

Задания лабораторных работ по курсу «Системы искусственного интеллекта»

ЛАБОРАТОРНАЯ РАБОТА 1	
Разведочный анализ данных (EDA) и подготовка к моделированию	
Цель	Научиться загружать данные, выявлять проблемы (пропуски, выбросы, типы данных), строить понятные визуализации, формулировать гипотезы и готовить датасет к моделированию.
Датасет	Использовать датасет согласно варианту (например, titanic.csv), структура столбцов стандартная: Survived, Pclass, Sex, Age, SibSp, Parch, Fare, Embarked (+ возможные доп. поля).
Форма отчетности	<ul style="list-style-type: none"> • Jupyter Notebook (.ipynb) с кодом, графиками и пояснениями. • Отчёт с основными выводами и скриншотами ключевых графиков.
Требования к оформлению	Читаемые графики (подписи осей, легенды, единицы измерения), пояснительный текст к каждому блоку, отсутствие «немых» ячеек.
Инструменты	Python, pandas, numpy, matplotlib, seaborn, scikit-learn
Частые ошибки	<ul style="list-style-type: none"> • пустые/неподписанные графики; • перечисление цифр без выводов; • бездумное удаление строк с пропусками (покажите альтернативы); • изменение типа Pclass на float без необходимости (категория!); • закодировали категории – и забыли сохранить mapping/one-hot колонки.
Оценка:	100 баллов
Срок сдачи	11 октября

СПИСОК РЕКОМЕНДОВАННЫХ ДАТАСЕТОВ		
ЕДА		
Грибы: съедобный или ядовитый	8124×22, бинарная	https://www.kaggle.com/datasets/uciml/mushroom-classification
Семена пшеницы	210×7, 3 класса	https://www.kaggle.com/datasets/muratkokludataset/seeds-dataset
Качество красного вина	1599×11, мультикласс	https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009
ДЕНЬГИ		
Доходы населения	48k×14, бинарная	https://www.kaggle.com/datasets/uciml/adult-census-income
Маркетинг банковских услуг	45k×16, бинарная	https://www.kaggle.com/datasets/henriqueyamahata/bank-marketing
Дефолты по кредитным картам	30k×23, бинарная	https://www.kaggle.com/datasets/uciml/default-of-credit-card-clients-dataset
Кредитный скоринг	1k×20, бинарная	https://www.kaggle.com/datasets/uciml/german-credit
Подлинность банкнот	1372×4, бинарная	https://www.kaggle.com/datasets/ritesaluja/banknote-authentication-uci-data
Спам-письма (Spambase)	4601×57, бинарная	https://www.kaggle.com/datasets/uciml/spambase-dataset
ЗДОРОВЬЕ		
Диабет	768×8, бинарная	https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database
Болезни сердца	303×13, бинарная	https://www.kaggle.com/datasets/cherngs/heart-disease-cleveland-uci
Рак груди	569×30, бинарная	https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data

Рубрика оценивания (на примере датасета «Титаник»)			
	Этап	Баллы	Содержание
1	Импорт и первичная диагностика	20	<p>1. Загрузка данных</p> <ul style="list-style-type: none"> ○ Прочитать CSV. ○ Показать первые/последние строки, размерность, типы столбцов. <p>2. Качество данных</p> <ul style="list-style-type: none"> ○ Посчитать число и долю пропусков по каждому признаку. ○ Найти дубликаты (если есть) и решить, что с ними делать. ○ Проверить корректность типов (int/float/category/datetime) и при необходимости привести к нужным. <p>В ноутбуке обязательно коротко объясняйте почему вы приняли то или иное решение (например, почему удалили дубликаты).</p>
2	Описательная статистика и однофакторный анализ	20	<p>3. Числовые признаки</p> <ul style="list-style-type: none"> ○ Для количественных столбцов (<code>Age</code>, <code>Fare</code>, ...) вывести: <code>count</code>, <code>mean</code>, <code>std</code>, <code>min</code>, <code>25%</code>, <code>50%</code>, <code>75%</code>, <code>max</code>. ○ Для каждого числового признака построить гистограмму (или KDE) и boxplot. ○ Коротко интерпретировать: где асимметрия, длинные хвосты, есть ли выбросы? <p>4. Категориальные признаки</p> <ul style="list-style-type: none"> ○ Частоты категорий для <code>Sex</code>, <code>Pclass</code>, <code>Embarked</code> (таблица + столбчатая диаграмма). ○ Комментарий: какие категории доминируют? есть ли несбалансированность?
3	Двухфакторный анализ	20	<p>5. Целевой признак и группы</p> <ul style="list-style-type: none"> ○ Исследовать связь <code>Survived</code> с ключевыми факторами: <ul style="list-style-type: none"> ▪ <code>Survived</code> VS <code>Sex</code> (stacked barplot/percentage barplot). ▪ <code>Survived</code> VS <code>Pclass</code>. ▪ <code>Survived</code> VS <code>Embarked</code>. ○ Написать 2–3 наблюдения (напр., «доля выживших у женщин существенно выше»). <p>6. Числовые признаки vs <code>Survived</code></p> <ul style="list-style-type: none"> ○ Сравнить распределения <code>Age</code> и <code>Fare</code> в группах <code>Survived=0/1</code> (boxplot/violinplot + краткий комментарий). ○ Простой статистический тест (например, U-критерий Манна–Уитни) с пояснением результата. <p>7. Корреляции</p> <ul style="list-style-type: none"> ○ Тепловая карта – только по числовым признакам; для категория–категория допустимы хи-квадрат/Cramér's V. ○ Коротко отметить 1–2 самые сильные связи (по модулю).
4	Подготовка данных	20	<p>8. Работа с пропусками</p> <ul style="list-style-type: none"> ○ Показать выбранную стратегию: удаление/импьютация (среднее/медиана/мода/по группам). ○ Обосновать выбор хотя бы для <code>Age</code> и <code>Embarked</code>. <p>9. Кодирование категорий</p> <ul style="list-style-type: none"> ○ Перевести <code>Sex</code>, <code>Embarked</code>, <code>Pclass</code> в числовой вид (Label/One-Hot Encoding). ○ Пояснить, почему выбран именно этот подход. <p>10. Масштабирование числовых признаков</p> <ul style="list-style-type: none"> ○ Применить <code>StandardScaler</code> или <code>MinMaxScaler</code>. ○ Показать, как изменились распределения/статистики. <p>Итог: получился чистый датасет + понятна структура признаков.</p>

5	Анализ выбросов	10	<p>Провести анализ выбросов по IQR</p> <ul style="list-style-type: none"> ○ Отсортируйте выборку. ○ Найдите квартели ○ Посчитайте межквартильный размах ○ Границы «усов» boxplot ○ Все точки ниже нижней или выше верхней границы считаются выбросами (outliers). <p>11. Не обязательно удалять; покажите влияние (до/после) и аргументируйте стратегию</p>
6	Результаты и гипотезы	10	<ul style="list-style-type: none"> • 5–7 пунктов: что важно для выживаемости, какие проблемы были в данных, какие решения приняли. • Сформулировать 2 гипотезы для следующего шага (моделирования). Примеры: <ul style="list-style-type: none"> ○ «Женщины выживают чаще мужчин при прочих равных». ○ «Пассажиры 1 класса имеют больше шансов на выживание, чем 3 класса». <p>«Увеличение тарифа (<code>Fare</code>) связано с большей вероятностью выживания».</p>

Мини чек-лист

- Показал размерность, типы колонок, пропуски, дубликаты
- Сделал гистограммы и boxplot для всех числовых
- Посчитал частоты и сделал barplot для категориальных
- Провёл двухфакторный анализ: таргет vs 2–3 признака (категориальные/числовые)
- Построил heatmap корреляций только для числовых
- Обосновал стратегию пропусков, кодирования и масштабирования
- Проанализировал выбросы по IQR и показал влияние (до/после), без обязательного удаления
- Сформулировал 5–7 выводов и 2 гипотезы
- Зафиксировал `random_state` и вывел версии библиотек

ЛАБОРАТОРНАЯ РАБОТА 2	
Задача классификации	
Цель	Построить и оценить модель классификации поверх очищенного датасета из ЛР-1; освоить базовые и продвинутые алгоритмы, метрики, интерпретацию, анализ ошибок.
Датасет	Использовать очищенный датасет из лабораторной работы 1
Форма отчетности	<ul style="list-style-type: none"> Jupyter Notebook (.ipynb) с кодом, графиками и пояснениями. Отчёт с основными выводами и скриншотами ключевых графиков.
Требования к оформлению	Читаемые графики (подписи осей, легенды, единицы), пояснительный текст к каждому блоку, отсутствие «немых» ячеек, фиксированный random state.
Инструменты	Python, pandas, numpy, matplotlib, seaborn, scikit-learn (по желанию: shap, xgboost / lightgbm / catboost)
Частые ошибки	<ul style="list-style-type: none"> Использование только accuracy (без F1/ROC-AUC/PR-AUC). Нестратифицированный train/valid/test, отсутствует фиксированный seed. Утечка таргета (препроцессинг по всему датасету или fit на test). One-Hot/Scaler обучены на всём датасете, а не только на train (нет Pipeline/ColumnTransformer). Отсутствует анализ ошибок.
Оценка:	100 баллов
Срок сдачи	8 ноября

Рубрика оценивания			
	Этап	Баллы	Содержание
1	Повторная загрузка и препроцессинг (Pipeline)	10	<ul style="list-style-type: none"> Возьмите очищенный датасет из ЛР-1. Соберите единый sklearn-pipeline (например, ColumnTransformer + StandardScaler для числовых + OneHotEncoder для категориальных). Все преобразования fit только на train (чтобы избежать утечки). Зафиксируйте список признаков и стратегию обработки пропусков (медиана/мода/константа – обосновать).
2	Разбиение и валидация	15	<ul style="list-style-type: none"> Стратифицированный train/valid/test = 60/20/20, random_state=42 (или заданный). Для дисбалансных задач — проверьте стратификацию. Допустимо CV (например, 5-fold) вместо фиксированного valid, но test оставить отложенным. Test-набор не используем до финального выбора модели и порога; все решения принимаются по train/valid (или CV). При CV используйте StratifiedKFold(shuffle=True, random_state=42).
3	Базовые модели	15	<ul style="list-style-type: none"> Обучите LogisticRegression и DecisionTreeClassifier (подберите max_depth). Отчёт по метрикам на valid: F1, ROC-AUC или PR-AUC (если дисбаланс заметный). Если доля позитивного класса < 30%, помимо F1 обязательно показать PR-AUC. Краткий вывод: что и почему лучше как baseline.
4	Продвинутая модель и поиск гиперпараметров	30	<ul style="list-style-type: none"> Одна из: CatBoost / LightGBM / XGBoost. Поиск 5–10 конфигураций (Grid/Random/SearchCV); сохраните лог результатов. Сравните с baseline, выберите «победителя».

5	Интерпретация и анализ ошибок	25	<ul style="list-style-type: none"> • Explainability: Permutation Importance (обязательно), SHAP – по желанию (если позволяет среда). • Error analysis: покажите минимум 3 FN и 3 FP с комментариями; сгруппируйте ошибки по одному признаку (например, sex/class/length_bin) и предложите 1–2 шага улучшения.
6	Репродуцируемость	5	<ul style="list-style-type: none"> • Зафиксировать random_state/seed. • Вывести версии библиотек (sys.version, pandas.__version__, sklearn.__version__). • Убедиться, что ноутбук запускается без ошибок.

Мини чек-лист

Данные и Pipeline

- Использую очищенный датасет из ЛР-1
- Все преобразования собраны в sklearn Pipeline/ColumnTransformer
- fit только на train (импютация/скейлер/one-hot) – утечек нет
- Список признаков и стратегия пропусков зафиксированы и описаны

Разбиение / валидация

- Stratified train/valid/test = 60/20/20, random_state=42
- При CV: stratifiedKFold(shuffle=True, random_state=42)
- Test не трогаю до финального выбора модели/порога

Базовые модели

- Обучены LogisticRegression и DecisionTree (подобран max_depth)
- Метрики на valid: F1 + ROC-AUC (и PR-AUC, если позитивов < 30%)
- Краткий вывод: какая baseline-модель лучше и почему

Продвинутая модель

- Одна из: CatBoost / LightGBM / XGBoost
- Поиск 5–10 конфигураций (Grid/Random/SearchCV) с логом результатов
- Сравнение с baseline — выбран «победитель»

Интерпретация и ошибки

- SHAP (summary + 1–2 waterfall) или permutation importance
- Показаны минимум 3 FN и 3 FP с кратким разбором причин и идеями улучшения

Репродуцируемость

- Зафиксировал random_state/seed
- Вывел версии библиотек
- Убедился, что ноутбук запускается без ошибок

ЛАБОРАТОРНАЯ РАБОТА 3

Задачи машинного обучения

Цель	Закрепить навыки применения методов машинного обучения. На выбор: <ul style="list-style-type: none"> • Построить и оценить модель регрессии. • Провести кластеризацию и интерпретировать её. • Построить простую рекомендательную систему.
Датасет	Любой открытый (по желанию можно взять тот же, что в ЛР-1/ЛР-2, если подходит).
Форма отчетности	<ul style="list-style-type: none"> • Jupyter Notebook (.ipynb) с кодом, графиками и пояснениями. • Отчёт с основными выводами и скриншотами ключевых графиков.
Требования к оформлению	читаемые графики (подписи осей, легенды, единицы), пояснительный текст к каждому блоку, отсутствие «немых» ячеек, фиксированный random state.
Инструменты	Python, pandas, numpy, matplotlib, seaborn, scikit-learn (по желанию: surprise, implicit, xgboost/lightgbm/catboost, umap, shap)
Частые ошибки	<ul style="list-style-type: none"> • Использование только одной метрики без пояснений. • Отсутствие baseline-модели. • Нет сравнения методов. • Отсутствие анализа ошибок или интерпретации. • Нет пояснений к визуализациям.
Оценка:	100 баллов
Срок сдачи	13 декабря

Рубрика оценивания

	Этап	Баллы	Содержание
1	Постановка задачи и данные	10	Кратко описать задачу (регрессия / кластеризация / рекомендации). Описание датасета, признаков и целевой переменной (если есть). Обоснование выбора.
2	Базовые методы	15	<ul style="list-style-type: none"> • Для регрессии: LinearRegression, Ridge/Lasso. • Для кластеризации: KMeans. • Для рекомендаций: baseline (например, популярность или user-average). Показать метрики/графики + краткий вывод.
3	Продвинутый метод	30	<ul style="list-style-type: none"> • Для регрессии: GradientBoosting/RandomForest/ XGBoost. • Для кластеризации: DBSCAN/Hierarchical. • Для рекомендаций: ALS/Matrix Factorization/Content-based. Сравнить с baseline, выбрать «победителя».
4	Метрики и сравнение	15	<ul style="list-style-type: none"> • Для регрессии: MSE/RMSE, MAE, R². • Для кластеризации: silhouette, inertia, ARI. • Для рекомендаций: Precision@10, Recall@10, MAP (или NDCG@10 по желанию). Краткий вывод о качестве.
5	Интерпретация и визуализация	15	Показать важность признаков (Permutation Importance / SHAP) или визуализацию (кластеров/сниженного пространства/топ-N рекомендаций). Объяснить результаты словами.
6	Анализ ошибок или результатов	10	<ul style="list-style-type: none"> • Для регрессии: примеры с наибольшими ошибками. • Для кластеризации: «неудачные» точки. • Для рекомендаций: примеры неудачных рекомендаций. Предложить идеи улучшения.
7	Репродуцируемость	5	Зафиксировать seed, вывести версии библиотек. Код запускается без ошибок.

Мини чек-лист

- Выбрал тип задачи (регрессия / кластеризация / рекомендации).
- Описал датасет и задачу.
- Реализовал baseline-метод.
- Реализовал продвинутый метод.
- Посчитал метрики, сделал сравнение.
- Добавил интерпретацию/визуализацию.
- Провёл анализ ошибок или ограничений.
- Зафиксировал seed, проверил воспроизводимость.