

Курсовой проект по дисциплине  
«Анализ данных и машинное обучение»

**Исходные данные.**

Исходные данные для выполнения курсового проекта представляют собой набор данных в формате csv. Набор данных содержит определенный набор признаков и целевой признак.

**Задание на курсовой проект.**

**1. Визуальный анализ данных.**

С использованием полученного набора данных необходимо провести визуальный анализ данных. Всего должно быть построено не менее 15 различных графиков. Обязательно должны быть использованы графики:

- Stackplot

С использованием pandas plot графики типов:

- pie;
- density;
- bar.

С использованием Seaborn:

- boxplot
- countplot
- distplot
- pairplot

Возможно использовать также любые графики из любых доступных библиотек.

Ссылка на хорошую статью с применением графиков: <https://habr.com/ru/post/468295/>

**2. Построение моделей машинного обучения.**

Используя полученный набор данных необходимо построить модели машинного обучения. Целевой признак отдельно указан для каждого набора данных. В качестве набора для обучения выбирается 80% данных, в качестве тестовой выборки 20% данных. Параметр random\_state при разбиении указывается 15.

Для полученного набора данных требуется построить следующие модели машинного обучения:

- дерево решений;
- метод ближайших соседей;
- логистическая регрессия;
- случайный лес;
- градиентный бустинг.

Обучение моделей провести в 3 этапа:

1. С минимальной подготовкой данных. Выполнить только преобразование значений только в случае если это требуется для обучения модели.

2. Выполнить чистку и предподготовку данных. Удалить «выбросы», заполнить недостающие значения и т.д. Удалить не более 5% данных.

3. Выполнить подбор гиперпараметров модели.

Для каждого из этапов записать результат работы модели на тестовой выборке используя метрики: accuracy, precision, recall, roc auc.

Jupyter Notebook с полным кодом анализа данных и построения моделей прилагается к курсовому проекту.

### Критерии оценивания курсового проекта.

**Отлично.** Выполнены все задачи курсового проектирования. Пояснительная записка написана грамотным языком и отражает все этапы выполнения курсового проекта.

**Хорошо.** Построено 15 графиков, но использовано только 5 из 8 обязательных. Применено 4 из 5 обязательных моделей машинного обучения, но выполнены все этапы обучения.

**Удовлетворительно.** Построено 10 графиков, но использовано только 3 из 8 обязательных. Применено 3 из 5 обязательных моделей машинного обучения, выполнено 2 из 3 этапов обучения.

### Варианты для выполнения курсового проекта.

Вариант для выполнения курсового проекта выбирается в соответствии с порядковым номером в списке группы.

№	Целевой признак	DataSet
1	income	<a href="https://www.kaggle.com/code/riccardobollati01/85-accuracy-ensemble-learning-model-data-viz/data">https://www.kaggle.com/code/riccardobollati01/85-accuracy-ensemble-learning-model-data-viz/data</a>
2	hospital_death	<a href="https://www.kaggle.com/datasets/mitishaagarwal/patient">https://www.kaggle.com/datasets/mitishaagarwal/patient</a>
3	class	<a href="https://www.kaggle.com/datasets/kukuroo3/body-performance-data">https://www.kaggle.com/datasets/kukuroo3/body-performance-data</a>
4	Target	<a href="https://www.kaggle.com/datasets/krantiswalke/bankfullcsv">https://www.kaggle.com/datasets/krantiswalke/bankfullcsv</a>
5	Star type	<a href="https://www.kaggle.com/datasets/deepu1109/star-dataset">https://www.kaggle.com/datasets/deepu1109/star-dataset</a>
6	is_safe	<a href="https://www.kaggle.com/datasets/mssmartypants/water-quality">https://www.kaggle.com/datasets/mssmartypants/water-quality</a>
7	depressed	<a href="https://www.kaggle.com/datasets/diegobabativa/depression">https://www.kaggle.com/datasets/diegobabativa/depression</a>
8	specie	<a href="https://www.kaggle.com/datasets/bertiemackie/sloth-species">https://www.kaggle.com/datasets/bertiemackie/sloth-species</a>
9	Survived	<a href="https://www.kaggle.com/datasets/yasserh/titanic-dataset">https://www.kaggle.com/datasets/yasserh/titanic-dataset</a>
10	label	<a href="https://www.kaggle.com/datasets/praveengovi/credit-risk-classification-dataset">https://www.kaggle.com/datasets/praveengovi/credit-risk-classification-dataset</a>
11	TravelInsurance	<a href="https://www.kaggle.com/datasets/tejashvi14/travel-insurance-prediction-data">https://www.kaggle.com/datasets/tejashvi14/travel-insurance-prediction-data</a>
12	satisfied	<a href="https://www.kaggle.com/datasets/mohamedharris/employee-satisfaction-index-dataset">https://www.kaggle.com/datasets/mohamedharris/employee-satisfaction-index-dataset</a>
13	final_payment_status	<a href="https://www.kaggle.com/datasets/vizdom/customer-orders">https://www.kaggle.com/datasets/vizdom/customer-orders</a>
14	status	<a href="https://www.kaggle.com/datasets/debasisdotcom/parkinson-disease-detection">https://www.kaggle.com/datasets/debasisdotcom/parkinson-disease-detection</a>
15	status	<a href="https://www.kaggle.com/datasets/qusaybtoush1990/banks-loan">https://www.kaggle.com/datasets/qusaybtoush1990/banks-loan</a>
16	target	<a href="https://www.kaggle.com/datasets/yasserh/heart-disease-dataset">https://www.kaggle.com/datasets/yasserh/heart-disease-dataset</a>
17	Class	<a href="https://www.kaggle.com/datasets/mssmartypants/rice-type-classification">https://www.kaggle.com/datasets/mssmartypants/rice-type-classification</a>
18	rating	<a href="https://www.kaggle.com/datasets/qusaybtoush1990/the-cars">https://www.kaggle.com/datasets/qusaybtoush1990/the-cars</a>
19	Your level of satisfaction in Online Education	<a href="https://www.kaggle.com/datasets/sujaradha/online-education-system-review">https://www.kaggle.com/datasets/sujaradha/online-education-system-review</a>
20	BAD	<a href="https://www.kaggle.com/datasets/ajay1735/hmeq-data">https://www.kaggle.com/datasets/ajay1735/hmeq-data</a>