

**Проектная работа по курсу**  
**«Системы искусственного интеллекта»**  
**на тему «Анализ данных социальных сетей»**

**ЭТАП 1**

**Разведочный анализ данных (EDA) и подготовка к моделированию на реальных данных**

<b>Цель</b>	<ul style="list-style-type: none"> <li>Научиться работать с «живыми» данными соцсетей: находить ботов, определять качество данных, выявлять закономерности в самораскрытии пользователей.</li> <li>Подготовить датасет для последующего анализа (кластеризация, классификация).</li> </ul>
<b>Датасет</b>	<ul style="list-style-type: none"> <li>Основной: VK dataset</li> <li>Дополнительно: 2 датасета реальных пользователей (550 аккаунтов с метаданными).</li> </ul>
<b>Задачи</b>	<ol style="list-style-type: none"> <li><b>Загрузка и первичный осмотр</b> <ul style="list-style-type: none"> <li>Вывести размерность, типы признаков.</li> <li>Посчитать долю пропусков по каждому признаку.</li> <li>Найти дубликаты и аномалии.</li> </ul> </li> <li><b>Удаление ботов / «пустых» профилей</b> <ul style="list-style-type: none"> <li>Предложить критерии (отсутствие друзей, фото, постов, подписок).</li> <li>Построить визуализации: распределение «живых» и «пустых» профилей.</li> <li>Сделать обоснованный выбор стратегии фильтрации (правила, кластеры, аномалии).</li> </ul> </li> <li><b>Анализ заполненности профиля</b> <ul style="list-style-type: none"> <li>Построить распределения по ключевым полям (образование, возраст, семейное положение, число фото, друзей).</li> <li>Сравнить мужчин и женщин: кто чаще указывает образование, семейное положение, кто активнее выкладывает фото.</li> <li>Сформулировать и проверить гипотезы, например:           <ul style="list-style-type: none"> <li>«Люди с образованием имеют более заполненный профиль»</li> <li>«Люди с большим числом фото чаще пишут посты»</li> </ul> </li> </ul> </li> <li><b>Гендерные различия в самораскрытии</b> <ul style="list-style-type: none"> <li>Провести дисперсионный или регрессионный анализ: зависят ли показатели активности от пола и образования.</li> <li>Визуализировать результаты (boxplot, barplot, heatmap).</li> </ul> </li> <li><b>Подготовка данных для моделирования</b> <ul style="list-style-type: none"> <li>Очистить признаки (убрать нерелевантные, обработать пропуски, привести категории к единому виду).</li> <li>Создать новые признаки (например, «плотность профиля» = число заполненных полей / общее число полей).</li> <li>Подготовить итоговый датасет для следующих этапов.</li> </ul> </li> </ol>
<b>Форма отчетности</b>	<ul style="list-style-type: none"> <li>Jupyter Notebook (.ipynb) с кодом, графиками и пояснениями.</li> <li>Отчёт с основными выводами и скриншотами ключевых графиков.</li> </ul>
<b>Требования к оформлению</b>	<ul style="list-style-type: none"> <li>Графики с подписями, легендами, читаемыми осями.</li> <li>Каждый блок сопровождается кратким пояснением.</li> <li>Должны быть сформулированы минимум <b>3 гипотезы</b> о поведении пользователей.</li> </ul>
<b>Инструменты</b>	Python, pandas, numpy, matplotlib, seaborn, scikit-learn, statsmodels.
<b>Частые ошибки</b>	<ul style="list-style-type: none"> <li>«Механическое» удаление строк с пропусками.</li> <li>Подсчёт статистик без пояснения.</li> <li>Неполные/неподписанные графики.</li> <li>Нет попытки сформулировать гипотезы.</li> </ul>
<b>Оценка</b>	<b>100 баллов</b>
<b>Срок сдачи</b>	<b>18 октября</b>

## ЭТАП 2

### Статистический и факторный анализ профилей пользователей социальных сетей

<b>Цель</b>	<ul style="list-style-type: none"> <li>• Освоить методы статистического анализа для поиска закономерностей в поведении пользователей.</li> <li>• Научиться выявлять скрытые факторы и строить кластеры пользователей по паттернам активности.</li> <li>• Проверить гипотезы о взаимосвязях между характеристиками профиля.</li> </ul>
<b>Датасет</b>	<ul style="list-style-type: none"> <li>• Подготовленный на Этапе 1 очищенный VK dataset.</li> <li>• Дополнительно: 2 датасета реальных пользователей (550 аккаунтов с метаданными).</li> </ul>
<b>Задачи</b>	<ol style="list-style-type: none"> <li><b>Дисперсионный анализ (ANOVA)</b> <ul style="list-style-type: none"> <li>◦ Проверить различия в заполненности профиля по полу, возрасту, образованию.</li> <li>◦ Сформулировать статистически подтверждённые выводы.</li> </ul> </li> <li><b>Регрессионный анализ</b> <ul style="list-style-type: none"> <li>◦ Построить модели зависимости активности пользователя (число постов, фото, подписчиков) от признаков профиля.</li> <li>◦ Интерпретировать коэффициенты регрессии: какие признаки вносят наибольший вклад.</li> </ul> </li> <li><b>Факторный анализ</b> <ul style="list-style-type: none"> <li>◦ Выделить скрытые факторы самораскрытия (например: «социальность», «медиа-активность», «академичность»).</li> <li>◦ Интерпретировать каждый фактор и дать названия.</li> </ul> </li> <li><b>Кластерный анализ</b> <ul style="list-style-type: none"> <li>◦ Разделить пользователей на группы по поведенческим признакам.</li> <li>◦ Визуализировать кластеры (heatmap, PCA/UMAP проекция).</li> <li>◦ Дать характеристику каждому кластеру («активные постеры», «фото-ориентированные», «минималисты» и т.п.).</li> </ul> </li> <li><b>Проверка гипотез</b> <ul style="list-style-type: none"> <li>◦ Найти и подтвердить/опровергнуть не менее 2–3 интересных взаимосвязей.</li> <li>◦ Примеры: <ul style="list-style-type: none"> <li>▪ «Люди, указывающие образование, на X% чаще заполняют семейное положение».</li> <li>▪ «Пользователи с большим числом фото публикуют больше постов».</li> </ul> </li> </ul> </li> <li><b>Обогащение датасета (по желанию, бонус)</b> <ul style="list-style-type: none"> <li>◦ Найти и подключить дополнительные данные (например, поле «образование», «город» или «работа» из VK API либо вспомогательных источников).</li> <li>◦ Описать, как это влияет на полноту анализа и качество проверяемых гипотез.</li> <li>◦ Примеры: <ul style="list-style-type: none"> <li>▪ добавить образование и проверить, как оно связано с самораскрытием;</li> <li>▪ добавить город и проверить, есть ли различия между мегаполисами и малыми городами.</li> </ul> </li> </ul> </li> </ol>
<b>Форма отчетности</b>	<ul style="list-style-type: none"> <li>• Jupyter Notebook (.ipynb) с кодом, графиками и пояснениями.</li> <li>• Отчёт с основными выводами и скриншотами ключевых графиков.</li> </ul>
<b>Требования к оформлению</b>	<ul style="list-style-type: none"> <li>• Каждая модель/тест сопровождается визуализацией и текстовой интерпретацией.</li> <li>• Выводы должны быть связаны с поставленными гипотезами.</li> <li>• Минимум 2–3 статистических вывода и 1 интерпретация кластеров.</li> </ul>
<b>Инструменты</b>	Python, pandas, numpy, matplotlib, seaborn, scikit-learn, statsmodels.
<b>Частые ошибки</b>	<ul style="list-style-type: none"> <li>• Проведение тестов без проверки условий применимости.</li> <li>• Интерпретация коэффициентов без учёта значимости.</li> <li>• Формальные кластеры без попытки их осмыслить.</li> <li>• Отсутствие гипотез или формулировка без численных доказательств.</li> </ul>
<b>Оценка</b>	<b>100 баллов</b>
<b>Срок сдачи</b>	<b>8 ноября</b>

### ЭТАП 3

#### Построение среднего профиля пользователя и сравнение с результатами анализа

<b>Цель</b>	<ul style="list-style-type: none"> <li>Сформировать «типичный» профиль пользователя по ключевым признакам.</li> <li>Сравнить усреднённые профили разных групп (мужчины/женщины, по образованию, возрасту).</li> <li>Сопоставить полученные средние профили с результатами статистического и кластерного анализа.</li> </ul>
<b>Датасет</b>	<ul style="list-style-type: none"> <li>Подготовленный очищенный датасет из Этапа 1.</li> <li>Результаты анализа и кластеризации из Этапа 2.</li> </ul>
<b>Задачи</b>	<ol style="list-style-type: none"> <li><b>Построение усреднённых профилей</b> <ul style="list-style-type: none"> <li>Рассчитать «средний профиль» отдельно для мужчин и женщин.</li> <li>По желанию: построить средний профиль по типу образования (технари, гуманитарии, естественники) и по возрастным группам (18–25, 26–40, 40+).</li> </ul> </li> <li><b>Сравнение с кластеризацией (Этап 2)</b> <ul style="list-style-type: none"> <li>Проверить: совпадают ли усреднённые профили с «центрами кластеров».</li> <li>Если есть различия — интерпретировать (например: «Средний женский профиль соответствует кластеру 2, но кластер 2 делится на активных и неактивных»).</li> </ul> </li> <li><b>Интерпретация профилей</b> <ul style="list-style-type: none"> <li>Описать в словах полученные средние портреты (например: «типичный мужчина 18–25 лет указывает образование, редко семейное положение, имеет 50–70 друзей и мало фото»).</li> <li>Сравнить с гипотезами, сформулированными ранее.</li> </ul> </li> <li><b>Визуализация</b> <ul style="list-style-type: none"> <li>Построить графики (барплоты, boxplot, spider/radar chart) для наглядного сравнения профилей.</li> <li>Сделать минимум 2 визуализации для разных групп.</li> </ul> </li> </ol>
<b>Форма отчетности</b>	<ul style="list-style-type: none"> <li>Jupyter Notebook (.ipynb) с кодом, графиками и пояснениями.</li> <li>Отчёт с основными выводами и скриншотами ключевых графиков.</li> </ul>
<b>Требования к оформлению</b>	<ul style="list-style-type: none"> <li>У профилей должны быть как числовые характеристики, так и словесная интерпретация.</li> <li>Сравнение должно содержать минимум 2–3 выводов.</li> <li>Визуализации должны быть читаемыми (подписи, легенды).</li> </ul>
<b>Инструменты</b>	Python, pandas, numpy, matplotlib, seaborn, scikit-learn, statsmodels.
<b>Частые ошибки</b>	<ul style="list-style-type: none"> <li>Формальное усреднение без интерпретации («средний возраст = 23.7» без комментариев).</li> <li>Игнорирование сравнения с результатами предыдущего этапа.</li> <li>Отсутствие визуализаций.</li> <li>Слишком общие выводы без цифр.</li> </ul>
<b>Оценка</b>	<b>100 баллов</b>
<b>Срок сдачи</b>	<b>29 ноября</b>

## ЭТАП 4

### Разработка программного продукта для анализа профиля пользователя

<b>Цель</b>	<ul style="list-style-type: none"> <li>Создать инструмент, который позволяет автоматически анализировать профиль пользователя ВКонтакте.</li> <li>Реализовать инструмент, доступный через интерфейс (бот или скрипт), который сравнивает конкретный профиль со «средним» и оценивает его особенности.</li> <li>Ввести критерий «удалённости профиля от среднего».</li> </ul>
<b>Датасет</b>	<ul style="list-style-type: none"> <li>Подготовленный и очищенный датасет (Этап 1).</li> <li>Средние профили пользователей (Этап 3).</li> </ul>
<b>Задачи</b>	<ol style="list-style-type: none"> <li><b>Разработка функционала анализа профиля</b> <ul style="list-style-type: none"> <li>Реализовать функцию/бот-команду <code>analyze_profile(id)</code>.</li> <li>Должно выполняться: <ul style="list-style-type: none"> <li>вывод заполненных полей профиля,</li> <li>сравнение с «средним» профилем,</li> <li>отображение отличий.</li> </ul> </li> </ul> </li> <li><b>Критерий «удалённости» от среднего профиля</b> <ul style="list-style-type: none"> <li>Определить и реализовать метрику (например, % совпадения, косинусное расстояние).</li> <li>Выводить её при анализе.</li> </ul> </li> <li><b>Интерфейс взаимодействия</b> <ul style="list-style-type: none"> <li>Вариант А: Jupyter Notebook или скрипт с функцией и текстовым выводом.</li> <li>Вариант В (для продвинутых): Telegram-бот с командами <code>/analyze &lt;id&gt;</code>, <code>/avg</code>, кнопками и визуализациями.</li> </ul> </li> <li><b>Визуализация результатов</b> <ul style="list-style-type: none"> <li>Добавить хотя бы один способ графической демонстрации различий: <ul style="list-style-type: none"> <li><code>radar chart</code> (сравнение с «средним»),</li> <li>таблица совпадений/отличий.</li> </ul> </li> </ul> </li> <li><b>Документация и тестирование</b> <ul style="list-style-type: none"> <li>Описать работу функции/бота.</li> <li>Проверить на нескольких примерах профилей.</li> <li>Сделать выводы о том, какие профили ближе или дальше от «среднего».</li> </ul> </li> </ol>
<b>Форма отчетности</b>	<ul style="list-style-type: none"> <li>Код (Notebook / скрипт / бот).</li> <li>Отчёт с основными выводами и скриншотами ключевых графиков.</li> </ul>
<b>Требования к оформлению</b>	<ul style="list-style-type: none"> <li>Код должен быть читаемым и воспроизводимым.</li> <li>Все шаги сопровождены пояснениями.</li> <li>Результаты визуализированы.</li> </ul>
<b>Инструменты</b>	Python, pandas, numpy, matplotlib, seaborn, scikit-learn. (По желанию: Streamlit/Flask для интерфейса).
<b>Частые ошибки</b>	<ul style="list-style-type: none"> <li>Формальная реализация без проверки на реальных профилях.</li> <li>Нет визуализации различий.</li> <li>Отсутствие метрики «удалённости» или её некорректная интерпретация.</li> <li>Неполная документация к коду.</li> </ul>
<b>Оценка</b>	<b>100 баллов</b>
<b>Срок сдачи</b>	<b>20 декабря</b>