

БНТУ

Научная библиотека



\* 8 0 1 2 5 5 0 8 5 \*

# Создание приложений машинного обучения

от идеи к продукту

Эммануэль Амейзен



Санкт-Петербург · Москва · Минск

2023

---

# Оглавление

<b>Предисловие .....</b>	<b>10</b>
Зачем нужны приложения на базе машинного обучения.....	10
Используйте МО для создания практических приложений .....	10
Дополнительные ресурсы .....	11
МО на практике .....	12
Что вы найдете в книге.....	13
Необходимая подготовка .....	14
Наш учебный пример: написание текстов с использованием МО .....	14
Процесс МО как он есть .....	15
Условные обозначения.....	16
Использование исходного кода примеров.....	17
Благодарности .....	18
От издательства.....	18

## ЧАСТЬ I

### Выбор правильного подхода к машинному обучению

<b>Глава 1. От цели продукта к разработке модели МО .....</b>	<b>21</b>
Оценка осуществимости модели .....	22
Модели .....	24
Данные.....	31
Формулировка задачи по созданию редактора на основе МО .....	35
Делаем все с помощью МО: подход «от начала до конца» .....	35
Простейший подход: «сам себе алгоритм» .....	37
Золотая середина: обучение на полученном опыте .....	38
Моника Рогати: как выбирать и приоритизировать МО-проекты.....	40
Заключение .....	43
<b>Глава 2. Составление плана.....</b>	<b>44</b>
Оценка успешности.....	44
Производительность с точки зрения бизнеса .....	45

Производительность модели .....	46
Актуальность данных и сдвиг распределения .....	50
Скорость .....	52
Оценка масштаба и возможных проблем .....	53
Накапливайте знания в предметной области .....	53
Используйте опыт предшественников .....	55
Составление плана разработки МО-редактора .....	59
Начальный план разработки редактора.....	59
Всегда начинайте с простой модели.....	60
Как обеспечить устойчивый прогресс? Начинайте с простейшего решения .....	61
Начинайте с простого пайплайна .....	61
Пайплайн для МО-редактора.....	63
Заключение .....	65

## **ЧАСТЬ II**

### **Создание рабочего пайплайна**

<b>Глава 3. Создание первого сквозного пайплайна .....</b>	<b>68</b>
Простейший «каркас» приложения .....	68
Прототип МО-редактора.....	70
Парсинг и очистка данных .....	70
Токенизация текста .....	71
Генерирование признаков .....	72
Оцените рабочий процесс .....	74
Пользовательский опыт .....	74
Результаты моделирования.....	75
Оценка прототипа МО-редактора .....	76
Модель.....	77
Пользовательский опыт .....	78
Заключение .....	78
<b>Глава 4. Получение исходного датасета .....</b>	<b>80</b>
Итеративная доработка датасета.....	80
Проведите анализ данных.....	81
Изучение первого датасета .....	82
Начните с малого .....	82
Инсайт vs продукты .....	83
Критерии оценки качества данных .....	84

---

Разметка для выявления трендов в данных.....	91
Сводная статистика .....	91
Исследуйте и размечайте эффективно.....	93
Станьте алгоритмом .....	110
Тренды в данных .....	112
Используйте данные для принятия решений о признаках и моделях.....	113
Создание признаков на основе закономерностей.....	113
Признаки МО-редактора .....	117
Роберт Манро: как находить, размечать и использовать данные .....	118
Заключение .....	120

### **ЧАСТЬ III** **Итеративная доработка моделей**

<b>Глава 5. Обучение и оценка модели .....</b>	<b>123</b>
Самая простая адекватная модель.....	123
Простая модель.....	124
От закономерностей к моделям.....	126
Разделение датасета .....	128
Разделение данных МО-редактора .....	135
Оценка производительности.....	137
Оценка модели: не ограничивайтесь точностью .....	140
Сравнение предсказаний с данными.....	140
Матрица ошибок.....	141
ROC-кривая.....	142
Калибровочная кривая.....	145
Снижение размерности для анализа ошибок.....	145
Метод первых k элементов .....	147
Другие модели .....	151
Оценка важности признаков .....	152
Непосредственная оценка классификатора.....	153
Интерпретаторы «черного ящика».....	154
Заключение .....	156
<b>Глава 6. Отладка МО-приложения .....</b>	<b>157</b>
Передовые практики разработки ПО .....	157
Передовые практики для МО-приложений .....	159
Отладка потока данных: визуализация и тестирование .....	160

Начните с одного образца.....	160
Тестирование МО-кода .....	167
Отладка обучения: заставьте модель учиться .....	172
Сложность задачи .....	173
Проблемы оптимизации.....	175
Отладка обобщения: модель должна быть полезной .....	178
Утечка данных .....	178
Переобучение .....	179
Изучение решаемой задачи .....	183
Заключение .....	183
<b>Глава 7. Использование классификаторов для выдачи рекомендаций ....</b>	<b>184</b>
Вывод рекомендаций.....	185
Чего мы можем добиться без модели? .....	185
Извлечение глобальной важности признаков .....	187
Использование балла модели .....	188
Извлечение локальной важности признаков .....	188
Сравнение моделей .....	191
Версия 1: простейший отчет.....	191
Версия 2: более мощная, менее понятная .....	192
Версия 3: понятные рекомендации.....	194
Создание рекомендаций по редактированию текста .....	195
Заключение .....	199

## **ЧАСТЬ IV** **Развертывание и мониторинг**

<b>Глава 8. Что еще учесть при развертывании модели .....</b>	<b>203</b>
Забота о данных .....	204
Право собственности на данные.....	204
Смещение данных .....	205
Систематическое смещение .....	207
Забота о модели .....	208
Циклы обратной связи.....	208
Инклюзивная производительность модели.....	210
О контексте.....	211
Мошенники .....	212
Риск злоупотребления и двойного назначения.....	213

---

Крис Харланд: опыт поставки продуктов .....	214
Заключение .....	217
<b>Глава 9. Выбор варианта развертывания .....</b>	<b>218</b>
Развертывание на сервере .....	218
Потоковое приложение или API .....	219
Пакетные предсказания .....	221
Развертывание на стороне клиента.....	223
Развертывание на устройстве .....	225
Развертывание в браузере.....	227
Федеративное обучение: комбинированный подход.....	227
Заключение .....	229
<b>Глава 10. Создание защитных механизмов для моделей .....</b>	<b>231</b>
Проектирование с учетом возможных сбоев .....	231
Проверка входных данных и результатов.....	232
Резервные варианты на случай сбоя модели.....	236
Проектирование для обеспечения высокой производительности .....	240
Масштабирование при возрастании числа пользователей .....	241
Управление жизненным циклом модели и данных.....	244
Обработка данных и DAG .....	247
Запрос обратной связи.....	249
Крис Муди: специалисты по данным отвечают за весь пайплайн моделирования .....	252
Заключение .....	254
<b>Глава 11. Мониторинг и обновление моделей.....</b>	<b>255</b>
Мониторинг спасает жизни .....	255
Мониторинг для определения частоты обновлений .....	256
Мониторьте, чтобы выявить злоупотребления .....	256
Что мониторить .....	257
Метрики производительности .....	258
Бизнес-метрики .....	260
CI/CD для МО .....	261
А/В-тестирование и эксперименты .....	263
Другие подходы .....	266
Заключение .....	268
<b>Об авторе .....</b>	<b>270</b>
<b>Иллюстрация на обложке .....</b>	<b>271</b>