



Министерство науки  
и высшего образования РФ  
ФГБОУ ВО «НИУ «МЭИ»  
Институт дистанционного  
и дополнительного образования



**АННОТАЦИИ РАБОЧИХ ПРОГРАММ ДИСЦИПЛИН (МОДУЛЕЙ)  
ДОПОЛНИТЕЛЬНОЙ ОБРАЗОВАТЕЛЬНОЙ ПРОГРАММЫ**

*повышения квалификации  
«Основы анализа текстовых данных»,*

Раздел(предмет) *Основы анализа данных*

| Наименование дисциплин (модулей)            | Содержание дисциплин (модулей)   | Форма ТК            | Количество часов |
|---|--|---------------------|------------------|
| <i>Постановка задачи машинного обучения</i> | Типы задач, решаемых методами машинного обучения (МО). Формальная постановка задачи машинного обучения. Задачи анализа тестовой информации, решаемые методами МО. Критерии качества: аккуратность, полнота, точность, F-мера, площадь под кривой ошибок, матрица неточностей. Способы формирования выборок. Обучающая, экзаменационная и тестовая выборки. Обучение моделей. Явление переобучения модели. Несбалансированные классы. Методы борьбы с несбалансированностью: oversampling, undersampling. Специальные стратегии сэмплинга в условиях несбалансированных классов | <i>Тестирование</i> | 8                |

Раздел(предмет) *Анализ текстовых данных (Text Mining)*

| Наименование дисциплин (модулей)                       | Содержание дисциплин (модулей)  | Форма ТК            | Количество часов |
|--|---|---------------------|------------------|
| <i>Особенности и задачи обработки текстовых данных</i> | Text Mining и особенности задач, связанных с анализом текстов. Онтологии и тезаурусы. Статистический подход к анализу текстовой информации. Проблема снижения размерности в задачах Text Mining. Предварительная обработка данных: стемминг, лемматизация. Слова, не несущие информации. Выявление информативных признаков. Взвешивание как способ выявления информативных терминов. Статистический подход к выявлению информативных терминов. Теоретико-информационный подход к выявлению информативных терминов. Модель представления текстовых данных в математическом виде. Модель «Мешок слов» (Bag of words). Частично- и полностью структурированные модели. | <i>Тестирование</i> | 8                |

Раздел(предмет) ***Задача классификации текстовых документов***

| Наименование дисциплин (модулей)   | Содержание дисциплин (модулей)   | Форма ТК            | Количество часов |
|------------------------------------|--|---------------------|------------------|
| <i>Методы классификации данных</i> | Систематизация и обзор методов классификации. Метод ближайших соседей. Метод деревьев решений. Метод опорных векторов. Метод логистической регрессии. Наивный байесовский метод. Профильные методы классификации. Ансамблевые методы классификации. Оценка | <i>Тестирование</i> | 8                |

| Наименование дисциплин (модулей) | Содержание дисциплин (модулей)   | Форма ТК | Количество часов |
|----------------------------------|--|----------|------------------|
|                                  | разнородности методов классификации. Бустинг, бэггинг. Метод случайного леса деревьев решений. |          |                  |

Раздел(предмет) *Другие задачи, решаемые в рамках Text Mining*

| Наименование дисциплин (модулей)                                       | Содержание дисциплин (модулей)  | Форма ТК            | Количество часов |
|--|---|---------------------|------------------|
| <i>Кластеризация данных, выявление дубликатов текстовых документов</i> | Постановка задачи кластеризации, ее особенности. Иерархическая кластеризация. EM-алгоритм кластеризации. Семейство алгоритмов k-means, другие алгоритмы кластеризации. Виды дубликатов: полные дубликаты, явные дубликаты, нечеткие дубликаты. Коэффициент ассоциативности Жаккара. Семейство методов шинглов. Методы выявления дубликатов Winnowing, SpotSigs, I-Match, коэффициент Джаро-Винклера | <i>Тестирование</i> | <i>10</i>        |

Руководитель  
ОДПО, ЦПП УВО  
(должность)

|   |  |                                 |
|---|--|---------------------------------|
|  | Подписано электронной подписью ФГБОУ ВО «НИУ «МЭИ» |                                 |
|   | Сведения о владельце ЦЭП МЭИ                       |                                 |
|   | Владелец   | Максимова А.А.                  |
|   | Идентификатор                                      | R6a033f13-VoroZhtsovaAA-daecd83 |

А.А.  
Максимова  
(расшифровка подписи)

Начальник ОДПО  
(должность)

|   |  |                              |
|---|--|------------------------------|
|  | Подписано электронной подписью ФГБОУ ВО «НИУ «МЭИ» |                              |
|   | Сведения о владельце ЦЭП МЭИ                       |                              |
|   | Владелец   | Крохин А.Г.                  |
|   | Идентификатор                                      | R6d4610d5-KrokhinAG-aa301f84 |

А.Г. Крохин  
(расшифровка подписи)