

Проектная работа по модулю “Препроцессинг данных”

1. Загрузить файл data_breast.csv. В данном файле собрана расчетная информация с обработанных изображений биоптата молочных желез женщин. Задача заключается в предсказании переменной “Diagnosis” - является ли содержимое биоптата доброкачественным (значение “B” – benign) либо злокачественным (значение “M” – malicious). Описание данных доступно на сайте <https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>
2. Рассчитать основные статистики для переменных (среднее, медиана, мода, мин/макс, сред. отклонение).
3. Выбрать стратегию для работы с пропущенными значениями.
4. Рассчитать и визуализировать корреляционную матрицу для переменных.
5. Визуализировать взаимосвязи между переменными.
6. С помощью статистических методов проверить взаимосвязи между переменными.
7. Выбрать стратегию Feature Selection – сокращение размерности либо генерация новых переменных. Какой из этих двух подходов даст лучший результат при классификации?
8. Провести стратегию Oversampling/Undersampling, проверить дает ли она улучшение результатов.
9. Сделать кросс-валидацию данных с использованием подхода K-fold (n_folds=10).
10. Рассчитать Feature Selection для выбранных переменных.
11. Решить задачу бинарной классификации и предсказать переменную “Diagnosis” протестировав как минимум 2 алгоритма. Использовать те алгоритмы, которые позволяют предсказать вероятность класса (proba). Рассчитать и вывести вероятность каждого класса.
12. Проверить качество классификации с использованием следующих метрик: Accuracy, F1-Score, Precision, Recall
13. Проверить качество вероятности класса с использованием метрики: Brier Score
14. * Осуществить запуск вашего скрипта с использованием Docker.
15. Загрузить результат (в формате .ipynb ноутбука либо докер реализации) в репозиторий. Разместить ссылку в лк