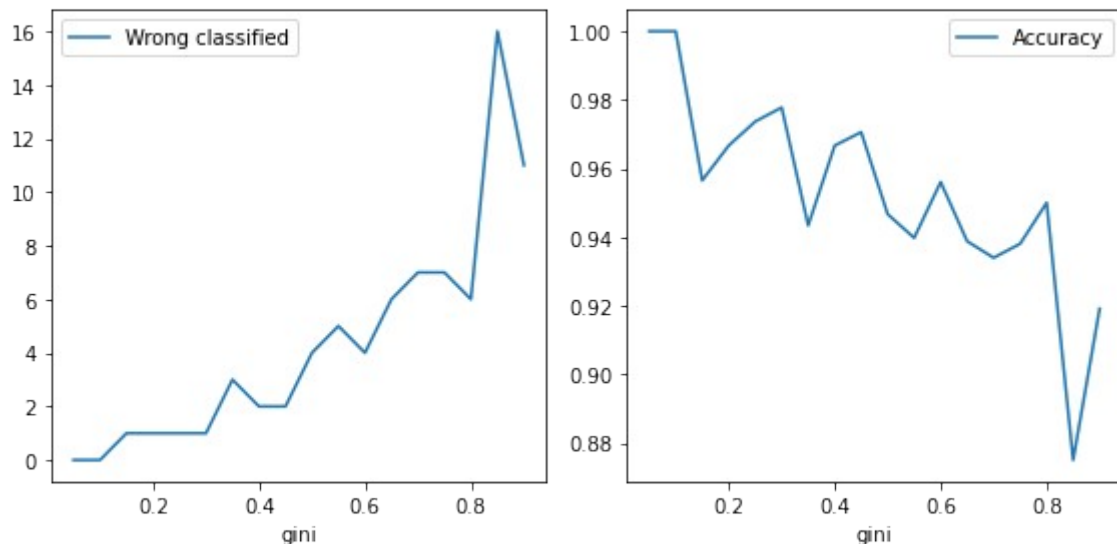
plt.subplots(1,1,figsize = (20,10))
tree.plot_tree(clf, filled = True)
plt.show()
```



На рисунке изображено дерево. Каждый узел — вопрос, который представлен первой строкой. В зависимости от ответа на вопрос (да или нет), данное идет по одной ветке или по другой, пока не дойдет до листа. Samples — количество измерений, которые прошли через узел. Gini — примесь gini. Value — количества узлов в каждом из классов (в данном случае класса 3, как и элементов массива).

Построим график зависимости неправильно классифицированных наблюдений и точности классификации от размера тестовой выборки. Размер тестовой выборки будет изменяться от 0.05 до 0.95 с шагом 0.05. Параметр `random_state` равен 630706.

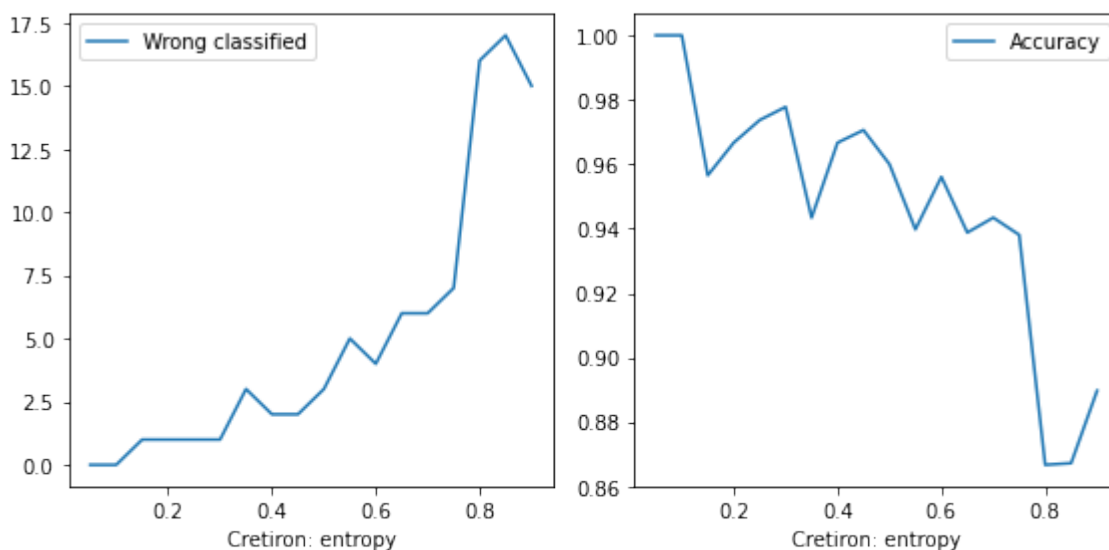
```
plot_clf(clf, "gini")
```



Исследуем работу классифицирующего дерева при различных параметрах `criterion`, `splitter`, `max_depth`, `min_samples_split`, `min_samples_leaf`.

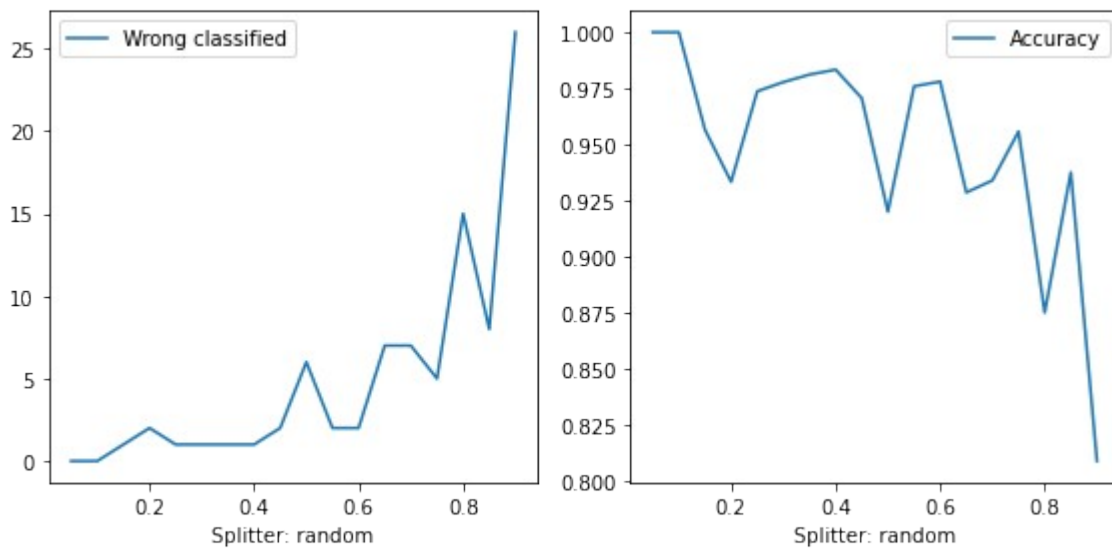
Entropy criterion

```
clf = tree.DecisionTreeClassifier(criterion="entropy")
plot_clf(clf, "Cretiron: entropy")
```



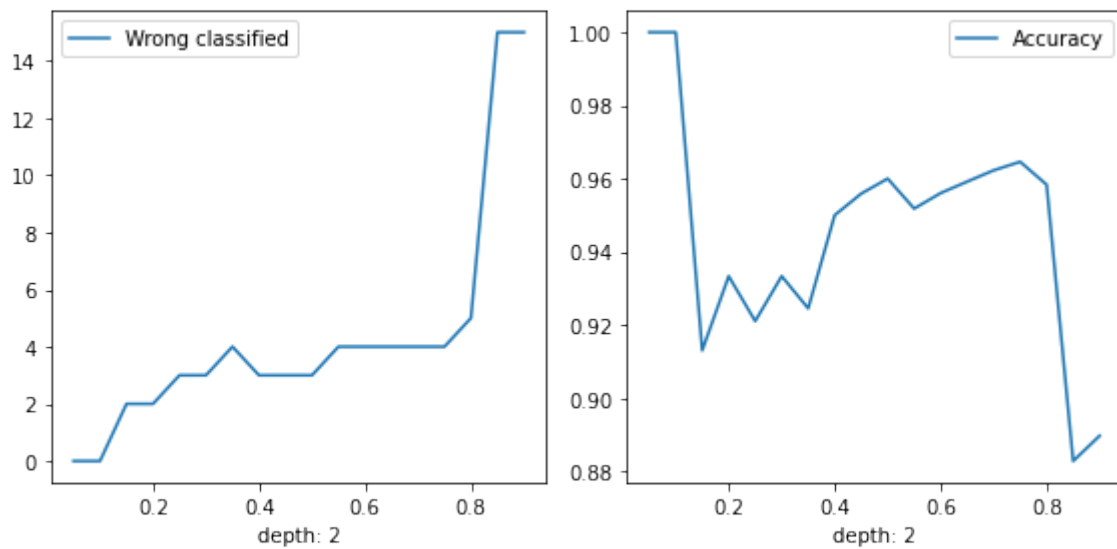
Random Splitter

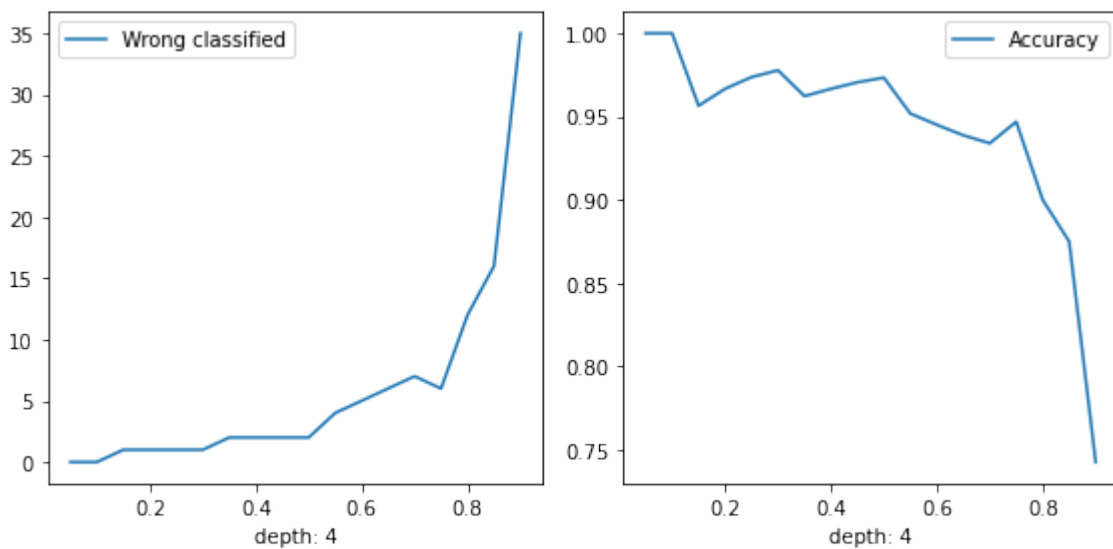
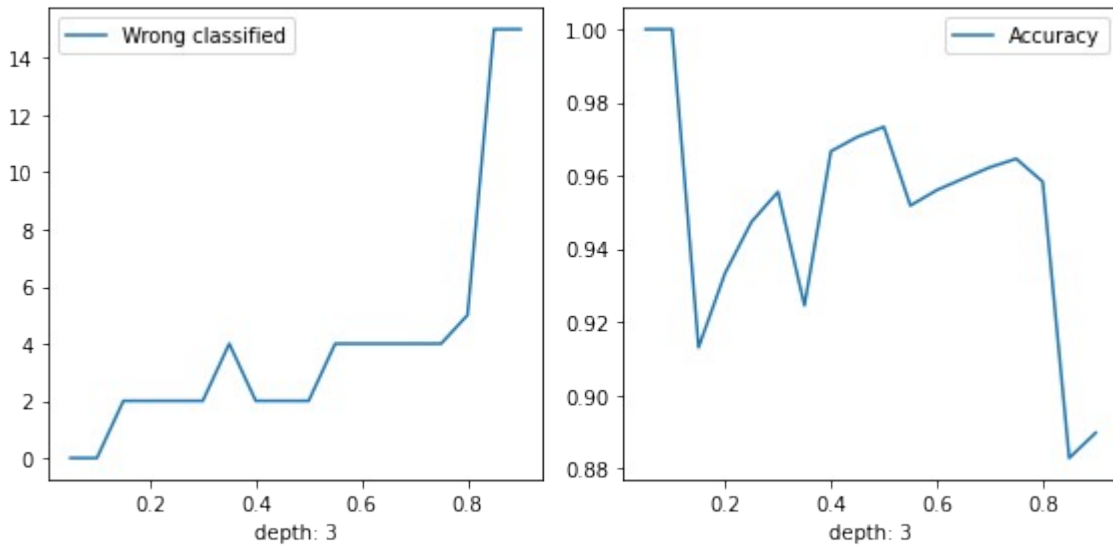
```
clf = tree.DecisionTreeClassifier(splitter="random")
plot_clf(clf, "Splitter: random")
```



Depth

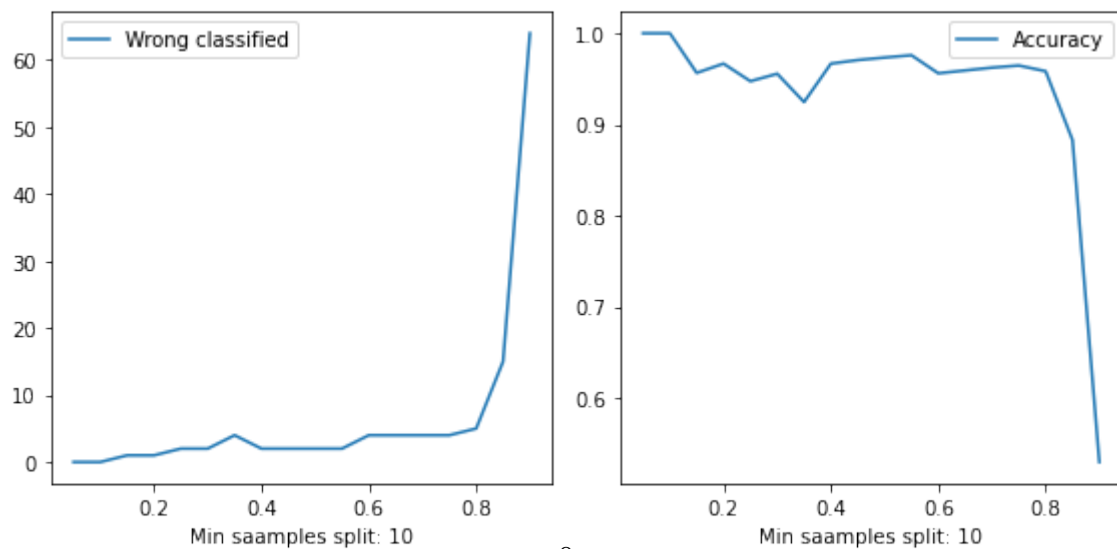
```
depths = np.arange(2, 5, 1)
for depth in depths:
    clf = tree.DecisionTreeClassifier(max_depth=depth)
    plot_clf(clf, f'depth: {depth}')
```

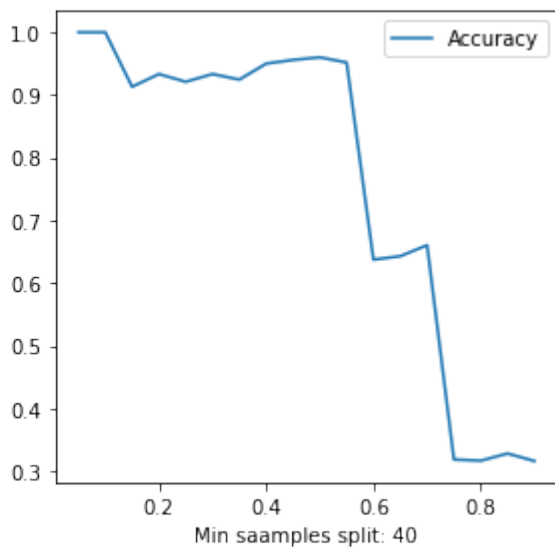
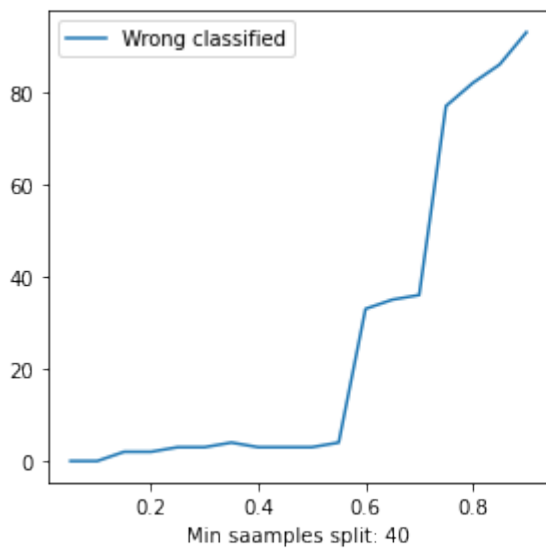
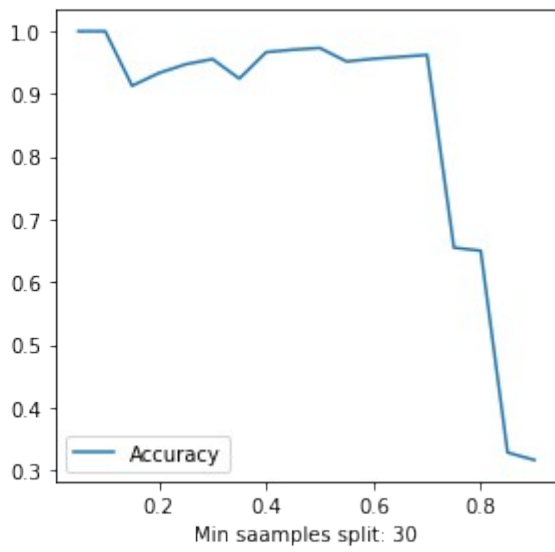
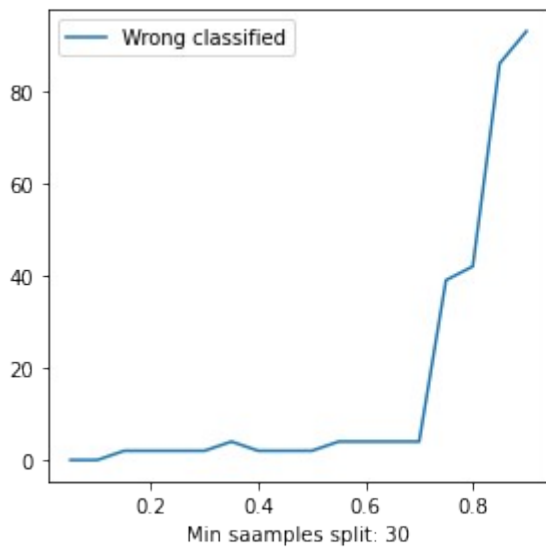
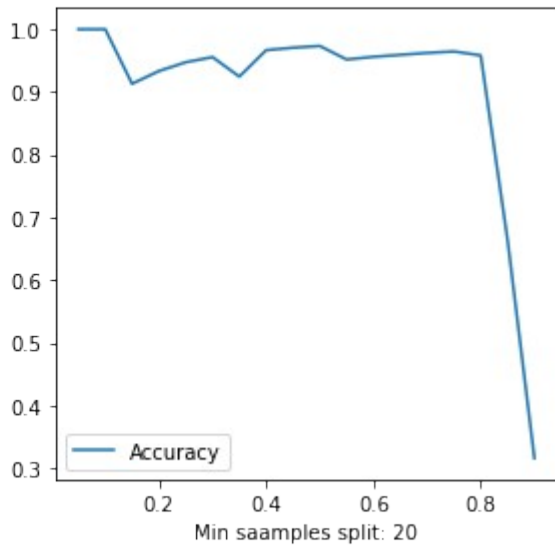
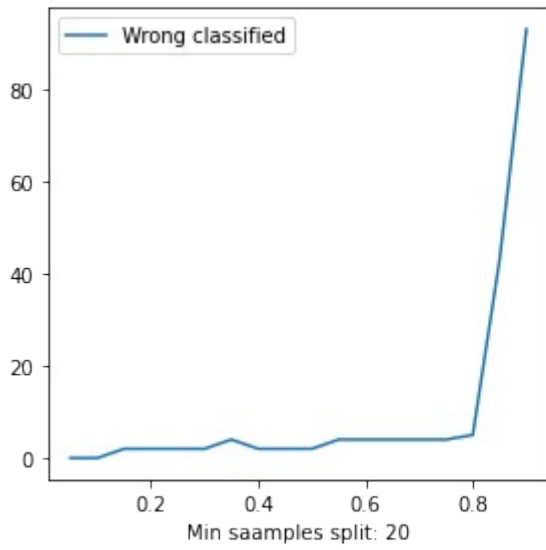


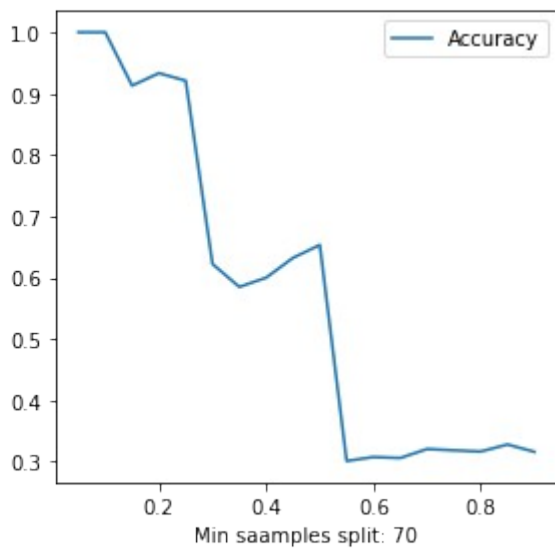
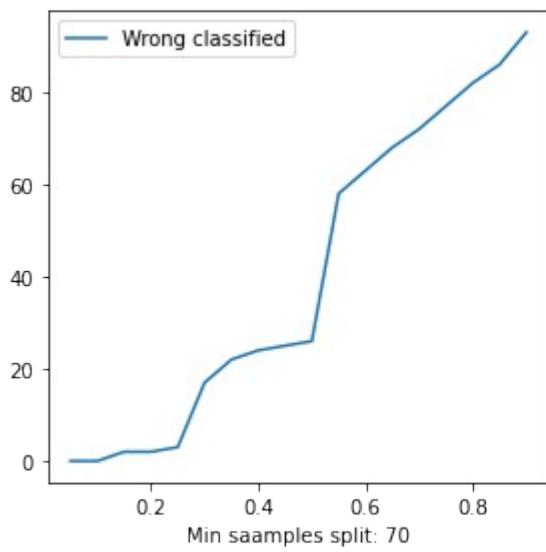
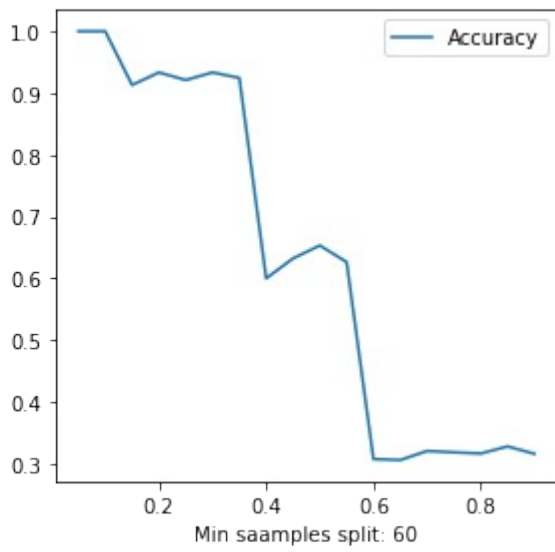
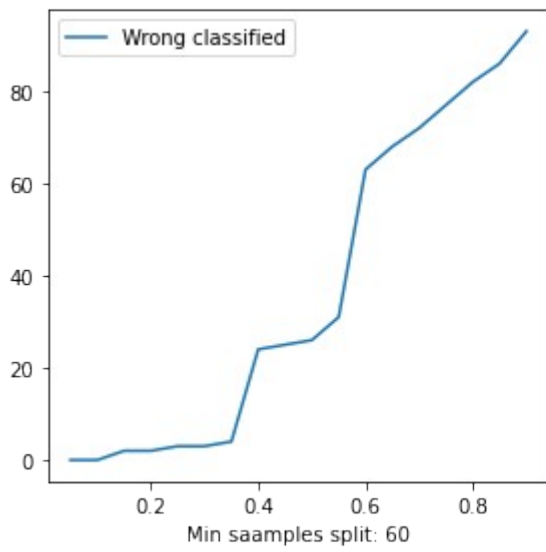
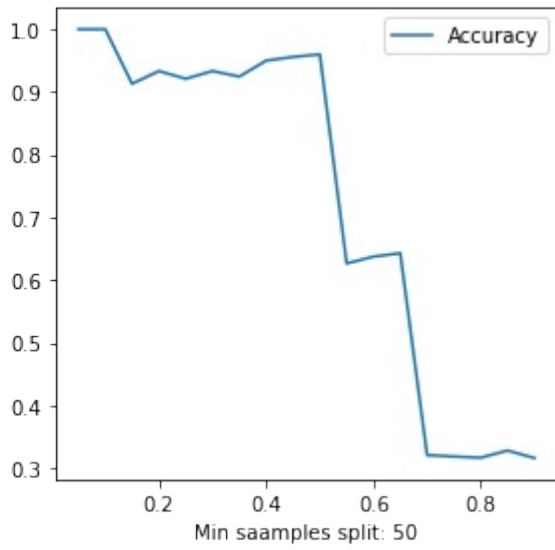
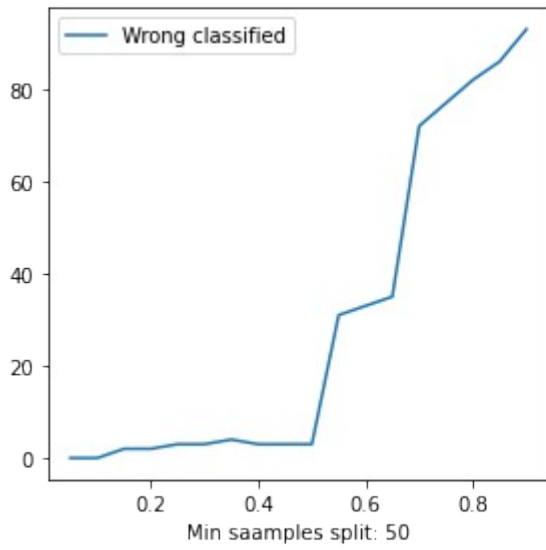


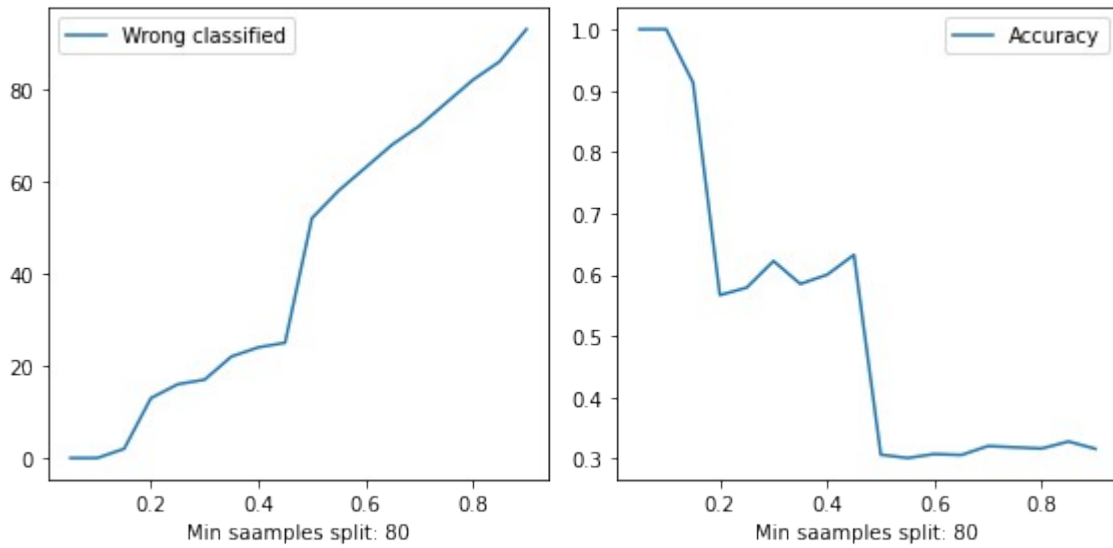
Min Samples Split

```
min_samples = np.arange(10, 90, 10)
for min_sample in min_samples:
    clf = tree.DecisionTreeClassifier(min_samples_split=min_sample)
    plot_clf(clf, f'Min saamples split: {min_sample}')
```



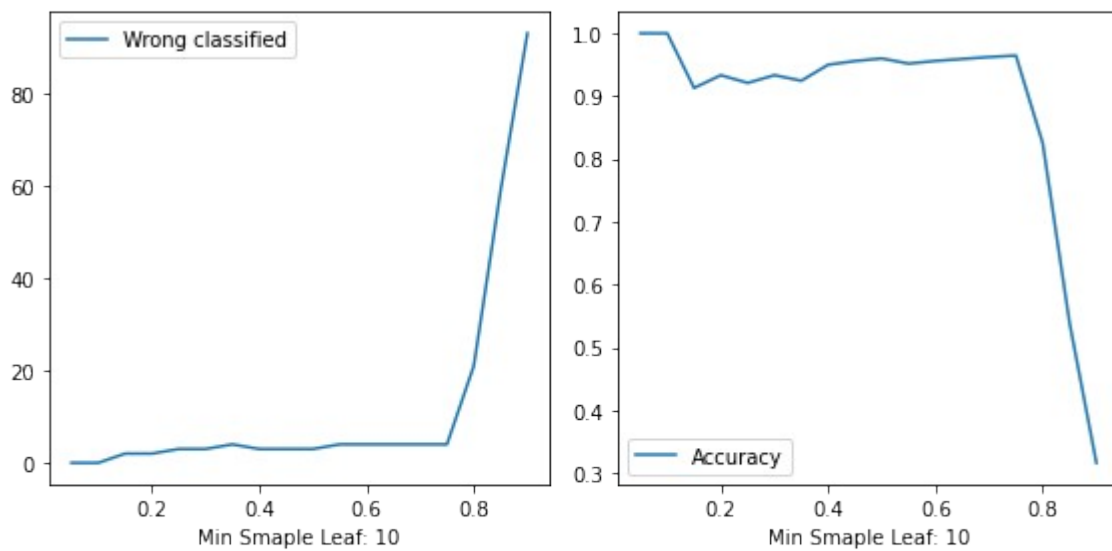


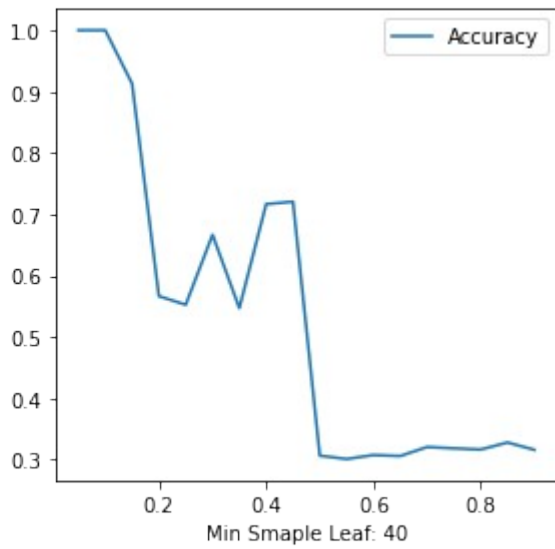
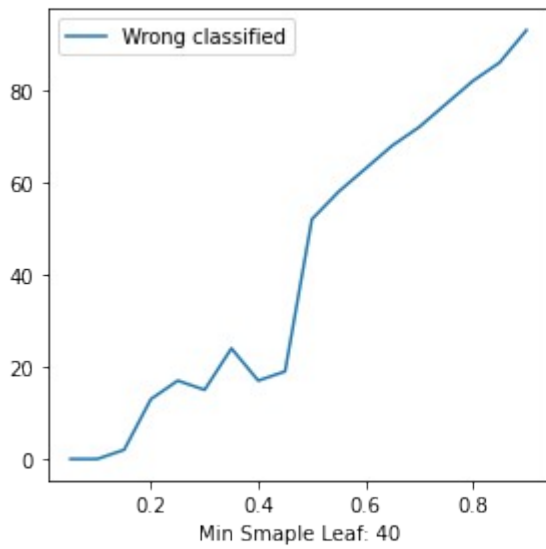
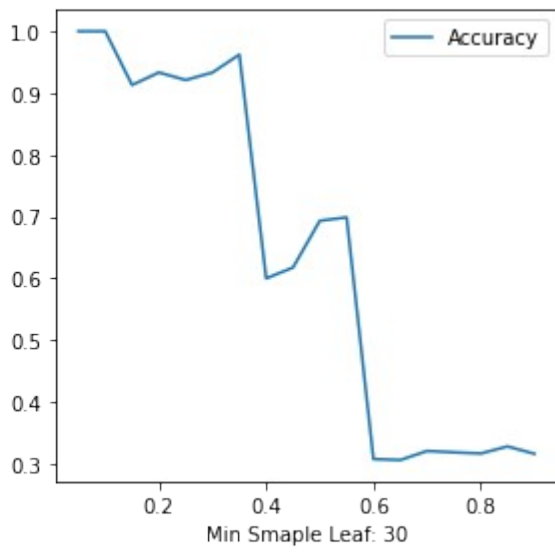
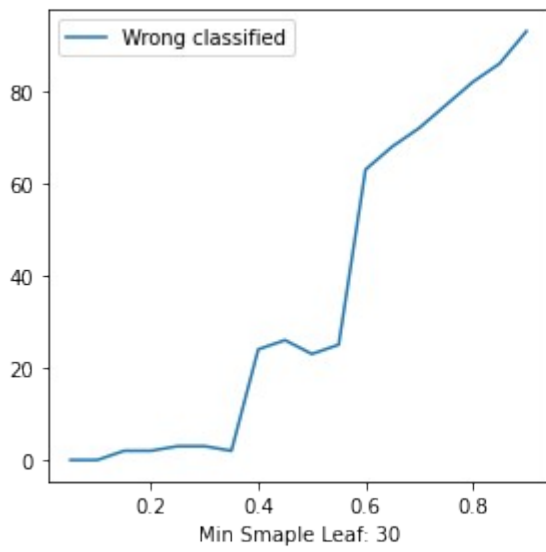
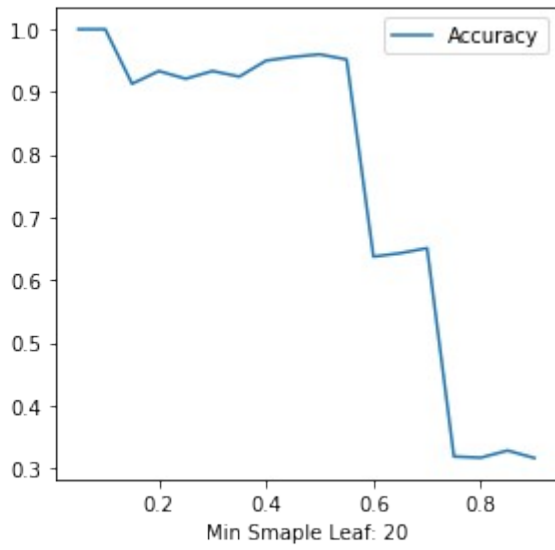
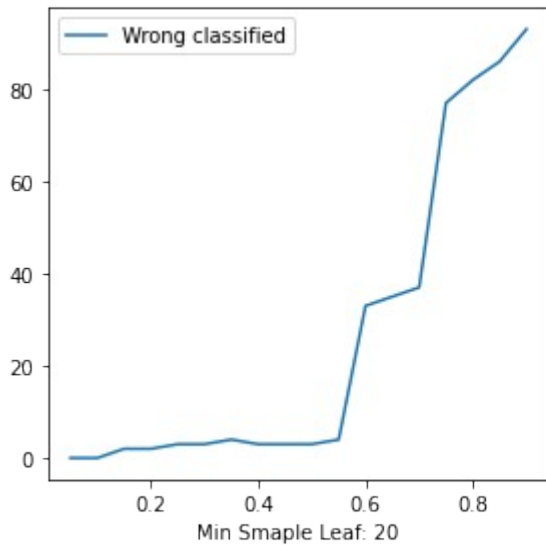


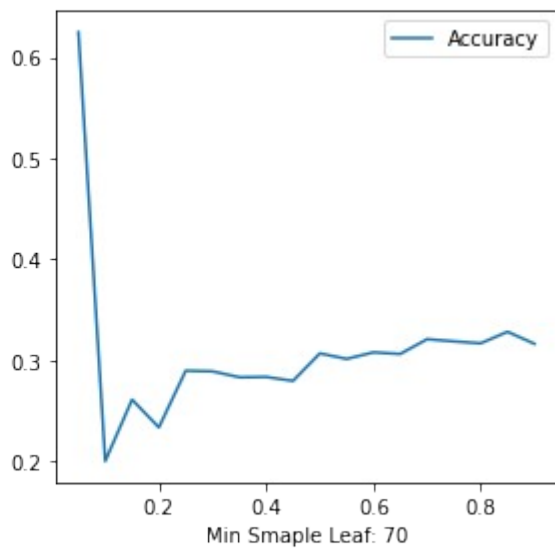
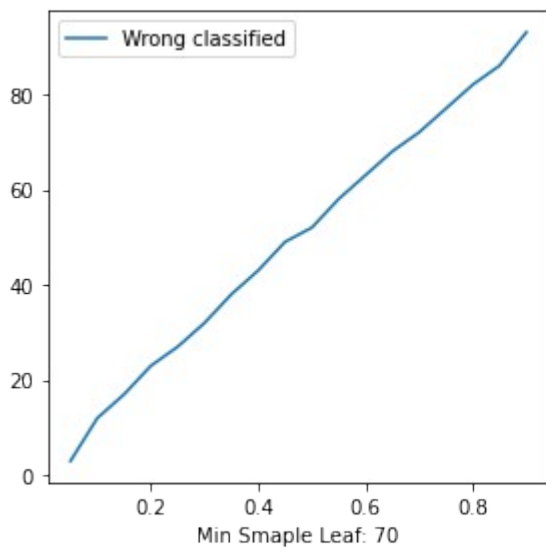
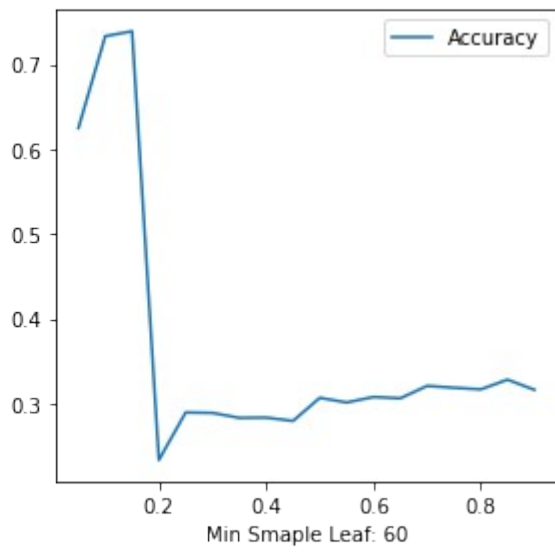
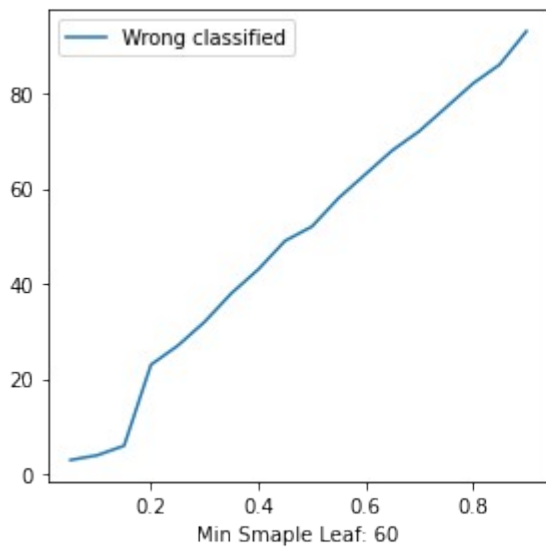
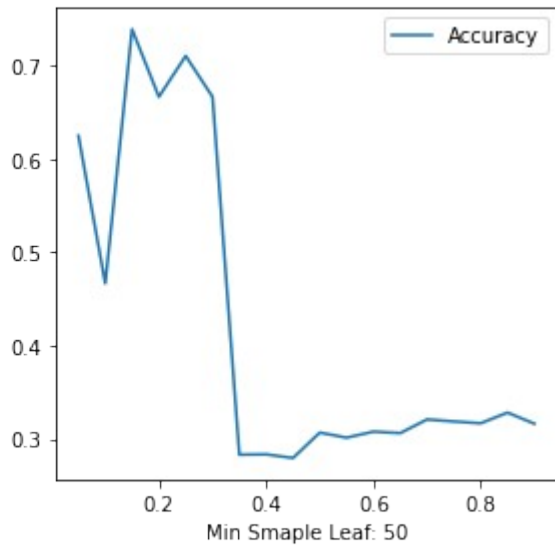
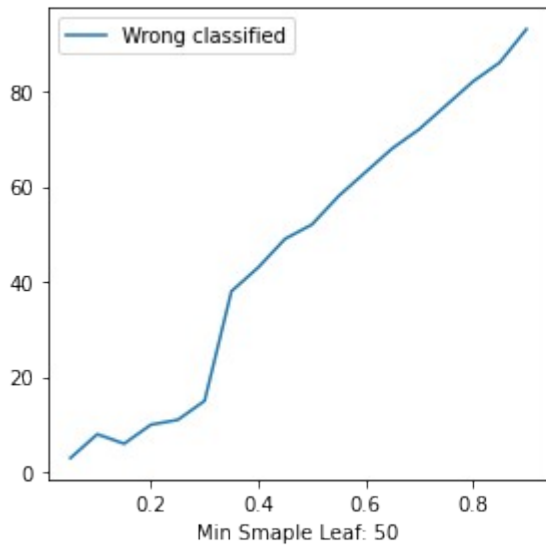


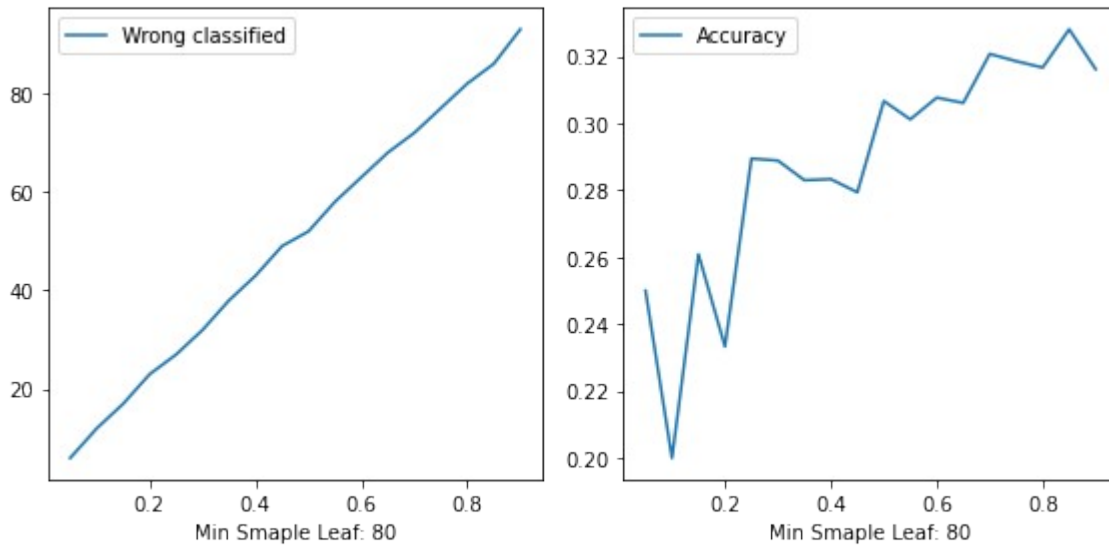
Min sample leaf

```
min_samples = np.arange(10, 90, 10)
for min_sample in min_samples:
    clf = tree.DecisionTreeClassifier(min_samples_leaf=min_sample)
    plot_clf(clf, f'Min Smaple Leaf: {min_sample}')
```









Объяснение параметров:

- Criterion — функция для измерения качества разбиения. Выбор значения повлиял на результат незначительно, и характер остался тем же.
- Splitter — стратегия выбора разделения. Поскольку один из вариантов — это рандом, то итоговый результат получился прежним, но с некоторым шумом.
- Max_depth — максимальная глубину дерева решений. При увеличении глубины, разумеется, точность повышается, правда, если перестараться, то может возникнуть оверфиттинг.
- Min_samples_split — минимальное число наблюдений необходимых для разбиения узла дерева. С увеличением значения наблюдается ухудшение классификации.
- Min_samples_leaf — минимальное число наблюдений, требующееся для листа. Чем больше параметр, тем хуже классификация.