

## Аннотация рабочей программы дисциплины

### Б1.В.12 «Технологии обработки больших данных»

Курс 4 Семестр 7 Количество з.е. 4

**Объем трудоемкости:** 4 зачетных единиц (144 ч., из них – 68 час. аудиторной нагрузки: лекционных 34ч., лабораторных работ - 34 ч., 36 часов самостоятельной работы, 4 часа КСР, 0,3 часа ИКР.), форма контроля – экзамен.

**Цель освоения дисциплины** Цель дисциплины - Изучение принципов обработки больших данных, технологий распределенных вычислений, облачных платформ и инструментов анализа данных.

#### Задачи дисциплины:

- Изучение архитектурных решений для работы с Big Data.
- Освоение методов обработки структурированных и неструктурированных данных.
- Применение распределенных вычислений (Hadoop, Spark).
- Разработка алгоритмов анализа данных в распределенных средах.
- Использование облачных платформ для обработки больших данных.

#### Требования к знаниям и навыкам:

- Умение работать с распределенными системами (Hadoop, Spark).
- Опыт обработки данных в облачных средах (Yandex Cloud).
- Навыки анализа данных с помощью Python (Pandas, PySpark, Dask).
- Понимание архитектуры Big Data-решений.

#### Место дисциплины в структуре ООП ВО:

Дисциплина «Технологии обработки больших данных» относится к Блок 1 Дисциплины, часть, формируемая участниками образовательного процесса.

Дисциплина изучается в 7-м семестре. Для успешного освоения необходимы знания, полученные в дисциплинах: «Алгебра и введение в тензорный анализ», «Теория вероятностей и математическая статистика», «Подготовка данных машинного обучения», «Технологии управления данными NoSQL», «Многомерный статистический анализ», и «Машинное обучение», «Программирование».

Преподавание ведется в виде лекций и лабораторных занятий с использованием интерактивных методов. Лабораторные работы направлены на практическое освоение методов и инструментов классификации на реальных данных.

Дисциплина формирует компетенции, необходимые для выполнения выпускной квалификационной работы и профессиональной деятельности в области вычислительных технологий.

#### Результаты обучения (знания, умения, опыт, компетенции):

BD-3	Способен организовывать хранения данных, выбирая адекватные технологические решения
------	---

BD-3.1	<p>Разрабатывает, отлаживает и тестирует прикладные решения с элементами ИИ с применением различных технологий хранения структурированных данных, оценивает качество.</p> <p>Пишет аналитические запросы к данным и анализирует план запроса. Умеет создавать представления, хранимые процедуры, функции и триггеры.</p> <p>Знание типов СУБД: реляционные, NoSQL, колоночные, документные - Архитектура распределенных систем хранения - Принципы ACID, CAP-теорема - Методы индексации и партиционирования - Принципы транзакционности</p> <p>Уметь: Выбирать тип СУБД под задачи ИИ - Проектировать схемы данных для ML-моделей - Анализировать и оптимизировать планы запросов - Разрабатывать сложные аналитические запросы - Создавать database objects (views, procedures, functions)</p> <p>Владение SQL (оконные функции, CTE, сложные джойны) - Навык чтения explain plan - Оптимизация запросов через индексы - Создание хранимых процедур для ETL - Работа с триггерами для поддержания целостности</p>
<b>BD-4</b>	<b>Способен применять различные модели и (или) технологии обработки данных</b>
BD-4.1	<p>Осуществляет выбор технологий обработки больших данных, приемлемых для создания прикладной системы ИИ с заданными требованиями</p> <p>Способен организовывать распределенное хранилище и параллельную обработку на базе современных технологий (Hadoop, Spark) больших данных:</p> <p>Знать: Архитектуру Hadoop ecosystem (HDFS, YARN) - Модели вычислений: MapReduce, DAG - Принципы RDD и DataFrame в Spark – принципы Стриминговой обработки vs batch processing - Паттерны Lambda/Карра архитектур</p> <p>Уметь: Выбирать стек технологий под задачи ИИ - Проектировать распределенные ETL-пайплайны - Оптимизировать производительность Spark-приложений - Организовывать шардирование и репликацию данных</p> <p>Владеть PySpark API - Настройка и администрирование Hadoop/Spark кластеров - Оптимизация через партиционирование, кэширование - Мониторинг производительности распределенных систем</p>
<b>BD-5</b>	<b>Способен применять технологии организации инфраструктуры БД</b>
BD-5.1	<p>Осуществляет выбор направления вспомогательных технологических решений для формирования единого стека работы с большими данными для решения поставленной задачи. Руководит проектами по организации инфраструктуры БД:</p> <p>Знать: принципы построения data lake, data warehouse - Методы оркестрации пайплайнов (Airflow, Prefect) - CI/CD для data projects - Мониторинг и observability в data-системах - Методологии управления data projects.</p> <p>Уметь: Формировать единый технологический стек - Управлять жизненным циклом data infrastructure - Выбирать инструменты</p>

	мониторинга и оркестрации - Оценивать TCO инфраструктурных решений. Владеть навыками проектного управления в data-проектах - Составление ТЗ на инфраструктуру - Ведение технической документации - Оценка рисков инфраструктурных решений.
<b>ML-1</b>	<b>Способен применять знания об истории развития и трендах современного ИИ для формулирования корректных постановок задач и поиска перспективных способов решения проблем с помощью ИИ</b>
ML-1.2	Определяет тенденции развития, оценивает новизну и практическую значимость своих решений с точки зрения современного искусственного интеллекта. Проектирует и внедряет комплексные пайплайны предварительной обработки данных с использованием современных методов ИИ, автоматизации и feature engineering в различных предметных областях. Знать: современные архитектуры ML-моделей (трансформеры, GAN, RL) - Методы feature engineering и feature selection - MLOps принципы и best practices - Современные фреймворки автоматического ML - Тренды в области ИИ (LLM, мультимодальные модели) Уметь: Проектировать end-to-end ML пайплайны - Применять автоматизированный feature engineering - Оценивать бизнес-ценность ML-решений - Адаптировать state-of-the-art подходы под задачи - Владение MLflow, Kubeflow - Создание воспроизводимых ML-экспериментов - Автоматизация пайплайнов предобработки - A/B тестирование ML-моделей
<b>PL-1</b>	<b>Способен применять язык программирования Python для решения задач в области ИИ</b>
PL-1.3	Разрабатывает и поддерживает системы обработки больших данных различной степени сложности. Способен строить архитектуру вычислений с использованием cloud-native инструментов, в том числе бессерверных решений (Yandex Cloud Functions): Знать: архитектурные паттерны big data систем - Принципы serverless computing - Cloud-native подходы (контейнеризация, orchestration) - Асинхронное программирование в Python - Мониторинг и отладка распределенных систем. Умеет: Проектировать масштабируемые data-приложения - Использовать бессерверные архитектуры для ETL - Оптимизировать производительность Python-кода - Интегрировать различные cloud-сервисы. Владеет – навыками разработки на PySpark, Dask, Ray - Созданием и деплом cloud functions - Контейнеризацией приложений (Docker) - Настройкой автоматического масштабирования - Оптимизация costs в cloud-среде.

### Содержание и структура дисциплины:

Распределение видов учебной работы и их трудоемкости по разделам дисциплины.  
Разделы дисциплины, изучаемые в 7 семестре (очная форма)

№	Наименование разделов (тем)	Количество часов			
		Всего	Аудиторная работа		Внеаудиторная работа
			Л	ЛР	СРС
1	2	3	4	6	7
1.	Введение в Big Data. Проблематика и базовые концепции.	6	2	2	2
2.	<b>Распределенная файловая система HDFS и объектное хранение. Hadoop и экосистема</b>	6	2	2	2
3.	<b>Модели вычислений. MapReduce и его эволюция.</b>	6	2	2	2
4.	<b>Введение в Apache Spark. Архитектура и RDD.</b>	6	2	2	2
5.	<b>Spark SQL и DataFrames.</b>	6	2	2	2
6.	<b>Оптимизация в Spark.</b>	6	2	2	2
7.	<b>Работа с Spark в облаке и кластерном режиме.</b>	12	4	4	4
8.	<b>Потоковая обработка данных со Structured Streaming.</b>	8	2	2	4
9.	<b>ETL-пайплайны на Spark. Best Practices..</b>	6	2	2	2
10.	<b>OLAP vs OLTP. Колоночные базы данных.</b>	6	2	2	2
11.	<b>Глубокое погружение в ClickHouse. Движки таблиц и партиционирование.</b>	6	2	2	2
12.	<b>Оптимизация запросов в ClickHouse.</b>	6	2	2	2
13.	<b>Интеграция Spark и ClickHouse.</b>	6	2	2	2
14.	<b>Инструменты оркестрации данных. Apache Airflow.</b>	6	2	2	2
15	<b>Архитектура Big Data-решений на практике: кейсы индустриальных партнеров</b>	12	4	4	4
<b>ИТОГО по разделам дисциплины</b>		<b>104</b>	<b>34</b>	<b>34</b>	<b>36</b>
Контроль самостоятельной работы (КСР)		<b>4</b>			
Промежуточная аттестация (ИКР)		<b>0,3</b>			
Подготовка к текущему контролю		<b>35,7</b>			
<b>Общая трудоемкость по дисциплине</b>		<b>144</b>			

*Примечание: Л – лекции, КСР – контрольные и самостоятельные работы, ЛР – лабораторные занятия, СРС – самостоятельная работа студента*

#### **Курсовые проекты или работы.**

Не предусмотрены учебным планом

**Вид аттестации:** ЛР, проект по кейсам индустриальных партнеров, экзамен.

Автор Приходько Т.А. – кандидат технических наук, доцент кафедры вычислительных технологий