

**Программа учебной дисциплины  
«Извлечение и анализ интернет-данных»**

Утверждена  
Академическим советом ОП  
Протокол № \_\_\_\_\_ от \_\_\_\_\_. \_\_\_\_\_. 20\_\_\_\_\_

Разработчики	Голубев Илья, Департамент больших данных и информационного поиска
Число кредитов	3
Контактная работа (час.)	32
Самостоятельная работа (час.)	82
Курс, Образовательная программа	2-4 (Б), Экономика и статистика
Формат изучения дисциплины	Без использования онлайн курса

**1. Цель, результаты освоения дисциплины и пререквизиты**

Цели:

1. Ознакомление студентов с основными способами извлечения информации из интернета и эффективного анализа этой информации
2. Формирование у студентов практических навыков анализа и извлечения данных и работы с ними

Планируемые результаты обучения (ПРО):

1. 1) Знать основные языковые конструкции и типы данных языка python
2. 2) Владение инструментарием pandas: уметь работать с табличными данными средствами python
3. 3) Знание основных типов графиков и инструментов визуализации для python. Уметь изобразить гистограмму, диаграмму рассеяния, поточечный график. Уметь добавлять описание графика
4. 4) Понимать внутреннюю структуру форматов xml/json/html. Уметь обрабатывать файлы таких форматов с помощью модулей python: json, BeautifulSoup и др.
5. 5) Владеть инструментами python для доступа к web. Библиотека request.
6. 6) Иметь представление о модулях python для статистического анализа и машинного обучения
7. 7) Научиться пользоваться документацией языка и его библиотек

Пререквизиты:

1. MOOC «Введение в Питон»
2. Основы теории вероятностей и статистики
3. Знаниями математики в объеме программы средней

## 2. Содержание учебной дисциплины

Тема (раздел дисциплины)	Объем в часах	Планируемые результаты обучения (ПРО), подлежащие контролю	Формы контроля
	лк		
	см		
	онл/ср		
Основы анализа данных в python	4	№: 1, 2, 7.	
	2		
	15		
Визуализация данных в python: библиотеки matplotlib, seaborn, plotly . Продвинутое инструменты для анализа данных	4	№: 3, 7.	
	2		
	15		
Парсинг открытых данных в различных форматах (xml/json/html)	2	№: 4, 7.	
	4		
	15		
Основы машинного обучения и практика применения	2	№: 6, 7.	
	4		
	15		
Извлечение и последующий анализ данных с ресурсов Википедия, Яндекс.Погода, Instagram	4	№: 5, 7.	
	6		
	20		
<b>Часов по видам учебных занятий:</b>	16		
	18		
	80		
<b>Итого часов:</b>	114		

### Содержание разделов дисциплины:

1. **Основы анализа данных в python**  
Повторение основных функций и объектов языка Python. Обзор библиотек numpy, pandas на основе данных из соревнований платформы kaggle.com.
2. **Визуализация данных в python: библиотеки matplotlib, seaborn, plotly . Продвинутое инструменты для анализа данных**  
Введение в визуальный анализ данных. Построение графиков, гистограмм, тепловых карт. Знакомство с порталом Открытых данных.

3. **Парсинг открытых данных в различных форматах (xml/json/html)**  
Изучение языков и библиотек для работы с xml/json/html
4. **Основы машинного обучения и практика применения**  
Основные термины, понятия и алгоритмы машинного обучения.  
Обсуждение самых популярных систем, основанных на машинном обучении (распознавание изображений, поиск, диалоговые системы).  
Алгоритмы машинного обучения: линейная и логистическая регрессии, градиентный бустинг и нейронные сети.
5. **Извлечение и последующий анализ данных с ресурсов Википедия, Яндекс.Погода, Instagram**  
Извлечение данных с перечисленных ресурсов, их последующий анализ и обработка. Визуализация полученных данных.

### 3. Оценивание

- **СЕМ1**, Не блокирующее, Работа на семинаре  
Основная функциональность Питона. Задачи в классе
- **ДЗ1**, Не блокирующее, Домашнее задание  
Основная функциональность Питона. Самостоятельная работа
- **СЕМ2**, Не блокирующее, Работа на семинаре  
Работа с текстами. Задачи в классе
- **ДЗ2**, Не блокирующее, Домашнее задание  
Работа с текстами. Самостоятельная работа
- **СЕМ3**, Не блокирующее, Работа на семинаре  
Извлечение данных с ресурсов. Задачи в классе
- **ДЗ3**, Не блокирующее, Домашнее задание  
Извлечение данных с ресурсов. Самостоятельная работа
- **СЕМ4**, Не блокирующее, Работа на семинаре  
Продвинутое извлечение данных с ресурсов. Задачи в классе
- **ДЗ4**, Не блокирующее, Домашнее задание  
Продвинутое извлечение данных с ресурсов. Самостоятельная работа
- **СЕМ5**, Не блокирующее, Работа на семинаре  
Извлечение и визуализация данных. Задачи в классе
- **ДЗ5**, Не блокирующее, Домашнее задание  
Извлечение и визуализация данных. Самостоятельная работа
- **ЭКЗ1**, Не блокирующее, Экзамен (письменный)  
Письменный экзамен: задачи и извлечение из популярных ресурсов и последующий анализ данных

**Формула округления:** Стандартное арифметическое округление

**Вид формулы оценивания:** Линейная

**Формула оценивания:**

Окончательная оценка = Округление( $0.35 * \text{ДЗ} + 0.35 * \text{Семинары} + 0.3 * \text{Экзамен}$ )

## 4. Примеры оценочных средств

### 5. Ресурсы

#### 5.1. Рекомендуемая основная литература

п/п	Наименование
1	<a href="#">Документация языка python</a>
2	<a href="#">Data Analysis in Python</a>
3	<a href="#">Документация Pandas</a>
4	<a href="#">Документация matplotlib</a>
5	<a href="#">Примеры визуализации в python</a>
6	<a href="#">Документация BeautifulSoup</a>
7	<a href="#">Мария Мансурова Web Scraping с помощью python</a>
8	<a href="#">Документация scikit-learn</a>
9	<a href="#">Статистические методы в python</a>

#### 5.2. Рекомендуемая дополнительная литература

*Не требуется*

#### 5.3. Программное обеспечение

п/п	Наименование	Условия доступа/скачивания
1	Microsoft Windows 7 Professional RUS Microsoft Windows 8.1 Professional RUS Microsoft Windows 10	<i>Из внутренней сети университета (договор)</i>
2	Microsoft Office Professional Plus 2010	<i>Из внутренней сети университета (договор)</i>
3	Anaconda - это бесплатный дистрибутив языка программирования Python для научных вычислений с открытым исходным кодом,	Бесплатно с <a href="https://www.anaconda.com/">https://www.anaconda.com/</a>

#### 5.4. Профессиональные базы данных, информационные справочные системы, интернет-ресурсы (электронные образовательные ресурсы)

п/п	Наименование	Условия доступа/скачивания
	<b><i>Профессиональные базы данных, информационно-справочные системы</i></b>	
1	Электронно-библиотечная система Юрайт	URL: <a href="https://biblio-online.ru/">https://biblio-online.ru/</a>
	<b><i>Интернет-ресурсы (электронные образовательные ресурсы)</i></b>	
1	Открытое образование	URL: <a href="https://openedu.ru/">https://openedu.ru/</a>

#### 5.5. Материально-техническое обеспечение дисциплины

Учебные аудитории для лекционных по дисциплине обеспечивают использование и демонстрацию тематических иллюстраций, соответствующих программе дисциплины в составе:

- ПЭВМ с доступом в Интернет (операционная система, офисные программы, антивирусные программы);

- мультимедийный проектор с дистанционным управлением.

Учебные аудитории для семинарских и самостоятельных занятий по дисциплине оснащены ПЭВМ, с возможностью подключения к сети Интернет и доступом к электронной информационно-образовательной среде НИУ ВШЭ.

## **6. Особенности организации обучения для лиц с ограниченными возможностями здоровья и инвалидов**

В случае необходимости, обучающимся из числа лиц с ограниченными возможностями здоровья (по заявлению обучающегося) а для инвалидов также в соответствии с индивидуальной программой реабилитации инвалида, могут предлагаться следующие варианты восприятия учебной информации с учетом их индивидуальных психофизических особенностей, в том числе с применением электронного обучения и дистанционных технологий:

6.1.1. *для лиц с нарушениями зрения:* в печатной форме увеличенным шрифтом; в форме электронного документа; в форме аудиофайла (перевод учебных материалов в аудиоформат); в печатной форме на языке Брайля; индивидуальные консультации с привлечением тифлосурдопереводчика; индивидуальные задания и консультации.

6.1.2. *для лиц с нарушениями слуха:* в печатной форме; в форме электронного документа; видеоматериалы с субтитрами; индивидуальные консультации с привлечением сурдопереводчика; индивидуальные задания и консультации.

6.1.3. *для лиц с нарушениями опорно-двигательного аппарата:* в печатной форме; в форме электронного документа; в форме аудиофайла; индивидуальные задания и консультации.