

Курс "Библиотеки Python для Data Science: продолжение"

Практическое задание к уроку 4

Тема "Оценка и интерпретация полученной модели"

1. Расскажите, как работает регуляризация в решающих деревьях, какие параметры мы штрафует в данных алгоритмах?

Регуляризация - это способ бороться с переобучением. Для добиваясь этого можно штрафовать слишком большую сложность модели.

Для деревьев решений основные способы борьбы с переобучением:

- наложить верхнее ограничение глубины или минимально числа объектов в узле
- стрижка дерева.

В XGBoost регуляризация производится через коэффициент λ в значаестве

похожести.

Если λ больше, то
похожесть в рамках одной
ветви становится меньше.

Если $\lambda = 0$, то похожесть
будет большой. В

дальнейшем, сравнимся
с γ (или

параметр γ Boost,
к-рый отвечает за
стрижку деревьев);

— если полученная идио
большая

этой γ , то мы
оставляем

это дерево

Т.е. регуляризация в XG
Boost используется при
подсчете похожести

2. По какому принципу рассчитывается "важность признака (feature_importance)" в ансамблях деревьев?

Признаки ранжируются в

соответствии с объясненной дисперсией, которую каждый признак вносит в модель. Признаки отбираются в зависимости от их относительной важности, то есть процентной важности от наиболее важного признака.