

Размер случайного леса

Данное задание основано на материалах лекций по логическим методам и направлено на знакомство со случайными лесами (Random Forests).

Вы научитесь:

- работать со случайным лесом — одним из наиболее распространенных семейств алгоритмов
- решать с его помощью задачи регрессии
- подбирать параметры случайного леса

Введение

Случайный лес — это модель классификации, объединяющая некоторое количество решающих деревьев в одну композицию, за счет чего улучшается их качество работы и обобщающая способность. Деревья строятся независимо друг от друга. Чтобы они отличались друг от друга, обучение проводится не на всей обучающей выборке, а на ее случайном подмножестве. Также, для дальнейшего уменьшения схожести деревьев, оптимальный признак для разбиения выбирается не из всех возможных признаков, а лишь из их случайного подмножества. Прогнозы, выданные деревьями, объединяются в один ответ путем усреднения.

Особенность случайного леса заключается в том, что он не переобучается по мере увеличения количества деревьев в композиции. Это достигается за счет того, что деревья не зависят друг от друга, и поэтому добавление нового дерева в композицию не усложняет модель, а лишь понижает уровень шума в прогнозах.

Реализация в Scikit-Learn

В библиотеке `scikit-learn` случайные леса реализованы в классах `sklearn.ensemble.RandomForestClassifier` (для классификации) и `sklearn.ensemble.RandomForestRegressor` (для регрессии). Обучение модели производится с помощью функции `fit`, построение прогнозов — с помощью функции `predict`. Число деревьев задается с помощью поля класса `n_estimators`.

Пример использования:

```
import numpy as np
from sklearn.ensemble import RandomForestRegressor
X = np.array([[1, 2], [3, 4], [5, 6]])
y = np.array([-3, 1, 10])
clf = RandomForestRegressor(n_estimators=100)
clf.fit(X, y)
predictions = clf.predict(X)
```

Также в этом задании вам понадобится вычислять качество предсказаний на тестовой выборке. Мы будем пользоваться метрикой R^2 — по сути, это среднеквадратичная ошибка (MSE), нормированная на отрезок $[0, 1]$ и обращенная так, чтобы ее наилучшим значением была единица. Ее можно вычислить с помощью функции `sklearn.metrics.r2_score`. Первым аргументом является список правильных ответов на выборке, вторым — список предсказанных ответов. Пример использования:

```
from sklearn.metrics import r2_score
print r2_score([10, 11, 12], [9, 11, 12.1])
```

Инструкция по выполнению

В этом задании вам нужно проследить за изменением качества случайного леса в зависимости от количества деревьев в нем.

1. Загрузите данные из файла `abalone.csv`. Это датасет, в котором требуется предсказать возраст ракушки (число колец) по физическим измерениям.
2. Преобразуйте признак `Sex` в числовой: значение `F` должно перейти в `-1`, `I` — в `0`, `M` — в `1`. Если вы используете `Pandas`, то подойдет

следующий код: `data['Sex'] = data['Sex'].map(lambda x: 1 if x == 'M' else (-1 if x == 'F' else 0))`

3. Разделите содержимое файлов на признаки и целевую переменную. В последнем столбце записана целевая переменная, в остальных — признаки.
4. Обучите случайный лес (`sklearn.ensemble.RandomForestRegressor`) с различным числом деревьев: от 1 до 50 (`random_state=1`). Для каждого из вариантов оцените качество работы полученного леса на кросс-валидации по 5 блокам. Используйте параметры `"random_state=1"` и `"shuffle=True"` при создании генератора кросс-валидации `sklearn.cross_validation.KFold`. В качестве меры качества воспользуйтесь коэффициентом детерминации (`sklearn.metrics.r2_score`).
5. Определите, при каком минимальном количестве деревьев случайный лес показывает качество на кросс-валидации выше 0.52. Это количество и будет ответом на задание.
6. Обратите внимание на изменение качества по мере роста числа деревьев. Ухудшается ли оно?

Ответ на каждое задание — текстовый файл, содержащий ответ в первой строчке. Обратите внимание, что отправляемые файлы не должны содержать пустую строку в конце.