

В этом задании вы сможете немного расслабиться после реализации случайного леса и градиентного бустинга по мотивам материалов прошлой недели. Все, что нужно будет делать — запускать методы из `sklearn`. Вам предлагается выяснить, какое распределение лучше использовать в наивном байесовском классификаторе в зависимости от вида признаков.

Загрузите датасеты `digits` и `breast_cancer` из `sklearn.datasets`. Выведите несколько строчек из обучающих выборок и посмотрите на признаки. С помощью `sklearn.cross_validation.cross_val_score` с настройками по умолчанию и вызова метода `mean()` у возвращаемого этой функцией `numpy.ndarray`, сравните качество работы наивных байесовских классификаторов на этих двух датасетах. Для сравнения предлагается использовать `BernoulliNB`, `MultinomialNB` и `GaussianNB`. Насколько полученные результаты согласуются с рекомендациями из лекций?

Два датасета, конечно, еще не повод делать далеко идущие выводы, но при желании вы можете продолжить исследование на других выборках (например, из UCI репозитория).

Для сдачи задания, ответьте на приведенные ниже вопросы.

Вопрос 1

Каким получилось максимальное качество классификации на датасете `breast_cancer`?

Вопрос 2

Каким получилось максимальное качество классификации на датасете `digits`?

Вопрос 3

Выберите верные утверждения и запишите их номера через пробел (в порядке возрастания номера):

1) На вещественных признаках лучше всего сработал наивный байесовский классификатор с распределением Бернулли

2) На вещественных признаках лучше всего сработал наивный байесовский классификатор с мультиномиальным распределением

3) Мультиномиальное распределение лучше показало себя на выборке с целыми неотрицательными значениями признаков

4) На вещественных признаках лучше всего сработало нормальное распределение