

УДК 004.415.2

**АРХИТЕКТУРНЫЕ ПАТТЕРНЫ ПОСТРОЕНИЯ
ВЫСОКОНАГРУЖЕННЫХ СИСТЕМ ПРЕДИКТИВНОГО
МОНИТОРИНГА OPEN-SOURCE ЗАВИСИМОСТЕЙ НА БАЗЕ
РЕКУРРЕНТНЫХ НЕЙРОСЕТЕЙ**

Брайловский А.В.

Магистрант,

МИРЭА – Российский Технологический Университет,

Москва, Россия

Иванова А.П.

к.ф.-м.н., доцент,

МИРЭА – Российский Технологический Университет,

Москва, Россия

Аннотация:

В статье рассматривается комплексная проблема проектирования масштабируемых систем для непрерывного предиктивного мониторинга жизненного цикла программного обеспечения с открытым исходным кодом. Обосновывается необходимость перехода от эпизодического статического анализа к проактивному потоковому мониторингу в связи с экспоненциальным ростом деревьев транзитивных зависимостей. Автором предлагается многоуровневая архитектура программного комплекса, базирующаяся на гибридном асинхронном конвейере данных. Описан алгоритм локального формирования еженедельных JSON-слепков репозитория с использованием утилиты `sls`, обогащаемых историческими социотехническими метриками через программные интерфейсы GraphQL и REST. Представлен механизм интеллектуальной маршрутизации и асинхронного вычисления прогнозов на базе рекуррентных нейронных сетей (LSTM), позволяющий оптимизировать использование лимитов API, изолировать аппаратные и социальные шумы

Дневник науки | www.dnevniknauki.ru | СМИ Эл № ФС 77-68405 ISSN 2541-8327

(CI/CD боты, Vendor Drops) и обеспечивать мгновенный отклик системы на запросы пользователей посредством жесткого разделения слоев сбора данных и тензорного инференса.

Ключевые слова: программная инженерия, open-source, транзитивные зависимости, предиктивная аналитика, фоновый инференс, рекуррентные нейросети, LSTM, конвейер данных, GitHub API.

***ARCHITECTURAL PATTERNS FOR BUILDING HIGH-LOAD SYSTEMS
FOR PREDICTIVE MONITORING OF OPEN-SOURCE DEPENDENCIES
BASED ON RECURRENT NEURAL NETWORKS***

Braylovskiy A.V.

Master's student,

MIREA - Russian Technological University

Moscow, Russia

Ivanova A.P.

PhD, Associate Professor,

MIREA – Russian Technological University,

Moscow, Russia

Abstract:

The article discusses the complex problem of designing scalable systems for continuous predictive monitoring of the open-source software lifecycle. The necessity of transitioning from episodic static analysis to proactive streaming monitoring due to the exponential growth of transitive dependency trees is substantiated. The author proposes a multi-level software architecture based on a hybrid asynchronous data pipeline. An algorithm for the local generation of weekly JSON snapshots of the repository using the cloc utility, enriched with historical sociotechnical metrics via GraphQL and REST programming interfaces, is described. A mechanism for

Дневник науки | www.dnevniknauki.ru | СМИ Эл № ФС 77-68405 ISSN 2541-8327

intelligent routing and asynchronous calculation of predictions based on recurrent neural networks (LSTM) is presented, allowing for the optimization of API limit usage, isolation of hardware and social noise (CI/CD bots, Vendor Drops), and ensuring immediate system response to user requests through strict separation of data collection and tensor inference layers.

Keywords: software engineering, open-source, transitive dependencies, predictive analytics, background inference, recurrent neural networks, LSTM, data pipeline, GitHub API.

Введение

Современная парадигма разработки программного обеспечения неразрывно связана с принципами переиспользования кода и интеграции сторонних библиотек. Экосистемы пакетов с открытым исходным кодом (open-source) являются фундаментальной основой для подавляющего большинства информационных систем. Однако такой уровень интеграции порождает критический вектор уязвимости, связанный с непредсказуемостью жизненного цикла внешних компонентов.

Проблема многократно усугубляется механизмом транзитивных зависимостей. При добавлении в проект одной базовой библиотеки разработчик неявно интегрирует в кодовую базу все зависимости этой библиотеки, а также зависимости ее зависимостей. Если активная поддержка даже одного, глубоко вложенного узла этого графа прекращается (проект переходит в стадию стагнации или забрасывается автором), это ставит под угрозу возможность обновления, масштабирования и обеспечения безопасности всего конечного продукта. Критичность ситуации обостряется феноменом низкого «Bus Factor» (фактора автобуса) — для многих фундаментальных open-source библиотек контроль над репозиторием сосредоточен в руках одного-двух независимых мейнтейнеров.

Своевременное выявление риска прекращения поддержки пакета до момента фактического отказа требует внедрения автоматизированных систем предиктивного мониторинга на базе архитектур глубокого обучения [1, 2]. Данные системы должны непрерывно оценивать «здоровье» репозитория на основе комплексного анализа нестационарных временных рядов технической эволюции кода и социальной вовлеченности сообщества. Разработка масштабируемой архитектуры подобной системы, способной функционировать в условиях жестких ограничений программных интерфейсов (API), представляет собой сложную инженерно-научную задачу.

1. Вычислительная сложность сбора данных и проблема «холодного старта»

При проектировании ядра системы мониторинга фундаментальным вызовом является колоссальный объем данных, генерируемых в процессе эволюции open-source проектов [5]. Стандартный подход к анализу истории репозитория предполагает процесс «холодного старта» — полного первичного копирования кодовой базы на локальный сервер (операция git clone).

Для проектов с многолетней историей разработка алгоритма, который бы «на лету» скачивал и пересчитывал всю статистику при каждом синхронном запросе пользователя, нежизнеспособна. Помимо затрат на дисковые операции, сбор дополнительных социальных метрик (issues, forks, активности участников) жестко лимитирован провайдерами исходного кода (GitHub, GitLab). Следовательно, архитектура предиктивного сервиса должна базироваться на принципах асинхронного, инкрементального накопления данных, где сложные вычисления выполняются в фоне, а пользовательский интерфейс оперирует лишь заранее подготовленными ответами [7].

2. Архитектура гибридного конвейера обработки данных (Data Pipeline)

Для решения описанных проблем разработана гибридная архитектура конвейера данных, разделяющая процессы локального анализа структуры кода

Дневник науки | www.dnevniknauki.ru | СМИ Эл № ФС 77-68405 ISSN 2541-8327

и удаленного обогащения метриками через API. Весь процесс выполняется служебными скриптами-планировщиками в фоновом режиме, независимо от действий пользователя.

2.1 Слой локального структурного анализа

Базовым источником данных для предиктивной модели является физическая структура проекта. Алгоритм выполняет первичное клонирование репозитория. Далее запускается итеративный процесс по исторической шкале времени: система разбивает историю проекта на недельные интервалы. Для каждой недели вычисляется хэш (SHA) последнего коммита. Чтобы получить состояние проекта на этот момент времени, скрипт выполняет команду `git archive | tar`, извлекая «чистый» слепок файлов во временную директорию.

К этой директории применяется специализированная утилита статического анализа `cloc` (Count Lines of Code), которая собирает метрики объема кодовой базы, количества комментариев и пустых строк. Параллельно высчитывается глубина директорий и распределение файлов. Итогом этого этапа является формирование отдельного JSON-файла для каждой обрабатываемой недели (например, `week_2024-01-01.json`). Этот подход решает проблему инкрементального обновления: при наступлении новой календарной недели или при фоновом обновлении проекта скрипт лишь высчитывает дельту и создает новый JSON-файл для новой недели, не пересчитывая исторические архивы [9].

2.2 Слой обогащения через API (REST и GraphQL)

Сформированные структурные JSON-отчеты требуют обогащения историческим контекстом. GitHub позволяет извлекать метрики ретроспективно, однако агрессивный парсинг быстро исчерпывает часовые лимиты. Предложенная архитектура решает эту задачу точечным гибридным обогащением. Используя язык запросов GraphQL и фильтры поиска (например, `is:issue state:closed closed:YYYY-MM-DD..YYYY-MM-DD`), алгоритм за один вызов запрашивает точное количество задач и форков за конкретную неделю

Дневник науки | www.dnevnika.ru | СМИ Эл № ФС 77-68405 ISSN 2541-8327

[3]. REST API используется для получения детальной статистики структуры по SHA коммита. Эти данные инкапсулируются в еженедельный JSON-файл, делая его самодостаточным слепком.

Схема работы гибридного конвейера данных и связанного с ним аналитического ядра представлена на рис. 1.

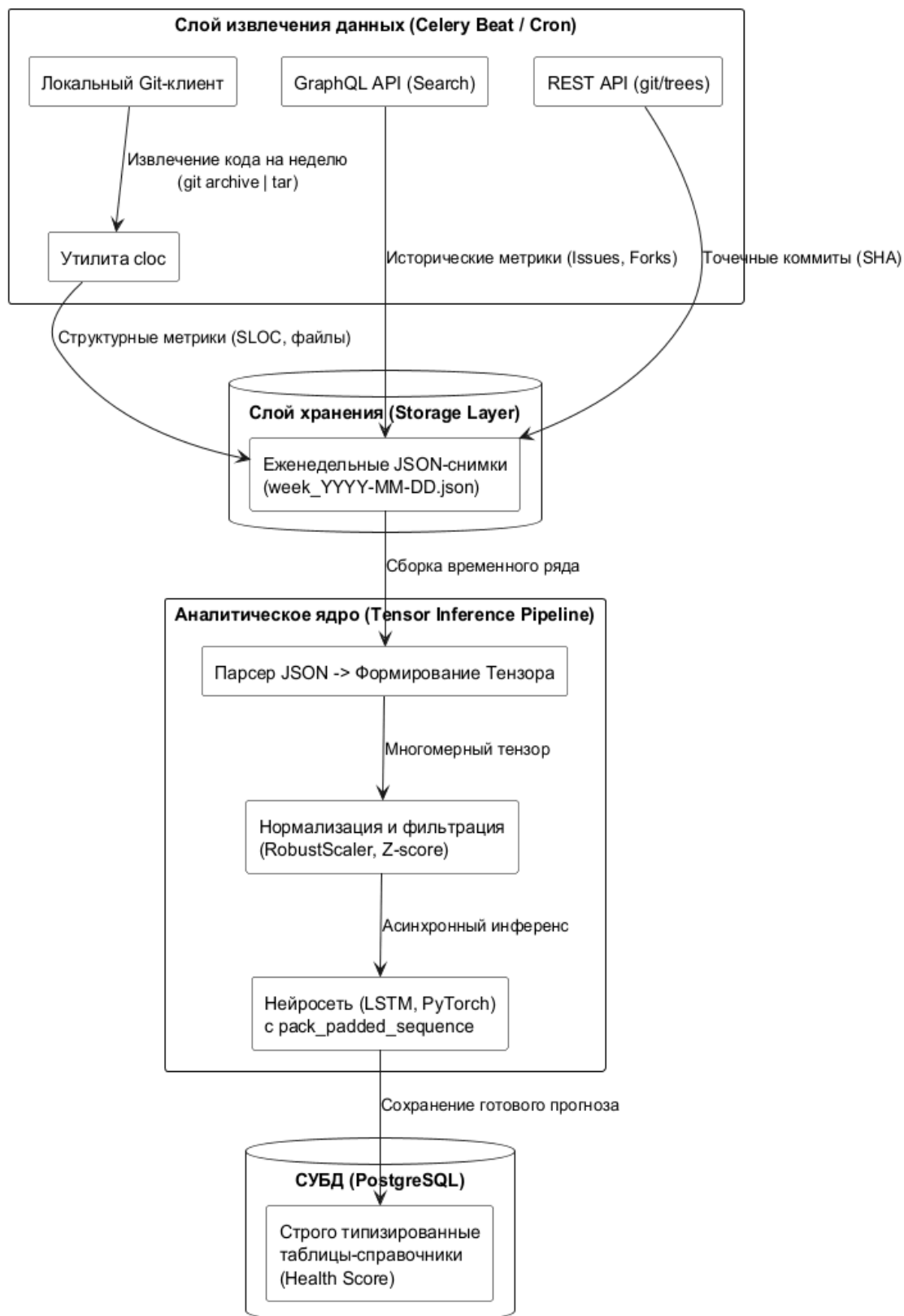


Рисунок 1 – Фоновый гибридный конвейер агрегации метрик. Авторская разработка

3. Асинхронный тензорный инференс и маршрутизация пользовательских запросов

Ключевым архитектурным решением, обеспечивающим работу системы в реальном времени, является полное отвязывание процесса инференса (предсказания нейросети) от синхронного пользовательского HTTP-запроса [6]. Модель глубокого машинного обучения не запускается в момент обращения пользователя к веб-интерфейсу.

Процесс подготовки прогноза выглядит следующим образом:

Как только скрипт сбора данных формирует новые JSON-файлы, специализированный парсер считывает всю накопленную историю проекта.

Вместо агрегации данных в плоские структуры с потерей временного контекста, алгоритм формирует многомерный тензор признаков $X_i = \{x_{i,1}, x_{i,2}, \dots, x_{i,T}\}$, где $x_{i,t}$ – вектор социотехнического состояния i -го репозитория в неделю t , а T – общая длина собранной истории. Перед подачей в вычислительный граф конвейер применяет эвристические фильтры пространственно-временного шума: выявляются аномальные всплески активности CI/CD ботов на основе стандартного отклонения времени коммитов, а также детектируются массовые вбросы легаси-кода (Vendor Drops) посредством алгоритмов динамического Z-score отклонения. Для нивелирования проблемы переменной длины жизненных циклов T применяется механизм аппаратного прерывания вычислений с использованием функции `pack_padded_sequence` фреймворка PyTorch. Данный подход позволяет нейросети игнорировать фиктивные нулевые элементы (Padding) на уровне архитектуры CUDA.

В фоновом режиме запускается рекуррентная нейронная сеть архитектуры Long Short-Term Memory (LSTM). Нейросеть последовательно обрабатывает временной ряд, аккумулируя историческую память в скрытом векторе состояний h_t [4].

Вероятность наступления целевого события (прекращения поддержки i -го проекта) вычисляется через сигмоидальную активацию полносвязного слоя над финальным скрытым состоянием:

$$P(y_i = 1|X_i) = \sigma(W \cdot h_T + b),$$

где W – матрица весов полносвязного слоя; h_T – вектор скрытого состояния на последнем временном шаге; b – смещение (bias).

Итоговый результат (Health Score) сохраняется напрямую в структурированные реляционные таблицы базы данных PostgreSQL. Отказ от хранения расчетных метрик в слабоструктурированных форматах (например, JSONB) в пользу строго типизированных полей обеспечивает максимальную скорость чтения готовых прогнозов клиентской частью и гарантирует предсказуемую производительность СУБД при масштабировании.

Алгоритм обработки пользовательского запроса предельно упрощен и сводится к мгновенному чтению готовых результатов. Важной особенностью алгоритма является паттерн идемпотентности, защищающий пул воркеров от проблемы «Громящего стада» (Thundering Herd), когда сотни пользователей одновременно запрашивают отсутствующий проект. Логика маршрутизации запросов проиллюстрирована на рис. 2.

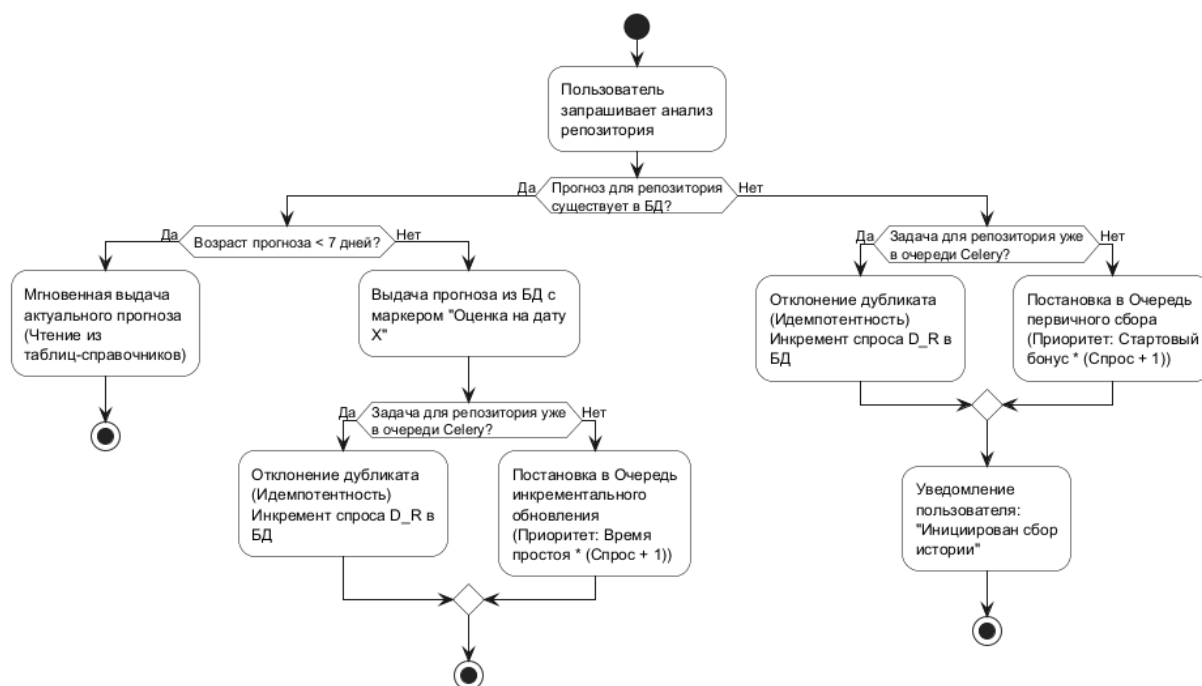


Рисунок 2 – Алгоритм обработки пользовательского запроса с защитой от DDoS-аномалий. Авторская разработка

Позиция репозитория в очереди обновления (представленной на рисунке 2) вычисляется динамически на этапе постановки задачи на основе степени устаревания данных и множителя спроса. Математически приоритет $Priority_R$ выражается как:

$$Priority_R = (t_{current} - t_{last_update}) \cdot (D_R + 1)$$

где $(t_{current} - t_{last_update})$ – время простоя в днях, а D_R – количество уникальных запросов пользователей к данному репозиторию за отчетный период [8]. Введение константы смещения гарантирует, что даже глубокие транзитивные зависимости ($D_R = 0$) со временем достигнут начала очереди исключительно за счет своего возраста. В свою очередь, востребованные пользовательские проекты получаюткратно больший приоритет при формировании новой фоновой задачи. В случае если задача уже ожидает в пуле Celery, алгоритм отклоняет дубликат, защищая брокер сообщений. При этом счетчик D_R инкрементируется в системной базе данных, что гарантирует Дневник науки | www.dnevnikaui.ru | СМИ Эл № ФС 77-68405 ISSN 2541-8327

автоматическое повышение приоритета проекта при последующих циклах обновления.

4. Оценка эффективности разработанной архитектуры

Для верификации предложенного паттерна была проведена сравнительная оценка временных затрат. В качестве эталона анализировался репозиторий среднего размера (500 файлов, 150 000 строк кода), типичный для базовых инфраструктурных open-source решений. Таблица 1 демонстрирует ресурсоемкость процесса сканирования и извлечения метрик в рамках одного недельного интервала.

Таблица 1 – Вычислительные затраты на обработку одного недельного интервала (для репозитория среднего размера). Авторская разработка

Этап фонового конвейера	Время первичного извлечения (1 итерация)	Время инкрементального обновления (новая неделя)
Извлечение снимка кода (git archive + tar)	12.0 – 15.0 секунд	12.0 – 15.0 секунд (выполняется 1 раз)
Анализ утилитой Cloudf	3.0 – 5.0 секунд	3.0 – 5.0 секунд
Обогащение API (GraphQL Search + REST)	4.0 – 6.0 секунд	1.0 – 2.0 секунды
Формирование JSON и тензорная агрегация	< 0.5 секунды	< 0.5 секунды
Итоговое время отклика для пользователя	Уведомление об очереди	~ 0.01 секунды (Чтение из БД)

Как видно из результатов, представленных в таблице 1, обработка даже одной недели истории для реального проекта занимает более 20 секунд. Таким образом, полноценный процесс «холодного старта» для репозитория со средней историей существования в 5 лет (около 260 недель) займет порядка 260×20 секунд \approx 86 минут синхронного серверного времени.

Подобная длительность математически доказывает невозможность использования классических подходов с вычислениями «на лету». Вынесение этих ресурсоемких операций (cloc, git archive, обращения к API) и процесса тензорного инференса нейросети в жесткий фоновый конвейер является не опциональным улучшением, а архитектурной необходимостью, позволяющей свести время отклика системы на пользовательский запрос к миллисекундам, требуемым лишь на обращение к СУБД.

Заключение

В рамках данного исследования спроектирована масштабируемая архитектура программного комплекса, позволяющая трансформировать тяжеловесные модели глубокого обучения (Deep Learning) в производительный инструмент предиктивного мониторинга.

Доказано, что жесткое разделение системы на независимые слои (фоновое клонирование и анализ cloc с генерацией недельных JSON-отчетов, обогащение через GraphQL/REST API и асинхронный тензорный инференс на базе LSTM с записью в строго типизированную реляционную БД) кардинально снижает вычислительную нагрузку в моменты пикового пользовательского спроса.

Внедренная система интеллектуальной приоритизации очередей и паттернов идемпотентности гарантирует защиту от аномальных всплесков трафика и обеспечивает рациональное распределение серверных ресурсов. Данный архитектурный паттерн позволяет осуществлять непрерывный мониторинг тысяч транзитивных зависимостей без блокировок пользовательского интерфейса, обеспечивая проактивное управление рисками в корпоративной разработке программного обеспечения.

Библиографический список:

1. Брайловский, А.В. Анализ факторов стабильности и прогнозирование жизненного цикла open-source проектов на основе данных платформы GitHub / Дневник науки | www.dnevnikaui.ru | СМН Эл № ФС 77-68405 ISSN 2541-8327

- А.В. Брайловский // «Парадигма»: научно-практический электронный журнал. – 2025. – № 6-1. – С. 38-46. EDN WRSEVM.
2. Брайловский, А.В. Методология сбора данных для анализа активности проектов с открытым исходным кодом на платформе GitHub / А.В. Брайловский // «Парадигма»: научно-практический электронный журнал. – 2025. – № 6-1. – С. 47-57. EDN LWCXPI.
3. Галигузова, Е.В. Язык запросов GraphQL как замена REST API / Е.В. Галигузова, Ю.Е. Илларионова // Символ науки: международный научный журнал. – 2023. – № 1-2. – С. 9-11.
4. Коротеев, М.В. Машинное обучение: учебник для вузов / М.В. Коротеев. – Москва: Издательство Юрайт, 2023. – 311 с.
5. Радченко, Г.И. Распределенные вычислительные системы: учебное пособие / Г.И. Радченко. – Москва: Издательство Юрайт, 2022. – 174 с.
6. Вьюгин, В.В. Математические основы теории машинного обучения и прогнозирования / В.В. Вьюгин. – Москва: МЦНМО, 2023. – 304 с.
7. Гагарина, Л.Г. Разработка и эксплуатация автоматизированных информационных систем: учебник / Л.Г. Гагарина. – Москва: ИНФРА-М, 2021. – 384 с.
8. Северинцев, Н.А. Проектирование баз данных: учебное пособие / Н.А. Северинцев. – Санкт-Петербург: Лань, 2022. – 208 с.
9. Чернышев, А.А. Инженерия программного обеспечения: учебник / А.А. Чернышев. – Москва: КноРус, 2021. – 232 с.