

Университет ИТМО

Практическая работа №5  
по дисциплине «Визуализация и моделирование»

**Автор:** Никулина Ксения Григорьевна

**Поток:** ВИМ 1.1

**Группа:** К3221

**Факультет:** ИКТ

**Преподаватель:** Чернышева А.В.

Санкт-Петербург, 2021 г.

# 1 Heart Failure Prediction

Сердечно-сосудистые заболевания являются причиной смерти номер 1 во всем мире, забирая в среднем 17.9 миллионов жизней в год, что составляет 31% от всех смертей по всему миру. Остановка сердца - результат сердечно-сосудистых заболеваний и данный датасет содержит 12 показателей, которые можно использовать для предсказания смертности от остановки сердца.

Название столбца	Описание данных	Тип данных	Шкала
Возраст	Возраст пациента	Integer	Интерв
Анемия	Снижение гемоглобина	Binary	Номин
Креатинкиназа	Уровень фермента	Integer	Относит
Диабет	Наличие диабета	Binary	Номин
Фракция выброса	% крови при сокращении	Integer	Относит
Высокое кровяное давление	Гипертония	Binary	Номин
Тромбоциты	Тромбоциты в крови	Float	Номин
Креатинин	Уровень креатинина в крови	Float	Номин
Натрий	Уровень натрия в крови	Integer	Относит
Пол	Пол пациента	Binary	Номин
Курение	Курит ли пациент	Binary	Номин
Время	Кол-во дней наблюдения	Integer	Относит
Смерть	Умер ли пациент	Binary	Номин

## 2 Анализ датасета

Обучение с учителем (supervised learning) предполагает наличие полного набора размеченных данных для тренировки модели на всех этапах ее построения.

Наличие полностью размеченного датасета означает, что каждому примеру в обучающем наборе соответствует ответ, который алгоритм и должен получить. Таким образом, размеченный датасет из фотографий цветов обучит нейронную сеть, где изображены розы, ромашки или нарциссы. Когда сеть получит новое фото, она сравнит его с примерами из обучающего датасета, чтобы предсказать ответ.

Я считаю, мой датасет подходит для обучения с учителем, так как на вход он будет получать входные данные - показатели пациента, а на выход его шансы на смерть, основываясь на предыдущем опыте осматриваемых пациентов. Таким образом, можно будет предсказать шансы на выживание у поступившего на осмотр пациента. Данная задача является задачей регрессии. Отклик - действительное число 0 или 1 (0 - жив, 1 - смерть)

## 3 CRISP-DM анализ

### 1. Понимание бизнеса

- Формулировка глобальной задачи: предсказание смерти на этапе обследования пациента, предотвращение ранней смерти пациента
- Задача: предотвращение ранней смерти пациента

### 2. Понимание данных

- Данные: датасет Heart Failure Prediction
- Описание данных: таблица
- Исследование и обработка данных: см. лаб 2-4

### 3. Подготовка данных

- Обработка: см. лаб 4

### 4. Моделирование

```
on", "high_blood_pressure", "platelets", "serum_creatinine", "serum_sodium", "sex", "smoking"]].to_numpy(), df.DEATH_EVENT.to_numpy(), train_size = 0.8)

[45] reg = LinearRegression().fit(X_train, y_train)
reg.score(X_train, y_train) # коэффициент детерминации
# Для линейной зависимости коэффициент детерминации равен квадрату коэффициента
# корреляции rxy: R2 = rxy^2 . Например, значение R2 = 0.28, означает, что в 28%
# случаев изменения x приводят к изменению y . Другими словами, точность подбора уравнения регрессии - ниже средней.
0.2485462768255866

[46] reg.score(X_test, reg.predict(X_test))
1.0

[47] mean_squared_error(y_test, reg.predict(X_test), squared=False) # MSE
0.44797895521567466

[48] pd.DataFrame.from_records(data=[list(y_test), list(reg.predict(X_test))])
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	1.000000	0.000000	1.000000	0.000000	0.000000	1.000000	0.000000	0.000000	1.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
1	0.550227	0.349626	0.314086	0.789167	0.559804	0.723177	0.281116	0.567216	0.80751	0.802468	0.379183	0.981089	0.218089	0.490811	0.175558	0.36136

Рис. 1 Моделирование

### 5. Оценка

- Благодаря анализу при помощи линейной регрессии мы заметили, что изменение наших данных практически не влияют на изменение результата. Это может быть связано с тем, что количество разнообразных данных недостаточно, так как в датасете всего около 200 строк и для более точной информации требует более расширенный датасет.