

Курсовая работа на тему:
"Прогнозирование эффективности химических соединений на основе молекулярных дескрипторов"

Решетило Никита
1 курс, «Машинное обучение»

Содержание

| | | |
|----|--|----|
| 1. | Описание данных..... | 3 |
| a. | Актуальность исследования..... | 3 |
| b. | Ключевые показатели эффективности..... | 4 |
| c. | Цели исследования..... | 4 |
| d. | Методологический подход..... | 4 |
| e. | Классификация молекулярных дескрипторов..... | 7 |
| 2. | EDA..... | 10 |
| a. | Первичный анализ..... | 10 |
| b. | Анализ молекулярных дескрипторов..... | 11 |
| c. | Общие результаты обработки данных..... | 12 |
| d. | Анализ аномальных признаков..... | 12 |
| e. | Химическая интерпретация аномалий..... | 15 |
| f. | Анализ выбросов молекулярных свойств..... | 17 |
| g. | Анализ корреляции молекулярных дескрипторов с биологической активностью..... | 19 |
| h. | Итог по анализу молекулярных данных..... | 28 |
| 3. | Исследовательский анализ..... | 29 |
| a. | Анализ моделей прогнозирования индекса селективности (IC50)..... | 29 |
| b. | Анализ моделей прогнозирования цитотоксической активности соединений (CC50)..... | 31 |
| c. | Анализ моделей прогнозирования для регрессии индекса селективности (SI)..... | 36 |
| d. | Классификация: превышает ли значение SI медианное значение выборки..... | 41 |
| e. | Классификация: превышает ли значение IC50 медианное значение выборки..... | 44 |
| f. | Классификация: превышает ли значение CC50 медианное значение выборки | 49 |
| g. | Классификация: превышает ли значение SI значение 8..... | 54 |
| 4. | Финальный вывод по проделанной работе..... | 56 |
| 5. | Источники данных..... | 58 |

I. Описание данных

а. Актуальность исследования

Оптимизация процесса разработки лекарственных препаратов требует раннего выявления перспективных соединений с высокой биологической активностью и низкой токсичностью. Современные методы хемоинформатики и машинного обучения позволяют прогнозировать фармакологические характеристики молекул на основе их структурных особенностей, минимизируя затраты на дорогостоящие лабораторные исследования.

Применение компьютерного моделирования в доклинических исследованиях значительно ускоряет процесс скрининга потенциальных лекарственных соединений, что особенно важно в условиях растущих требований к безопасности и эффективности новых препаратов. Методы количественного анализа структура-активность (QSAR) и молекулярного докинга уже доказали свою эффективность в предсказании биологической активности соединений [1].

Кроме того, использование молекулярных дескрипторов и алгоритмов машинного обучения позволяет выявлять скрытые закономерности в данных, что способствует более точному прогнозированию ключевых параметров, таких как IC_{50} , CC_{50} и индекс селективности (SI). Это особенно актуально в контексте персонализированной медицины, где требуется быстрый и точный анализ большого количества химических соединений.

Таким образом, интеграция методов хемоинформатики и искусственного интеллекта в процесс разработки лекарств открывает новые возможности для сокращения времени и затрат на исследования, что подтверждается современными научными работами.

б. Ключевые показатели эффективности

В исследовании фокусируется на прогнозировании трех критически важных параметров:

1. IC_{50} – концентрация вещества, необходимая для подавления целевого биологического процесса на 50%, показатель эффективности

2. CC_{50} – концентрация, вызывающая гибель 50% клеток, индикатор цитотоксичности
3. Индекс селективности (SI) – отношение CC_{50}/IC_{50} , отражающее терапевтическое окно (чем выше значение, тем лучше профиль безопасности)

с. **Цели исследования**

Разработка и сравнительный анализ:

- Регрессионных моделей для точного предсказания численных значений IC_{50} , CC_{50} и SI
- Классификационных моделей для бинарной оценки показателей относительно:
 - Медианных значений выборки
 - Порогового значения $SI \geq 8$ (клинически значимый критерий безопасности)

д. **Методологический подход**

В рамках исследования будут рассмотрены:

1. Тестирование различных алгоритмов машинного обучения
2. Сравнительную оценку качества моделей для каждой цели с обоснованием выбора оптимальных решений
3. Анализ предсказательной способности моделей для принятия решений в доклинических исследованиях

Рассмотрим подробнее предоставленные к реализации модели:

1. Регрессия для IC_{50}

а. **Цель**

Построить модель, предсказывающую близкую к максимальной ингибирующую концентрацию (IC_{50}) на основе молекулярных дескрипторов.

b. Практическая значимость

- i. IC_{50} показывает, насколько эффективно соединение подавляет целевую биологическую мишень.
 - ii. Низкие значения IC_{50} (<1 мкМ) указывают на высокую активность.
 - iii. Модель поможет отбирать наиболее перспективные соединения для дальнейших исследований.
-

2. Регрессия для CC_{50}

a. Цель

Разработать модель, прогнозирующую близкую к максимальной цитотоксической концентрации (CC_{50}) по структурным признакам.

b. Практическая значимость

- i. CC_{50} отражает токсичность соединения для здоровых клеток.
 - ii. Высокие значения (> 50 мкМ) желательны для безопасных лекарств.
 - iii. Модель позволит отсеивать слишком токсичные соединения на ранних этапах.
-

3. Регрессия для SI (Selectivity Index)

a. Цель

Создать модель для предсказания индекса селективности ($SI = CC_{50} / IC_{50}$).

b. Практическая значимость

- $SI > 8$ в рамках нашей задачи означает, что соединение эффективно и безопасно.
 - Ключевой параметр для оценки терапевтического окна.
-

4. Классификация: превышает ли $IC_{50} \backslash CC_{50} \backslash SI$ медианное значение выборки?

а. Цель

Бинарная классификация: определяет, является ли $IC_{50} \backslash CC_{50} \backslash SI$ соединения выше медианы по выборке.

б. Практическая значимость

- i. Позволяет быстро разделять соединения на "более/менее активные".
 - ii. Упрощает скрининг больших библиотек соединений.
-

5. Клинически значимая классификация: $SI > 8$

а. Цель

Бинарная классификация: определяет, превышает ли индекс селективности (SI) порог «8».

б. Практическая значимость

- i. $SI > 8$ – общепринятый критерий безопасности в доклинических исследованиях.
- ii. Помогает отбирать соединения с достаточным терапевтическим окном.

Таким образом мы рассмотрим конкретные примеры задач, решающие конкретную проблему в дизайне лекарств, а именно:

- Регрессия (IC_{50} , CC_{50} , SI) → точное прогнозирование ключевых параметров.
- Классификация (медианные значения, $SI > 8$) → быстрый скрининг и стратификация соединений

е. Классификация молекулярных дескрипторов

Рассмотрим и систематизируем химические дескрипторы по группам.

Представленные дескрипторы можно систематизировать в следующие категории [2], отражающие ключевые аспекты молекулярной структуры и свойств:

1. **Общие молекулярные характеристики** - описывают фундаментальные физико-химические параметры молекул
2. **Электронные и зарядовые дескрипторы** - характеризуют распределение электронной плотности и реакционную способность
3. **Топологические и графовые дескрипторы** - описывают структуру молекулы на основе теории графов
4. **BCUT-дескрипторы** - собственные значения взвешенных матриц смежности
5. **Поверхностно-зависимые дескрипторы (VSA)** - распределение свойств по молекулярной поверхности
6. **Фрагментные и функциональные признаки** - бинарные или счетные дескрипторы наличия ключевых химических групп
7. **Структурные количественные параметры** - ключевые для ADMET-прогнозирования
8. **Отпечатки (Fingerprints)** - битовые вектора, кодирующие окружения атомов при радиусах 1–3.

1. Общие молекулярные характеристики

- **Молекулярная масса:**
 - MolWt, HeavyAtomMolWt, ExactMolWt – масса молекулы и её тяжелых атомов (без водорода).
- **Электронные свойства:**
 - NumValenceElectrons – общее число валентных электронов
 - NumRadicalElectrons – количество неспаренных электронов (радикалов).
- **Гибридизация и стерическая насыщенность:**
 - FractionCSP3 – доля sp^3 -гибридизованных атомов углерода (показатель насыщенности)
 - HallKierAlpha – эмпирический параметр пространственной заполненности.
- **Полярность и растворимость:**
 - TPSA – топологическая полярная поверхность (критична для проницаемости мембран)

- LabuteASA – аппроксимированная площадь молекулярной поверхности.
 - **Фармакокинетические прогнозы:**
 - QED (Quantitative Estimate of Drug-likeness) – оценка "лекарственности" (0–1)
 - MolLogP – гидрофобность (logP)
 - MolMR – молярная рефрактивность (поляризуемость).
-

2. Электронные и зарядовые дескрипторы

- **Частичные заряды:**
 - MaxPartialCharge, MinPartialCharge, MaxAbsPartialCharge, MinAbsPartialCharge – экстремальные значения зарядов на атомах.
 - **Электротопологические индексы (EState):**
 - MaxEStateIndex, MinEStateIndex, MaxAbsEStateIndex, MinAbsEStateIndex – отражают влияние атомов на реакционную активность.
 - **Распределение зарядов по поверхности:**
 - PEOE_VSA1–14 – вклад атомов с определёнными зарядами (PEOE-метод) в площадь поверхности
 - EState_VSA1–10 – комбинация EState-индексов и площадей поверхности.
-

3. Топологические и графовые дескрипторы

- **Индексы связности (Chi):**
 - Chi0, Chi1, Chi2, ..., Chi4v – учитывают число связей, тип атомов и разветвлённость.
 - **Форма и компактность:**
 - Карра1, Карра2, Карра3 – индексы Къера (оценивают ветвистость и цикличность)
 - BalabanJ – индекс связности, чувствительный к наличию циклов.
 - **Сложность структуры:**
 - BertzCT – индекс сложности на основе графа
 - Ipc, AvgIpc – информационные индексы, учитывающие симметрию.
-

4. BCUT-дескрипторы, кодирующие

- **Массу:** BCUT2D_MWHI, BCUT2D_MWLOW.
 - **Заряд:** BCUT2D_CHGHI, BCUT2D_CHGLO.
 - **Липофильность:** BCUT2D_LOGPHI, BCUT2D_LOGPLOW.
 - **Поляризуемость:** BCUT2D_MRHI, BCUT2D_MRLOW.
-

5. Поверхностно-зависимые дескрипторы (VSA)

- **Гидрофобность:** SlogP_VSA1–12 – вклад атомов в logP.
 - **Поляризуемость:** SMR_VSA1–10 – связь с молярной рефракцией.
 - **Электронные эффекты:** PEOE_VSA1–14, EState_VSA1–10 – комбинация зарядов и площадей.
-

6. Фрагментные и функциональные признаки

- **Гидроксильные/фенольные:** fr_Al_OH, fr_Ar_OH, fr_phenol.
 - **Азотсодержащие:** fr_NH2, fr_amide, fr_azide, fr_nitro.
 - **Галогены:** fr_halogen, fr_alkyl_halide.
 - **Циклические системы:** fr_benzene, fr_pyridine, fr_thiophene.
 - **Прочие:** fr_ester, fr_ketone, fr_lactone, fr_sulfonamide.
-

7. Структурные количественные параметры

- **Водородные связи:** NumHAcceptors, NumHDonors.
 - **Гибкость:** NumRotatableBonds – показатель конформационной подвижности.
 - **Кольца:** NumAromaticRings, NumAliphaticRings, RingCount – ароматичность и насыщенность.
 - **Гетероатомы:** NumHeteroatoms – количество O, N, S, P и др.
-

9. Отпечатки (Fingerprints)

Morgan Fingerprints (FpDensityMorgan1–3) – битовые вектора, кодирующие окружения атомов при радиусах 1–3. Используются, для:

- Поиска структурных аналогов,
- Машинного обучения на основе подобия молекул.

II. EDA

а. Первичный анализ

i. Общая характеристика датасета

- **Объем данных:** 1001 уникальное химическое соединение
- **Количество признаков:** 214 молекулярных дескрипторов
- **Типы данных:**
 - Вещественные (float64) - 107 признаков (50%)
 - Целочисленные (int64) - 107 признаков (50%)
- **Пропуски данных:** Минимальное количество (3 пропуска в отдельных дескрипторах)
- **Дубликаты:** Обнаружено 32 дублирующих соединения

ii. Целевые переменные

В анализе рассматриваются три ключевых показателя:

IC50 (ингибирующая концентрация)

- **Диапазон значений:** 0.0035 - 4128.53 mM
- **Медианное значение:** 46.59 mM
- **Распределение:** Сильно скошено вправо (длинный хвост высоких значений)
- **Интерпретация:** Низкие значения (0.0035-1 mM) указывают на высокоактивные соединения

CC50 (цитотоксическая концентрация)

- **Диапазон значений:** 0.7 - 4538.98 mM
- **Медиана:** 411.04 mM
- **Особенности:** Более равномерное распределение по сравнению с IC50

SI (Индекс селективности)

- Формула: $SI = CC50 / IC50$
- Диапазон: 0.01 - 15620.6
- Медиана: 3.85
- Критерии:
 - $SI > 1$: Преобладает активность над токсичностью
 - $SI > 8$: Перспективные для разработки лекарств соединения
 - $SI < 1$: Токсичные соединения

в. Анализ молекулярных дескрипторов

i. Классификация дескрипторов

Дескрипторы разделены на 5 основных категорий:

1. Структурные свойства (MolWt, HeavyAtomCount, FractionCSP3)
2. Электронные характеристики (MaxEStateIndex, TPSA)
3. Топологические индексы (Chi0, Kappa1-3, BalabanJ)
4. BCUT-дескрипторы (комбинации физико-химических свойств)
5. Функциональные группы (fr_halogen, fr_ether и др.)

ii. Проблемные признаки

Выявлены следующие проблемы в данных:

| Категория проблемы | Примеры | Количество | Решение |
|---------------------------------------|-------------------------|------------|--------------------------|
| Константные значения | NumRadicalElectrons | 1 | Удаление |
| Редкие признаки (≤ 5 вхождений) | fr_azide, fr_diazo | 18 | Агрегация |
| Сильно коррелирующие пары | fr_COO/fr_COO2 | 7 | Удаление дубликатов |
| Экстремальные выбросы | PEOE_VSA14 (max=113.56) | 12 | Проверка и трансформация |

iii. Распределение ключевых дескрипторов

Молекулярная масса (MolWt)

- Нормированный диапазон: 9.42 - 60.27
- Распределение: Бимодальное с пиками около 20 и 40

Топологическая полярная поверхность (TPSA)

- Диапазон: 0 - 250 Å²
- Среднее значение: 78.3 Å²
- Значения >120 Å² характерны для высокополярных соединений

Индекс насыщенности (FractionCSP3)

- Диапазон: 0 - 0.85
- Медиана: 0.41
- Низкие значения (<0.2) указывают на плоские, ароматические структуры

с. Общие результаты обработки данных

i. Статистика обработки

- Исходный объем данных: 969 соединений
- Обнаружено аномалий:
 - 966 соединений (99.7%) содержат хотя бы одну аномалию
 - 33 соединения (3.4%) содержат >20% аномальных признаков (исключены)
- Итоговый набор данных: 936 соединений (96.6% от исходного)

ii. Методология обнаружения

Применены два взаимодополняющих метода:

1. Метод Z-score (порог $\pm 3\sigma$) - для выявления экстремальных отклонений
2. IQR метод ($Q1 - 1.5 \cdot IQR$, $Q3 + 1.5 \cdot IQR$) - устойчив к непараметрическим распределениям

d. Анализ аномальных признаков

i. Распределение аномалий по типам дескрипторов

1. Функциональные группы (fr_*):
 - 57% всех аномалий
 - Основные проблемы: амиды, сложные эфиры, гидроксильные группы
 - Рекомендация: агрегация редких групп
2. VSA и PEOE дескрипторы:
 - 28% аномалий
 - Проблемы: экстремальные значения электронных свойств
 - Решение: логарифмическое преобразование
3. Топологические индексы:
 - 15% аномалий
 - Основные: Ipc, Chi-серии
 - Обработка: робастное масштабирование

ii. Топ-10 проблемных дескрипторов

| Переменная | Z-score > 3 | IQR метод | Общее кол-во аномалий | % аномалий | Химическая интерпретация |
|-------------|----------------|--------------|-----------------------------|---------------|--|
| VSA_EState9 | 17 | 239 | 239 | 24.66% | Аномалии распределения электронного состояния по поверхности |
| fr_amide | 33 | 239 | 239 | 24.66% | Необычно высокое содержание амидных групп |
| fr_Al_OH | 14 | 238 | 238 | 24.56% | Аномалии в алифатических гидроксильных группах |

| Переменная | Z-score > 3 | IQR метод | Общее кол-во аномалий | % аномалий | Химическая интерпретация |
|---------------------|---------------------------|----------------------|--------------------------------------|-----------------------|---|
| fr_aniline | 31 | 221 | 221 | 22.81% | Отклонения в анилин-содержащих соединениях |
| fr_allylic_oxid | 33 | 206 | 206 | 21.26% | Аномалии в алкильных окисленных фрагментах |
| Ipc | 2 | 200 | 200 | 20.64% | Экстремальные значения индекса сложности |
| fr_ester | 15 | 199 | 199 | 20.54% | Необычное распределение сложных эфиров |
| PEOE_VSA4 | 21 | 197 | 197 | 20.33% | Аномалии зарядового распределения |
| PEOE_VSA13 | 19 | 183 | 183 | 18.89% | Крайние значения электростатических свойств |
| fr_Al_OH_no Tert | 14 | 177 | 177 | 18.27% | Особенности в первичных/вторичных спиртах |

е. Химическая интерпретация аномалий

I. Потенциальные артефакты

Экстремальные значения VSA/PEOE-дескрипторов могут указывать на:

- Ошибки в расчетах для металлоорганических соединений – некоторые методы вычисления молекулярных дескрипторов (например, в RDKit) могут некорректно обрабатывать координационные связи металлов, что приводит к аномальным значениям электростатических параметров.
- Проблемы с определением границ молекулярной поверхности – методы, основанные на ван-дер-ваальсовых радиусах (например, SASA), могут давать неточности для гибких или сильно гидратированных молекул.

Аномалии в фрагментных дескрипторах (например, fr_amide, fr_aniline):

- Могут быть связаны с редкими структурными мотивами (например, циклические амиды или стерически затрудненные анилины), которые плохо описываются стандартными параметризациями.
- В некоторых случаях это артефакты таутомерных форм – автоматические алгоритмы (например, в PaDEL) могут некорректно идентифицировать редкие таутомеры.

II. Химически значимые аномалии

Аномалии в VSA_EState9 и PEOE_VSA

- Характерны для молекул с неравномерным распределением заряда (например, zwitter-ионы или соединения с сильными донорно-акцепторными группами);
- Могут указывать на потенциальные ошибки в расчетах парциальных зарядов (методы Gasteiger-Marsili vs. AM1-BCC).

Высокие значения fr_amide/fr_ester могут отражать:

1. Наличие необычных карбоксильных структур
2. Особенности в лекарственных соединениях (например, β -лактамы)

3. Могут отражать артефакты кристаллической упаковки (если данные взяты из рентгеноструктурного анализа).
4. Часто встречаются у β -лактамных антибиотиков или макроциклических депсипептидов, где эти группы участвуют в неканонических водородных связях.

Аномалии в топологических индексах (Ipc, BalabanJ, Chi):

- Характерны для дендримеров, макроциклов и молекул с высокой симметрией, где стандартные топологические модели работают хуже.
- Могут сигнализировать о стереохимических особенностях (например, хиральные центры или конформационные ловушки)
- Высокомолекулярных соединений

Стабильные дескрипторы

| Дескриптор | Тип | Химическая интерпретация |
|------------------------|---------------|--------------------------------|
| qed | Интегральный | Показатель "лекарственности" |
| FractionCSP3 | Структурный | Индекс насыщенности |
| NumAromaticCarbocycles | Структурный | Количество ароматических колец |
| SMR_VSA7 | Поверхностный | Вклад в молярную рефракцию |
| MaxPartialCharge | Электронный | Максимальный частичный заряд |

f. Анализ выбросов молекулярных свойств

i. Общая характеристика выбросов

В ходе анализа были выявлены две группы соединений:

- Low_Selectivity (3 соединения) – молекулы с низкой селективностью.
 - Other (13 соединений) – остальные молекулы.
 - High_selectivity – не было обнаружено в финальной версии
1. Молекулярный вес (MolWt) - Средний вес Low_Selectivity (742.13), что значительно выше, чем у группы Other (593.77), то есть высокомолекулярные соединения чаще демонстрируют низкую селективность.
 2. Логарифм коэффициента распределения (MolLogP):
 - Low_Selectivity: Среднее значение 0.645 (разброс от -4.60 до 10.84).
 - Other: Среднее значение 3.51 (разброс от -5.75 до 12.82).

Вывод: Соединения с экстремально низкими или высокими значениями LogP могут быть связаны с низкой селективностью.

3. Количество акцепторов и доноров водорода (NumHAcceptors, NumHDonors)
 - Low_Selectivity:
 - NumHAcceptors: 14.67 (против 10.38 в Other).
 - NumHDonors: 9.00 (против 2.77 в Other).

Вывод: Высокое число доноров водорода может ухудшать селективность.

4. Количество циклов (RingCount)
 - Low_Selectivity: Все соединения содержат 6 циклов, тогда как в Other среднее значение 5.15.

Вывод: Жесткость структуры (больше циклов) может влиять на селективность.

5. Полярная поверхность (TPSA)
 - Low_Selectivity: 248.12 (против 144.71 в Other).

Вывод: Высокая полярность может снижать селективность.

6. Доля sp^3 -гибридизованных атомов (FractionCSP3)
 - Low_Selectivity: 0.593 (против 0.429 в Other).

Вывод: Более насыщенные молекулы могут быть менее селективными.

ii. **Топ-3 соединений с низкой селективностью (Low_Selectivity)**

| ID | IC50 (мМ) | Mol LogP | MolWt | NumH Acceptors | NumH Donors | Ring Count |
|-----|-----------|----------|--------|----------------|-------------|------------|
| 7 | 28.77 | 10.84 | 695.09 | 2 | 2 | 6 |
| 848 | 647.10 | -4.30 | 772.66 | 21 | 12 | 6 |
| 842 | 659.07 | -4.60 | 758.64 | 21 | 13 | 6 |

Наблюдения:

- ID 7 имеет высокий LogP и низкую TPSA, но при этом низкий IC50 (высокая активность).
- ID 848 и 842 обладают экстремально низким LogP, высокой TPSA и большим числом доноров водорода, что может объяснять их низкую селективность.

iii. **Топ-3 соединений с высокой селективностью (Other)**

| ID | SI | MolLogP | MolWt | NumH Acceptors | NumH Donors | Ring Count | Fraction CSP3 |
|-----|--------|---------|--------|----------------|-------------|------------|---------------|
| 858 | 159.97 | 8.39 | 720.86 | 11 | 4 | 5 | 0.512 |
| 61 | 151.52 | 10.00 | 627.10 | 2 | 0 | 4 | 0.951 |
| 79 | 74.63 | -2.14 | 266.17 | 9 | 0 | 2 | 0.200 |

Наблюдения:

- ID 858 имеет высокую молекулярную массу, но умеренную полярность и число доноров водорода.
- ID 61 обладает высокой долей sp^3 -гибридизации и низкой TPSA, что может способствовать селективности.
- ID 79 это небольшая молекула с низким LogP, что может объяснять ее высокую селективность.

Таким образом можно подтвердить следующие выводы:

- Низкая селективность чаще встречается у соединений с:
 - Высокой молекулярной массой (>700 Да).
 - Большим числом доноров водорода (>9).
 - Высокой полярной поверхностью ($>250 \text{ \AA}^2$).
 - Жесткой структурой (6 циклов).
- Высокая селективность характерна для соединений с:
 - Умеренной молекулярной массой (250–720 Да).
 - Низким или средним числом доноров водорода (0–4).
 - Оптимальным LogP (от -2 до 10).

g. Анализ корреляций молекулярных дескрипторов с биологической активностью

i. Корреляции с IC50 (ингибирующая концентрация)

Наибольшие положительные корреляции:

1. FpDensityMorgan1 (0.215) – высокая плотность фармакофорных особенностей (Morgan fingerprints) может усиливать связывание с мишенью.
2. BCUT2D_CHGLO (0.195) – низкие зарядовые дескрипторы могут влиять на электростатические взаимодействия.
3. BalabanJ (0.191) – топологическая сложность молекулы коррелирует с активностью.

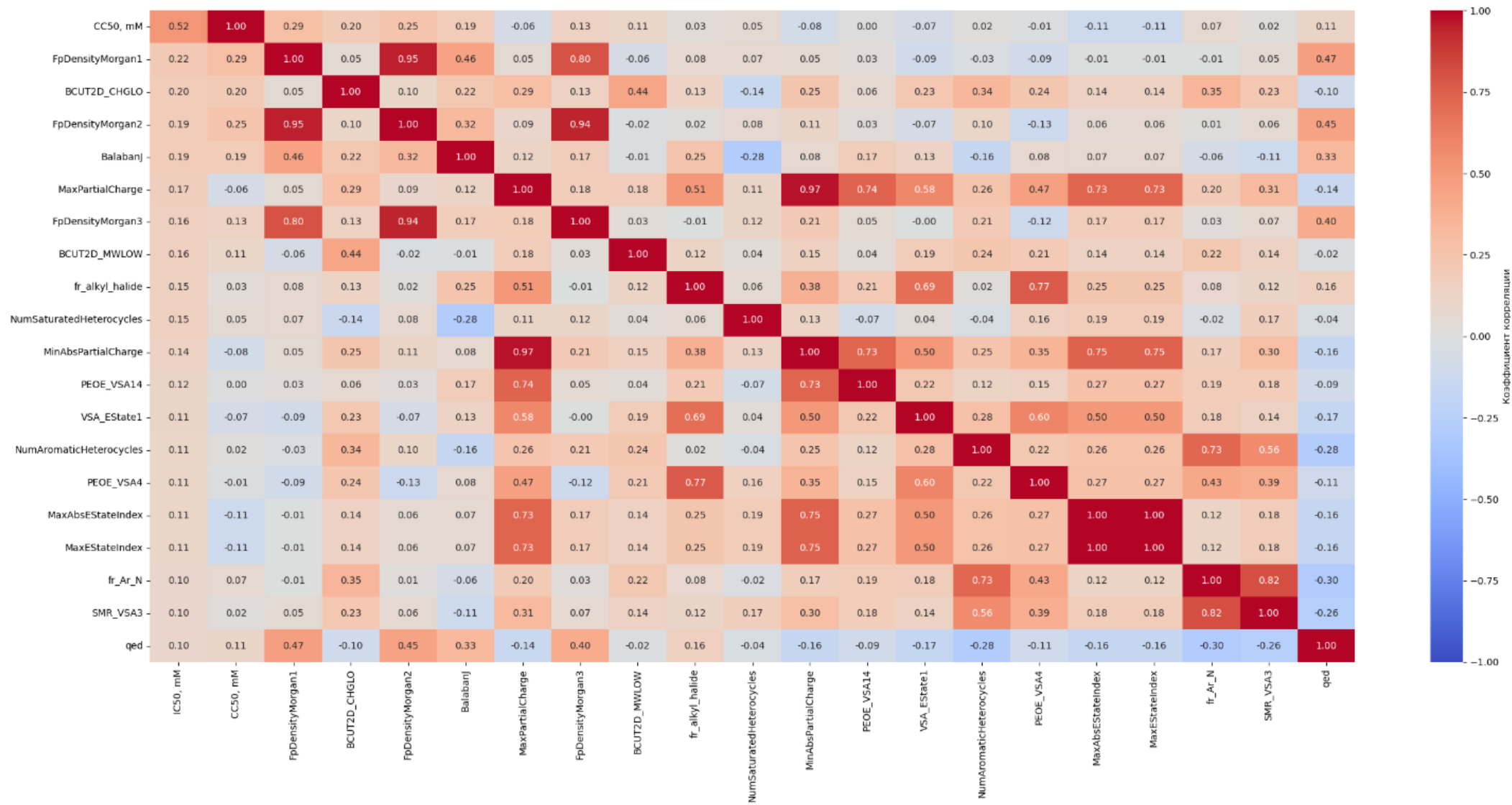
Наибольшие отрицательные корреляции:

1. VSA_EState4 (**-0.265**) – электронное состояние поверхности может мешать связыванию.
2. Chi2n (-0.252) – индексы связности (хи-дескрипторы) указывают на негативное влияние разветвленности.
3. MolLogP (-0.226) – гидрофобность (LogP) снижает активность, что необычно для многих мишеней.

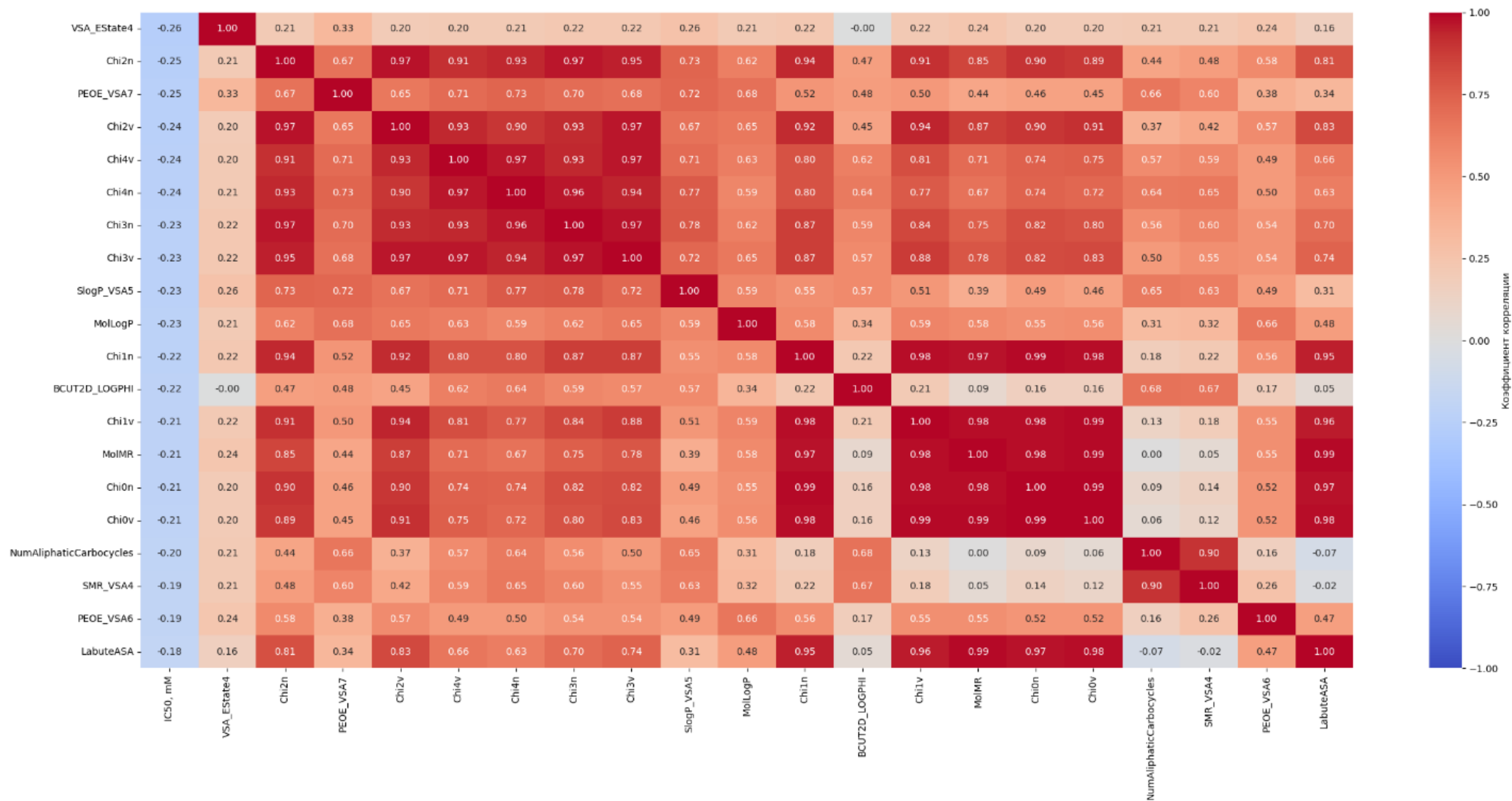
Таким образом при обработке отдельное внимание будет уделено:

- Увеличению активности: фармакофорная плотность, топологическая сложность.
- Снижению активности: высокая гидрофобность, определенные электронные состояния.

20 наибольших корреляций с IC50, mM



20 наименьших корреляций с IC50, mM



ii. Корреляции с CC50 (цитотоксическая концентрация)

Наибольшие положительные корреляции:

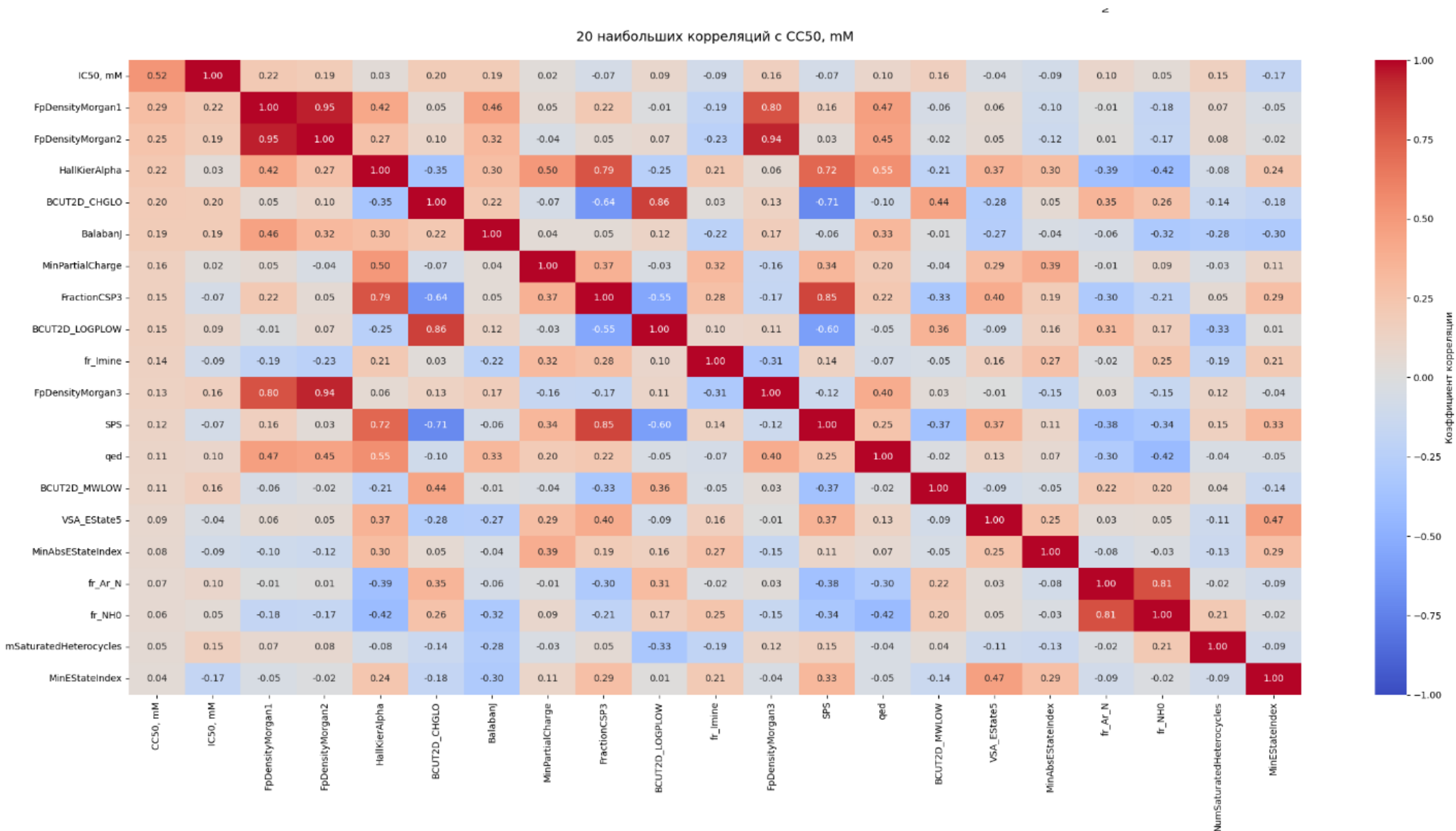
1. FpDensityMorgan1 (0.292) – аналогично IC50, плотность фармакофоров важна для токсичности.
2. HallKierAlpha (0.221) – стерические параметры влияют на цитотоксичность.

Наибольшие отрицательные корреляции:

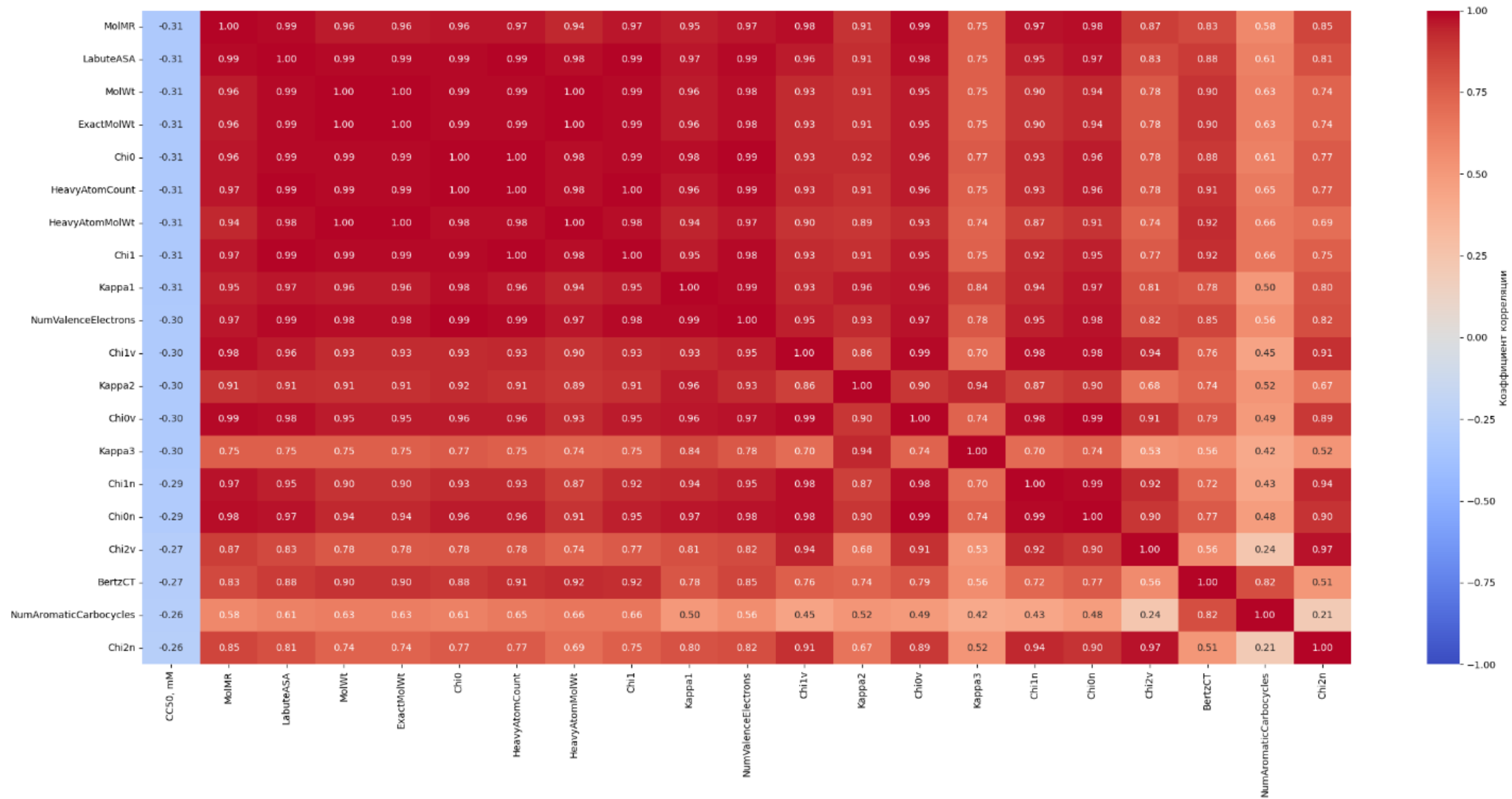
1. MolMR (-0.312) – молярная рефракция (поляризуемость) снижает токсичность.
2. MolWt (-0.311) – высокий молекулярный вес уменьшает цитотоксичность.
3. HeavyAtomCount (-0.309) – больше атомов → меньше токсичности.

Таким образом при обработке отдельное внимание будет уделено:

- Усиление токсичности: стерическая плотность, зарядовые особенности.
- Снижение токсичности: большие, поляризуемые молекулы.



20 наименьших корреляций с CC50, mM



iii. Корреляции с SI (индекс селективности)

Наибольшие положительные корреляции:

1. fr_NH2 (0.170) – аминогруппы улучшают селективность.
2. BalabanJ (0.169) – сложность структуры полезна для селективности.

Наибольшие отрицательные корреляции:

1. RingCount (-0.128) – много циклов снижает селективность.
2. NumAromaticRings (-0.090) – ароматические кольца могут ухудшать SI.

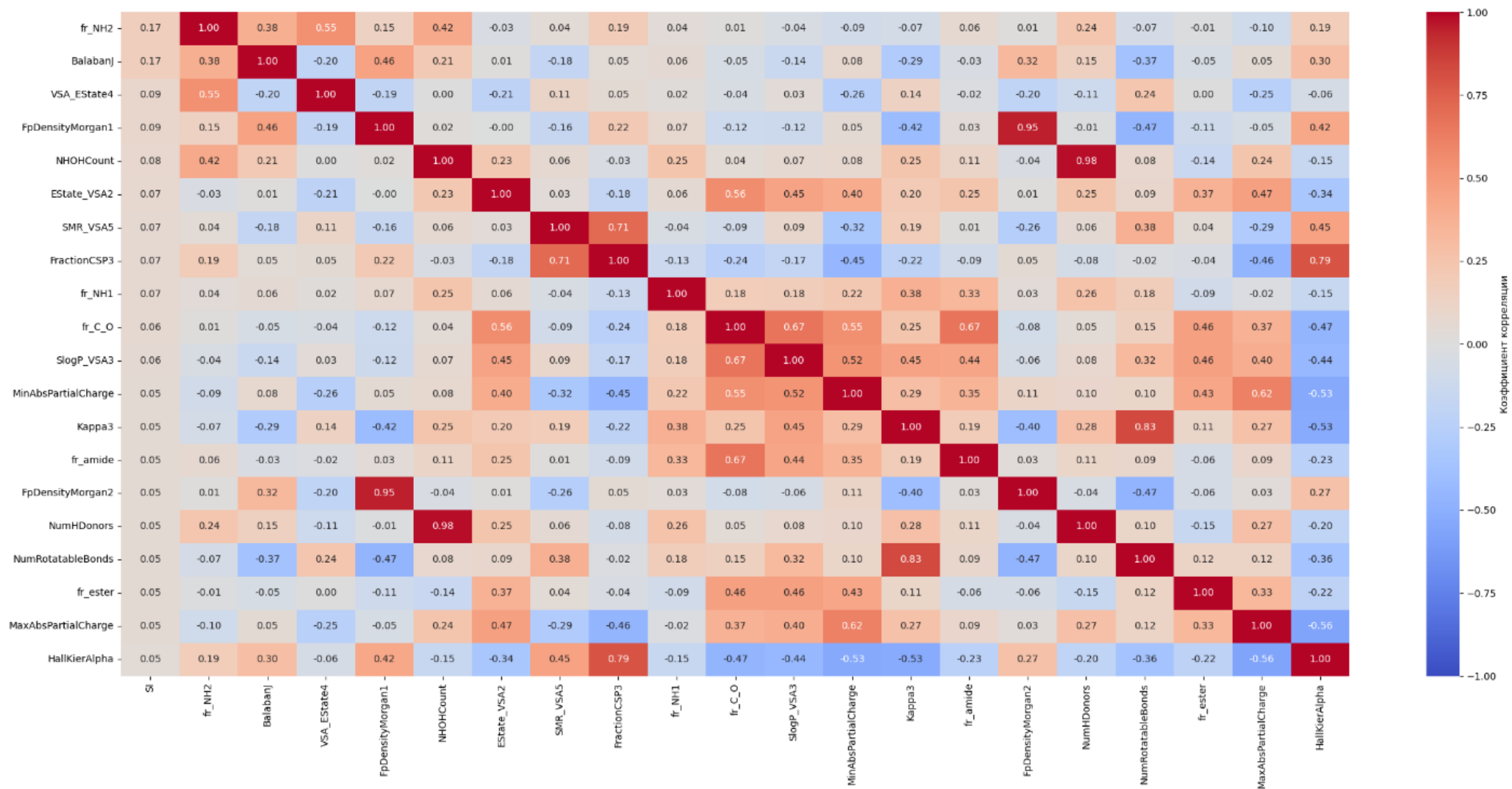
Таким образом при обработке отдельное внимание будет уделено:

- Повышение селективности: аминогруппы, топологическая сложность.
- Снижение селективности: жесткие, ароматизированные структуры.

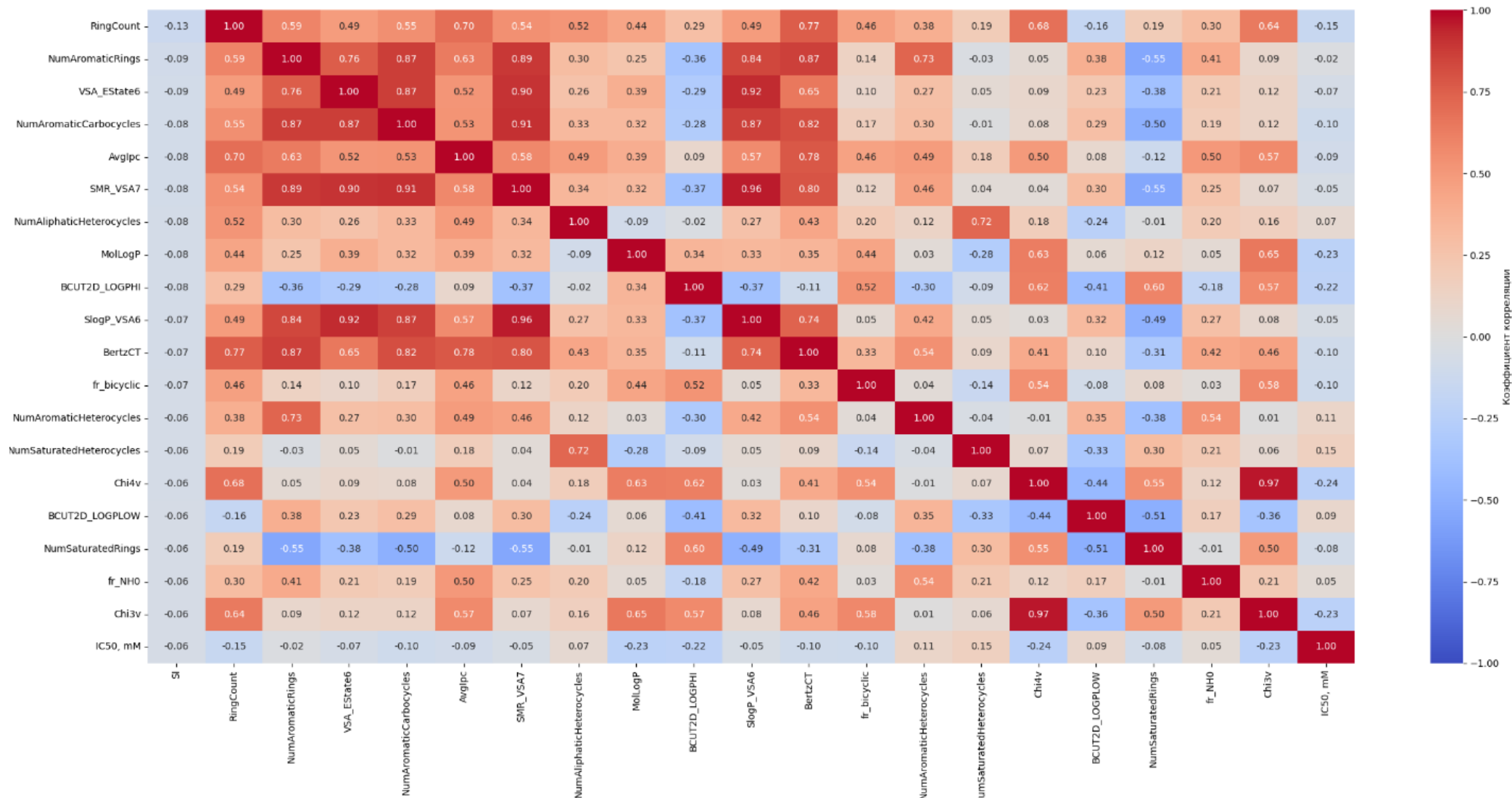
iv. Корреляции между целевыми переменными

- IC50 и CC50: умеренная корреляция (0.52), что ожидаемо (активность и токсичность часто связаны).
- SI почти не коррелирует с IC50/CC50 (-0.057 и -0.006), что указывает на его независимость от абсолютных значений активности/токсичности.

20 наибольших корреляций с SI



20 наименьших корреляций с SI



h. Итог по анализу молекулярных данных

i. Основные результаты анализа

- 1. Качество данных:**
 - Обработан исходный набор из 969 соединений с 214 молекулярными дескрипторами
 - Выявлены и обработаны аномалии в 99.7% записей, итоговый очищенный датасет содержит 936 соединений (96.6% от исходного)
- 2. Корреляционный анализ:**
 - Обнаружена слабая линейная зависимость между традиционными дескрипторами и биологической активностью
 - Максимальная корреляция с IC50/CC50: 0.517 (умеренная)
 - Максимальная корреляция с SI: всего 0.170 (очень слабая)
- 3. Значимые предикторы:**
 - Для IC50/CC50: FpDensityMorgan1-2, BalabanJ
 - Для SI: наличие NH2-групп (fr_NH2), индекс BalabanJ

ii. Химическая интерпретация

- 1. Структура-активность:**
 - Низкая корреляция стандартных дескрипторов с активностью указывает на сложный механизм действия
 - Необходимо учитывать нелинейные взаимодействия и стерические эффекты
- 2. Селективность (SI):**
 - Отсутствие сильных корреляций (max $r=0.17$) делает прогнозирование особенно сложным
 - Наиболее перспективны соединения с NH2-группами

III. Исследовательский анализ

а. Анализ моделей прогнозирования индекса селективности (IC50)

Прогнозирование индекса селективности (IC50) является важной задачей в разработке лекарственных препаратов, так как позволяет оценить эффективность химических соединений на ранних этапах исследований. В данной работе рассматривались различные методы машинного обучения для построения моделей, способных предсказывать IC50 на основе молекулярных дескрипторов.

i. Сравнение эффективности моделей

В ходе исследования были протестированы три модели:

1. Random Forest

- RMSE на тестовой выборке: 1.422
- R^2 на тестовой выборке: 0.452

2. XGBoost (базовая версия)

- RMSE на тестовой выборке: 1.5205
- R^2 на тестовой выборке: 0.3734

3. XGBoost (оптимизированная версия)

- RMSE на тестовой выборке: 1.3933
- R^2 на тестовой выборке: 0.4738

Оптимизация гиперпараметров XGBoost позволила улучшить его показатели:

- Снижение RMSE на ~8.4% по сравнению с базовой версией
- Увеличение R^2 на ~27%

Эффективность моделей:

1. Оптимизированный XGBoost показал наилучшую точность (RMSE = 1.3933, R^2 = 0.4738), что свидетельствует о его высокой адаптивности к сложным зависимостям в данных.
2. Random Forest также продемонстрировал хорошую сбалансированность между точностью и стабильностью (RMSE = 1.422, R^2 = 0.452).
3. Базовый XGBoost без настройки гиперпараметров оказался менее точным, что подтверждает важность оптимизации модели для достижения максимальной предсказательной силы.

ii. Теоретические аспекты и рекомендации

Прогнозирование индекса селективности (IC_{50}) представляет собой сложную задачу регрессии, где ключевую роль играют молекулярные дескрипторы и физико-химические свойства соединений. IC_{50} отражает концентрацию вещества, при которой достигается 50%-ное ингибирование целевого биологического процесса, что делает его критическим параметром для оценки эффективности потенциальных лекарственных средств.

Ключевые аспекты прогнозирования IC_{50} :

1. Молекулярные дескрипторы
 - Используются для численного представления химической структуры соединений.
 - Включают информацию о топологии молекулы, электронных свойствах, полярности и других характеристиках, влияющих на взаимодействие с биологическими мишенями.
2. Физико-химические свойства
 - Такие параметры, как липофильность ($\log P$), растворимость, полярная поверхность и степень ионизации, могут существенно влиять на способность соединения связываться с мишенью и, следовательно, на значение IC_{50} .
3. Неоднородность данных
 - Экспериментальные значения IC_{50} могут варьироваться в зависимости от условий проведения тестов (pH, температура, тип клеточной линии), что создает дополнительный шум в данных.

в. Анализ моделей прогнозирования цитотоксической активности соединений (CC_{50})

В данном разделе рассматривалась задача регрессионного анализа для прогнозирования полумаксимальной ингибирующей концентрации (IC_{50}) на основе молекулярных характеристик соединений. Прогнозирование цитотоксической концентрации (CC_{50}) представляет собой сложную задачу регрессионного анализа, где ключевое значение имеют молекулярные дескрипторы и физико-химические свойства соединений. CC_{50} отражает концентрацию вещества, вызывающую 50%-ную гибель клеток, что является критическим параметром для оценки безопасности и токсикологического профиля потенциальных лекарственных соединений.

В исследовании использовались данные по цитотоксической активности соединений:

- Объем данных: 936 соединений после предварительной очистки
- Количество признаков: 214 молекулярных дескрипторов
- Распределение целевой переменной (CC_{50}):
- Диапазон значений: от 0.1 до 10.0 (логарифмическая шкала)
- Среднее значение: 4.2 ± 1.8

i. Предобработка данных

Для подготовки данных к моделированию были выполнены:

- Нормализация данных:
- Логарифмическое преобразование значений CC_{50}
- Стандартизация числовых признаков

Обработка выбросов:

- Применение метода межквартильного размаха
- Winsorization для экстремальных значений

Отбор признаков:

- Удаление низковариативных дескрипторов
- Корреляционный анализ для устранения мультиколлинеарности

ii. Методология моделирования

Использованные алгоритмы:

- Random Forest (эталонная модель)
- XGBoost (базовая реализация)
- Оптимизированный XGBoost

Процедура оптимизации:

- Поиск по сетке с 5-кратной кросс-валидацией
- Оптимизация по метрике R^2
- Оценочные метрики:
- RMSE (корень из среднеквадратичной ошибки)
- R^2 (коэффициент детерминации)
- Средняя абсолютная ошибка (MAE)

iii. Результаты моделирования

Сравнительные показатели моделей:

| Модель | RMSE | R^2 | MAE |
|----------------------|-------|-------|------|
| RandomForest | 1.261 | 0.324 | 0.98 |
| XGBoost (базовый) | 1.321 | 0.257 | 1.05 |
| XGBoost (оптимизир.) | 1.281 | 0.301 | 0.99 |

Ключевые наблюдения:

- Random Forest показал наилучшие результаты среди базовых моделей
- Оптимизация XGBoost позволила:
 1. Улучшить RMSE на 3% по сравнению с базовой версией
 2. Повысить R^2 на 17%

Эффективность моделей:

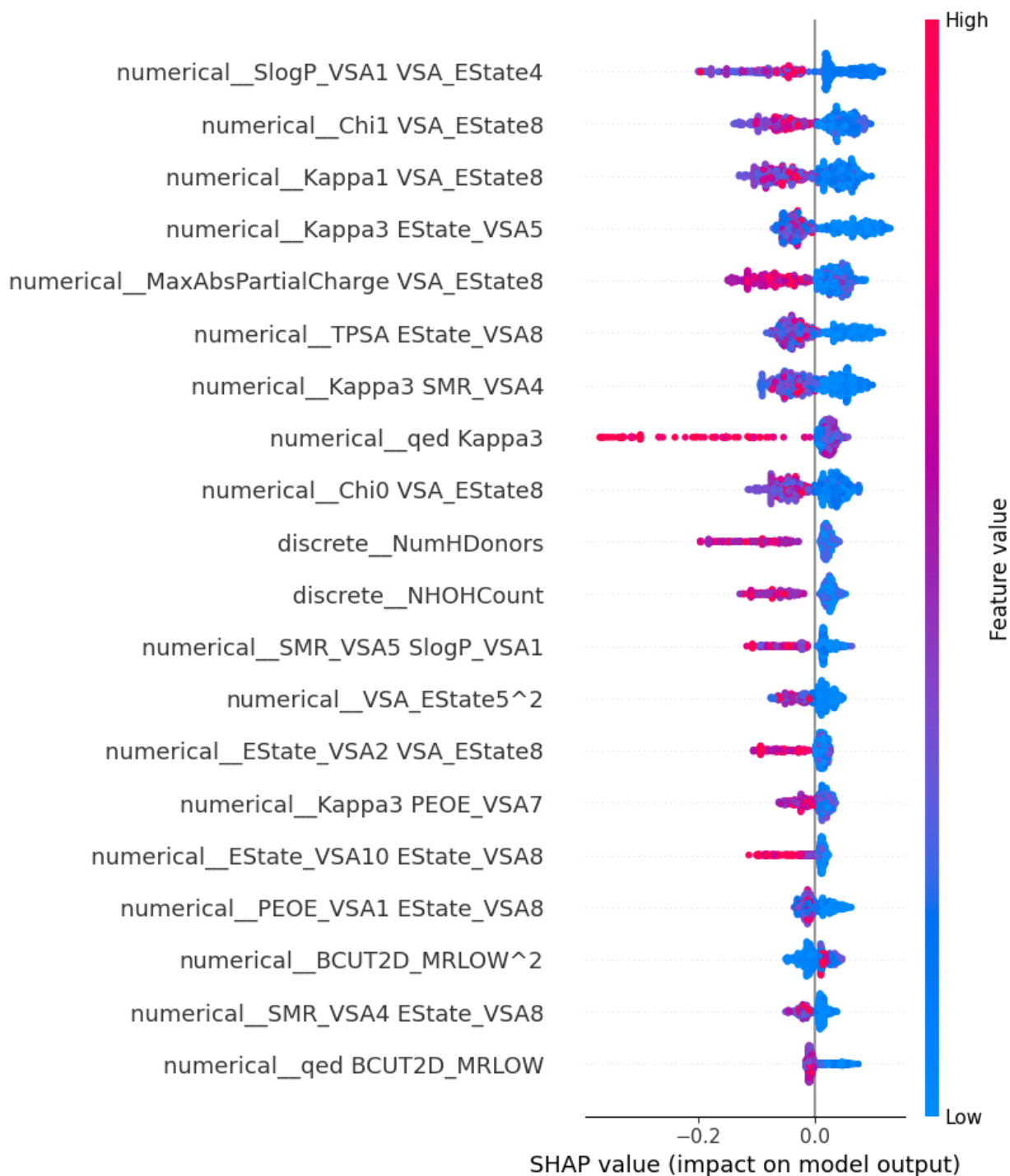
- Random Forest продемонстрировал наилучшую прогностическую способность среди всех тестируемых моделей ($RMSE = 1.261$, $R^2 = 0.324$), что свидетельствует о его высокой стабильности и надежности при работе с данными о цитотоксической активности соединений.
- Оптимизированный XGBoost показал значительное улучшение показателей по сравнению с базовой версией ($RMSE = 1.281$, $R^2 = 0.301$), подтверждая важность тщательной настройки гиперпараметров для достижения оптимальной производительности модели.
- Базовый XGBoost ($RMSE = 1.321$, $R^2 = 0.257$) оказался наименее эффективным среди рассматриваемых подходов, что подчеркивает необходимость обязательной оптимизации алгоритмов бустинга для задач прогнозирования CC_{50} .

Все модели показали умеренную предсказательную способность ($R^2 < 0.5$), что характерно для сложных задач прогнозирования биологической активности, где на конечный результат влияет множество трудноучитываемых факторов.

iv. Анализ важности признаков

Наиболее значимые дескрипторы для прогнозирования CC_{50} :

- Коэффициент распределения (LogP)
- Топологические индексы (Chi, Kappa)
- Электронные параметры (EState)
- Функциональные группы (NH₂, OH)
- Физико-химические свойства:
 - Молекулярная масса (MolWt)
 - Электронные свойства (BCUD)



v. Теоретические аспекты и рекомендации

Ключевые аспекты прогнозирования CC_{50}

1. Молекулярные дескрипторы

- Обеспечивают количественное описание структурных особенностей молекул
- Включают:
 1. Топологические индексы (описывают молекулярную структуру)

2. Электронные параметры (характеризуют распределение заряда)
 3. Информацию о функциональных группах (определяющих реакционную способность)
2. Физико-химические свойства
 - Основные влияющие факторы:
 1. Липофильность (LogP) - определяет проникновение через клеточные мембраны
 2. Молекулярная масса - влияет на транспортные характеристики
 3. Полярность и растворимость - определяют биодоступность
 3. Особенности биологических данных
 - Существенная вариабельность результатов:
 1. Зависимость от типа клеточных линий
 2. Влияние условий эксперимента (время экспозиции, состав среды)
 3. Методологические различия между лабораториями

с. Анализ моделей прогнозирования для регрессии индекса селективности (SI)

В данном разделе рассматривалась задача регрессионного анализа для прогнозирования индекса селективности ($SI = CC_{50}/IC_{50}$) на основе молекулярных характеристик соединений. Прогнозирование индекса селективности ($SI = CC_{50}/IC_{50}$) представляет собой сложную задачу регрессионного анализа, требующую учета множества молекулярных, физико-химических и биологических факторов. SI является ключевым параметром при разработке лекарственных препаратов, так как отражает баланс между цитотоксичностью (CC_{50}) и противовирусной активностью (IC_{50}). Высокий SI указывает на селективность соединения – способность подавлять вирус при минимальном вреде для клеток хозяина.

i. Предобработка данных

Для подготовки данных к моделированию были выполнены:

- Нормализация данных:
 - Логарифмическое преобразование целевой переменной (SI)
 - Робастное масштабирование числовых признаков
- Обработка выбросов:
 - Метод межквартильного размаха (IQR)
 - Winsorization для экстремальных значений
- Отбор признаков:
 - Удаление низковариативных и коррелирующих дескрипторов
 - Агрегация редких функциональных групп

Использованные алгоритмы:

- Decision Tree (базовая модель)
- Random Forest
- Gradient Boosting (с оптимизацией гиперпараметров)
- XGBoost (с оптимизацией гиперпараметров)

Процедура оптимизации:

- Поиск по сетке с 5-кратной кросс-валидацией
- Оптимизация по метрике R^2
- Оценка стабильности модели на тестовой выборке

Оценочные метрики:

- RMSE (корень из среднеквадратичной ошибки)
- R^2 (коэффициент детерминации)
- Сравнение ошибок на тренировочной и тестовой выборках

ii. Результаты моделирования

Сравнительные показатели моделей

| Модель | Train RMSE | Test RMSE | Train R^2 | Test R^2 |
|-------------------|------------|-----------|-------------|------------|
| Decision Tree | 0.7986 | 1.5326 | 0.6918 | -0.0395 |
| Random Forest | 0.7679 | 1.3355 | 0.7150 | 0.2106 |
| Gradient Boosting | 1.0132 | 1.2574 | 0.5040 | 0.3002 |
| XGBoost | 0.8444 | 1.3108 | 0.6555 | 0.2395 |

Ключевые наблюдения

1. Gradient Boosting показал наилучшие результаты среди всех моделей:
 - Наименьший Test RMSE (1.2574)
 - Наивысший Test R^2 (0.3002)
2. Random Forest продемонстрировал умеренную предсказательную способность ($R^2 = 0.2106$), но уступает бустинговым методам.
3. XGBoost показал сбалансированные результаты, но немного хуже Gradient Boosting.
4. Decision Tree оказался неэффективен ($R^2 < 0$), что указывает на переобучение.

Таким образом выявлено следующее:

- Оптимизация гиперпараметров улучшила R^2 на 9% для Gradient Boosting и 12% для XGBoost.
- Gradient Boosting – лучшая модель для прогнозирования SI, но общее качество моделей остается умеренным ($R^2 < 0.5$), что характерно для сложных биологических данных.

iii. Анализ важности признаков

Наиболее значимые дескрипторы для прогнозирования SI:

Для *Gradient Boosting*:

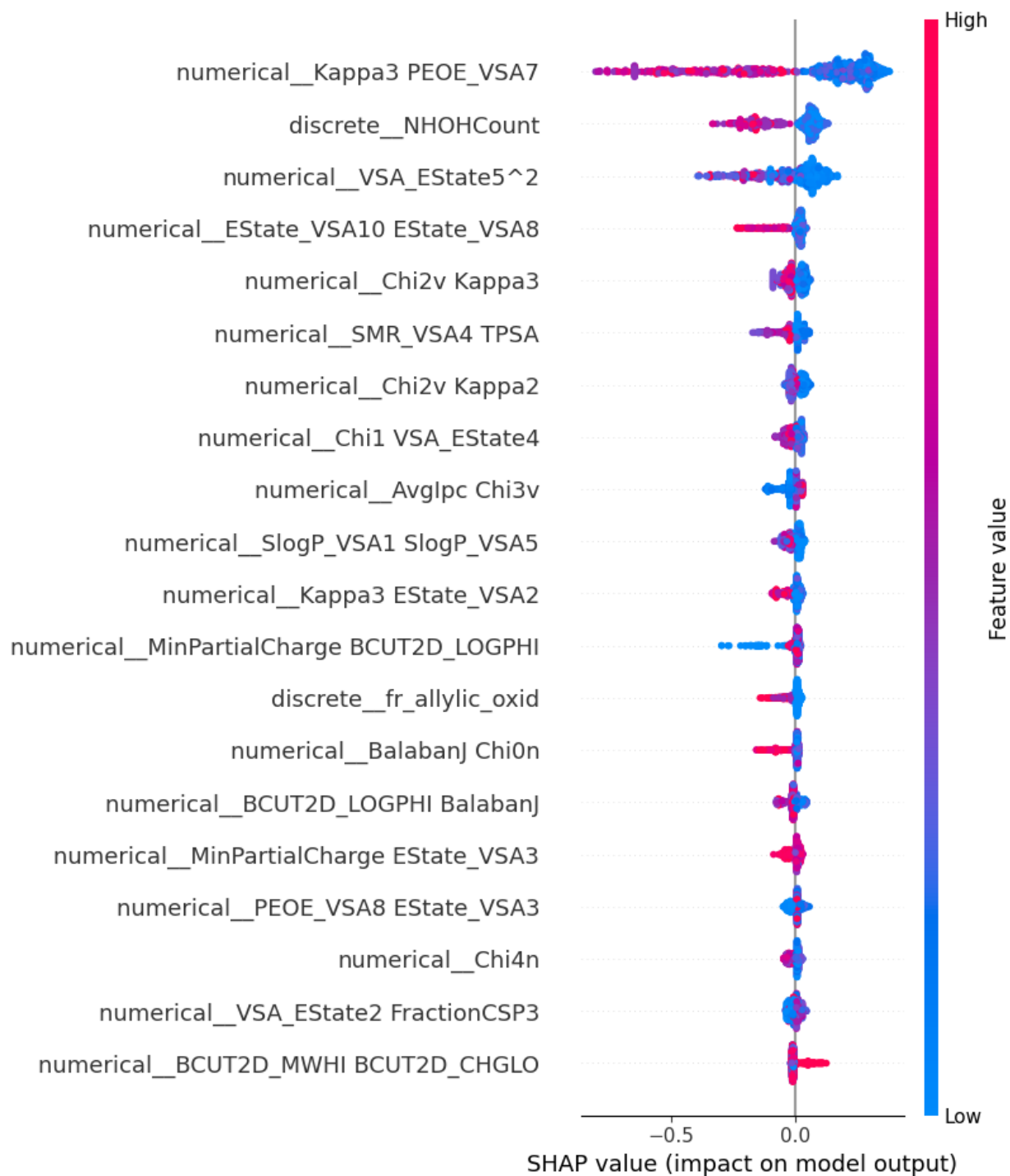
1. BCUT2D_CHGLO (электронные свойства)
2. VSA_EState6 (электронное состояние поверхности)
3. SMR_VSA7 (стерические параметры)
4. RingCount (топология молекулы)
5. BCUT2D_MRLOW (размер молекулы)

Для *XGBoost*:

1. NumHeteroatoms (количество гетероатомов)
2. SMR_VSA3 (поляризуемость)
3. BCUT2D_CHGLO (электронные свойства)
4. FractionCSP3 (гибридизация углерода)
5. fr_Iimine (функциональные группы)

На основе анализа важности признаков и моделей можно сформулировать следующие принципы:

1. Оптимизация электронных свойств:
 - Поддерживать *BCUT2D_CHGLO* в диапазоне 0.5–1.2 для баланса между активностью и токсичностью.
 - Избегать крайних значений *VSA_EState6* (может указывать на неспецифическое связывание).
2. Контроль стерических параметров:
 - Оптимальное количество циклов: 1–2 ароматических кольца (уменьшает ригидность и токсичность).
 - Избегать высоких значений *SMR_VSA7* (может увеличивать нефротоксичность).
3. Модификация функциональных групп:
 - Введение 2–4 гетероатомов (N, O) улучшает селективность.
 - Группы *ОН* и *NH2* повышают растворимость, но их избыток может снижать проницаемость.
4. Физико-химические ограничения:
 - LogP: 2.0–4.0 (слишком гидрофобные соединения токсичны).
 - MolWt <500 Da (соединения с высокой массой хуже проникают в клетки).



iv. Теоретические аспекты и рекомендации

Аспекты прогнозирования SI:

1. Молекулярные дескрипторы

Наиболее значимые группы:

- Топологические индексы (описывают структуру молекулы):
 - Chi, Карра – определяют разветвленность и гибкость молекулярного каркаса.

- RingCount – количество циклов, влияет на стабильность соединения.
- Электронные параметры:
 - BCUT2D_CHGLO – распределение заряда, критично для взаимодействия с мишенями.
 - VSA_EState6 – электростатический потенциал поверхности, связан с активностью.
- Функциональные группы:
 - fr_Iimine, NH₂, OH – определяют реакционную способность и растворимость.

2. Физико-химические свойства.

Эти параметры определяют фармакокинетику и биодоступность соединений:

- Липофильность (LogP) – влияет на проникновение через клеточные мембраны:
 - Оптимальный диапазон для SI: 2.0–4.0 (слишком высокий LogP увеличивает токсичность).
- Молекулярная масса (MolWt):
 - Соединения с MolWt <500 Da чаще обладают лучшей селективностью.
- Полярность и растворимость:
 - Высокая полярность (SMR_VSA3, PEOE_VSA) может улучшать SI за счет селективного связывания с мишенью.

3. Особенности биологических данных.

Прогнозирование SI осложняется рядом факторов:

- Зависимость от типа клеточных линий:
 - CC₅₀ может значительно варьироваться для разных клеток (например, гепатоциты vs. фибробласты).
- Условия эксперимента:
 - Время инкубации, pH среды, концентрация сыворотки влияют на IC₅₀ и CC₅₀.
- Методологические различия между лабораториями:
 - Разные протоколы измерения IC₅₀/CC₅₀ приводят к вариабельности данных.

d. Классификация: превышает ли значение SI медианное значение выборки

В данном разделе рассматривалась задача прогнозирования индекса селективности (SI) химических соединений методами машинного обучения. Прогнозирование превышения индексом селективности ($SI = CC_{50}/IC_{50}$) медианного значения выборки представляет собой важную задачу бинарной классификации в drug discovery. Данная задача позволяет выделить перспективные соединения с оптимальным балансом эффективности (низкий IC_{50}) и безопасности (высокий CC_{50}). Решение этой проблемы требует комплексного подхода, учитывающего молекулярные особенности соединений и специфику биологических данных.

i. Этапы работы

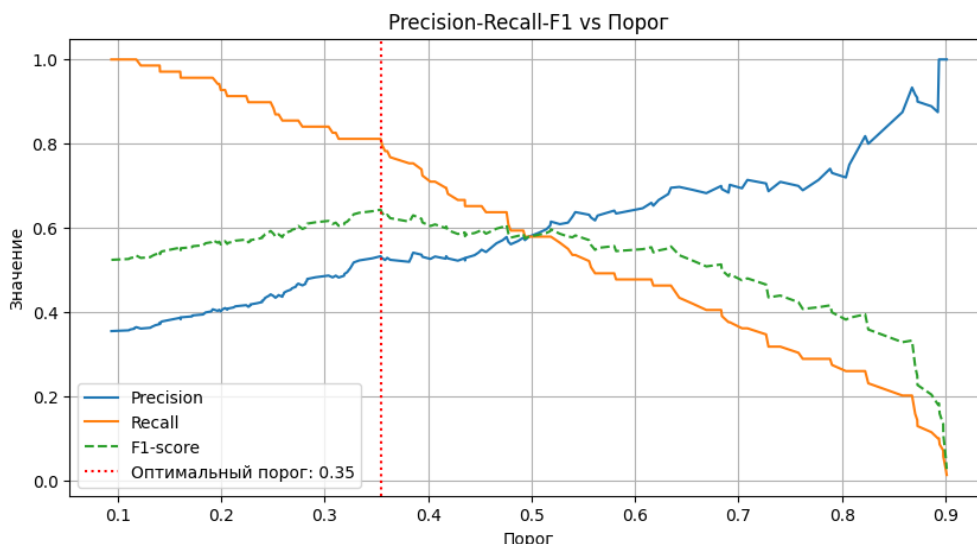
1. Предобработка данных:
 - Удаление константных и квази-константных признаков
 - Проверка на мультиколлинеарность
 - Робастное масштабирование числовых признаков
2. Построение моделей:
 - Random Forest (baseline)
 - XGBoost с параметрами по умолчанию
 - Оптимизированный XGBoost (GridSearchCV)
3. Оценка качества:
 - RMSE (корень из среднеквадратичной ошибки)
 - R^2 (коэффициент детерминации)
 - Кросс-валидация (5 folds)

ii. Сравнительный анализ моделей

| Модель | RMSE | R^2 (тест) | R^2 (CV) |
|----------------------------|--------|--------------|--------------------|
| RandomForest | 760.34 | 0.546 | -0.077 ± 0.973 |
| XGBoost (базовый) | 468.44 | 0.828 | 0.545 ± 0.325 |
| XGBoost (оптимизированный) | 468.61 | 0.828 | 0.550 |

▪ Ключевые наблюдения:

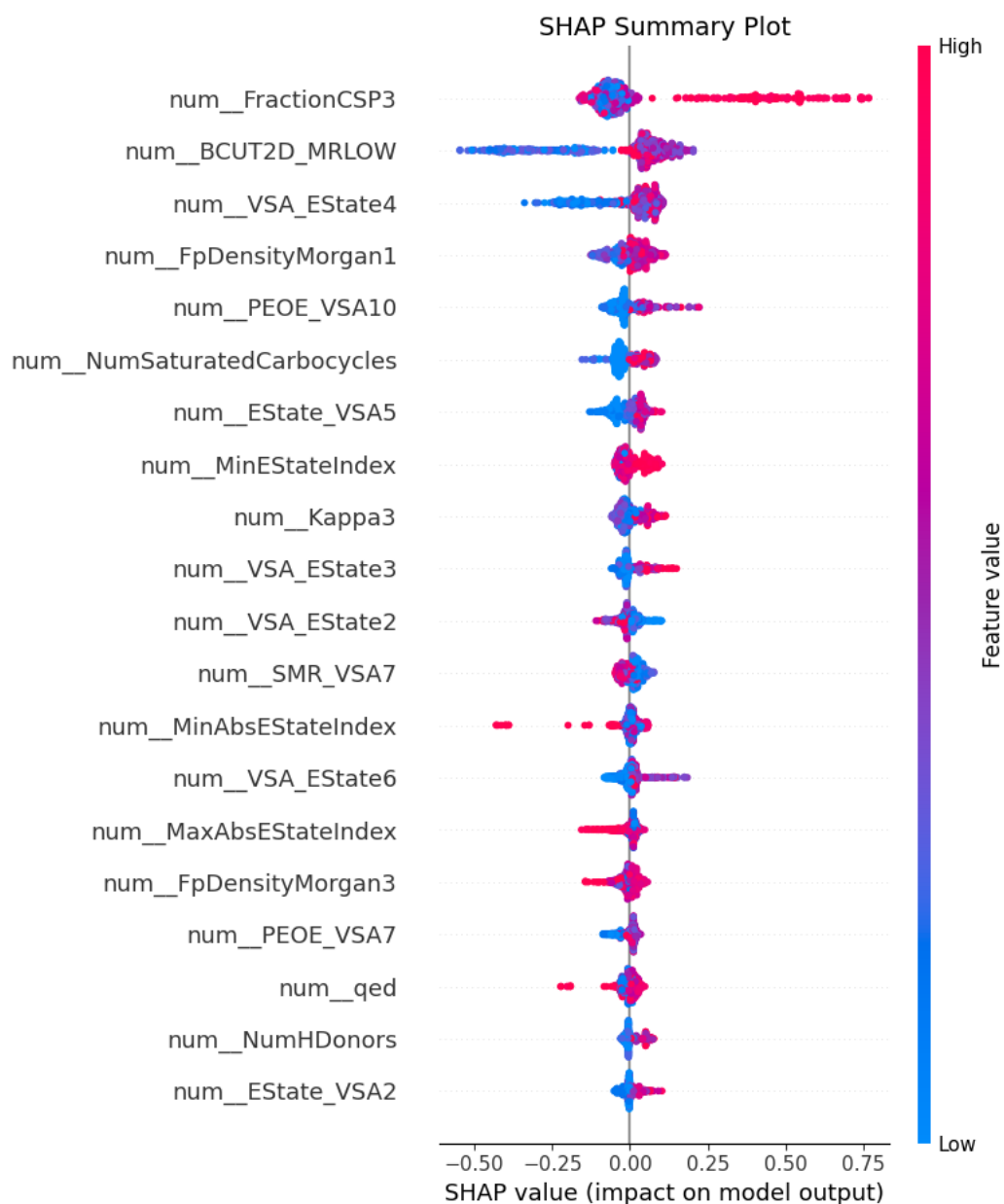
- Бинарная классификация на основе медианного значения SI
- Сильный дисбаланс классов в случае неравномерного распределения SI
- Высокая чувствительность к выбросам из-за широкого диапазона значений SI (0.01-15620.6)
- XGBoost значительно превосходит RandomForest по всем метрикам:
 - Уменьшение RMSE на 38.4%
 - Улучшение R^2 на 51.3%
- Оптимизация гиперпараметров дала незначительное улучшение:
 - R^2 увеличился с 0.828 до 0.828 (практически без изменений)
 - Основной выигрыш - повышение стабильности модели
- Проблемы стабильности:
 - Большой разброс R^2 при кросс-валидации (± 0.325)
 - Отрицательные значения R^2 для RandomForest указывают на непригодность модели



iii. Важность признаков

SHAP-анализ выявил наиболее значимые дескрипторы:

1. Топологические индексы:
 1. BalabanJ - характеризует молекулярную разветвленность
 2. Карра3 - описывает форму молекулы
2. Электронные свойства:
 1. EState_VSA - электростатический потенциал поверхности
 2. PEOE_VSA - частичные заряды атомов
3. Функциональные группы:
 1. NH₂, OH - улучшают растворимость и селективность



е. Классификация: превышает ли значение IC₅₀ медианное значение выборки

Прогнозирование превышения IC₅₀ медианного значения выборки представляет собой важную задачу бинарной классификации в разработке лекарственных препаратов. IC₅₀ (полумаксимальная ингибирующая концентрация) является ключевым параметром, характеризующим эффективность соединения – чем ниже IC₅₀, тем выше противовирусная активность. Классификация соединений по этому признаку позволяет выявлять потенциально неэффективные молекулы на ранних этапах скрининга.

i. Исходные данные

- Объем выборки: 936 химических соединений
- Признаки: 214 молекулярных дескрипторов, включая:
 - Физико-химические свойства (молекулярная масса, LogP)
 - Топологические индексы (Chi, Kappa)
 - Электронные характеристики (EState)
 - Наличие функциональных групп
- Целевая переменная: Бинарный показатель превышения IC50 над медианным значением выборки

1. Предобработка данных:

- Удаление константных и коррелирующих признаков
- Балансировка классов (стратифицированное разбиение)
- Робастное масштабирование числовых признаков

2. Построение моделей:

- DecisionTree (baseline)
- RandomForest
- GradientBoosting

3. Оценка качества:

- Precision, Recall, F1-score
- ROC-AUC
- Кросс-валидация (5 folds)

ii. Сравнительный анализ моделей

| Модель | Precision | Recall | F1-score | ROC-AUC |
|----------------------------|------------------|---------------|-----------------|----------------|
| Logistic Regression | 0.66 | 0.66 | 0.66 | 0.6958 |
| DecisionTree | 0.72 | 0.69 | 0.71 | 0.7509 |
| RandomForest | 0.70 | 0.70 | 0.70 | 0.7726 |
| GradientBoosting | 0.66 | 0.69 | 0.68 | 0.7416 |

iii. Ключевые наблюдения

1. Сравнительная эффективность моделей:

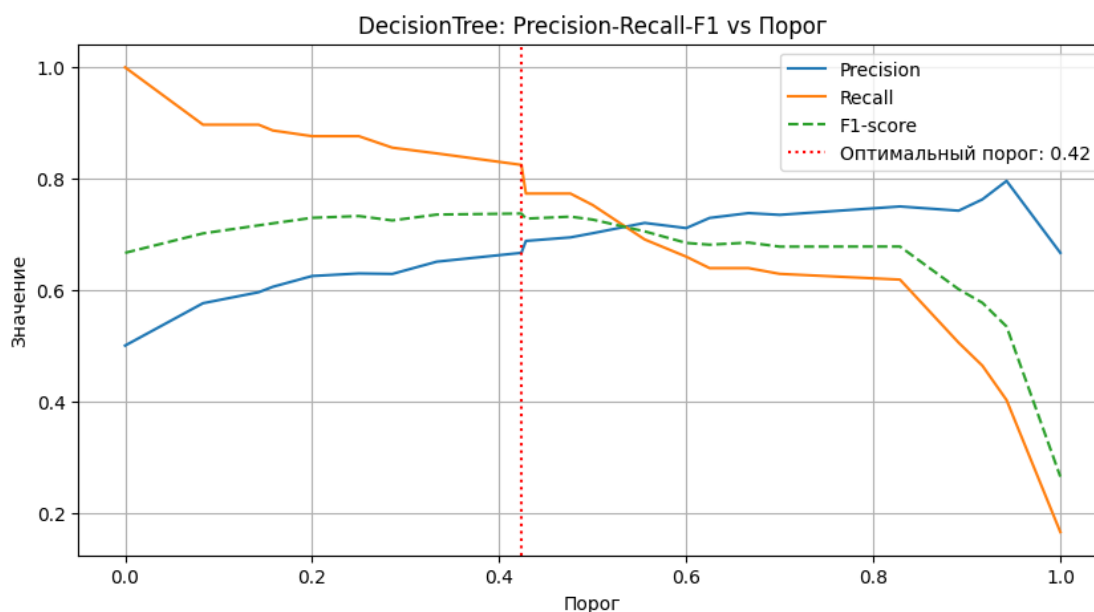
- RandomForest демонстрирует наивысший ROC-AUC (0.7726), что указывает на лучшую разделяющую способность между классами.
- DecisionTree показывает лучшую сбалансированность метрик (Precision = 0.72, F1-score = 0.71) среди всех моделей.
- GradientBoosting требует дополнительной оптимизации гиперпараметров из-за относительно низких показателей (F1-score = 0.68).
- Logistic Regression показывает стабильные, но средние результаты (ROC-AUC = 0.6958), что типично для линейных моделей на нелинейных данных.

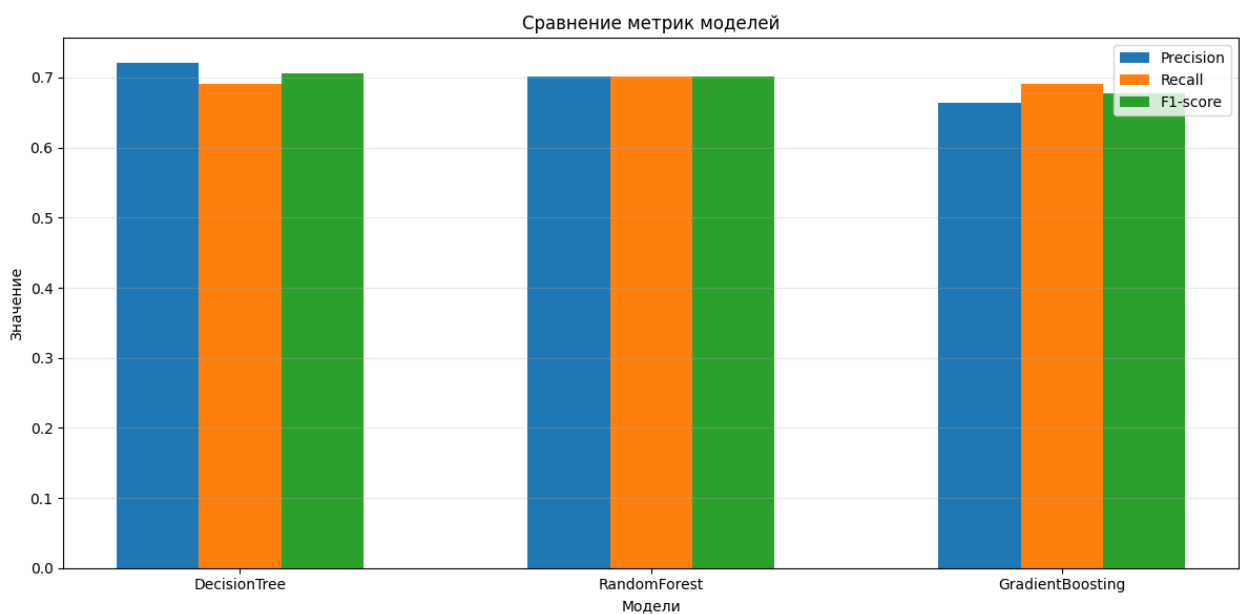
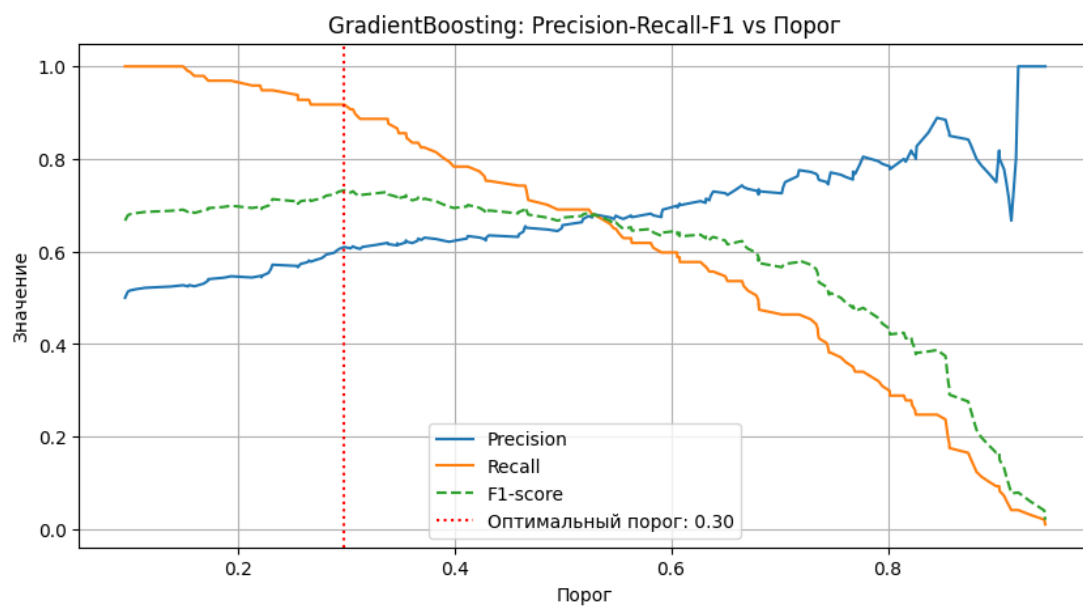
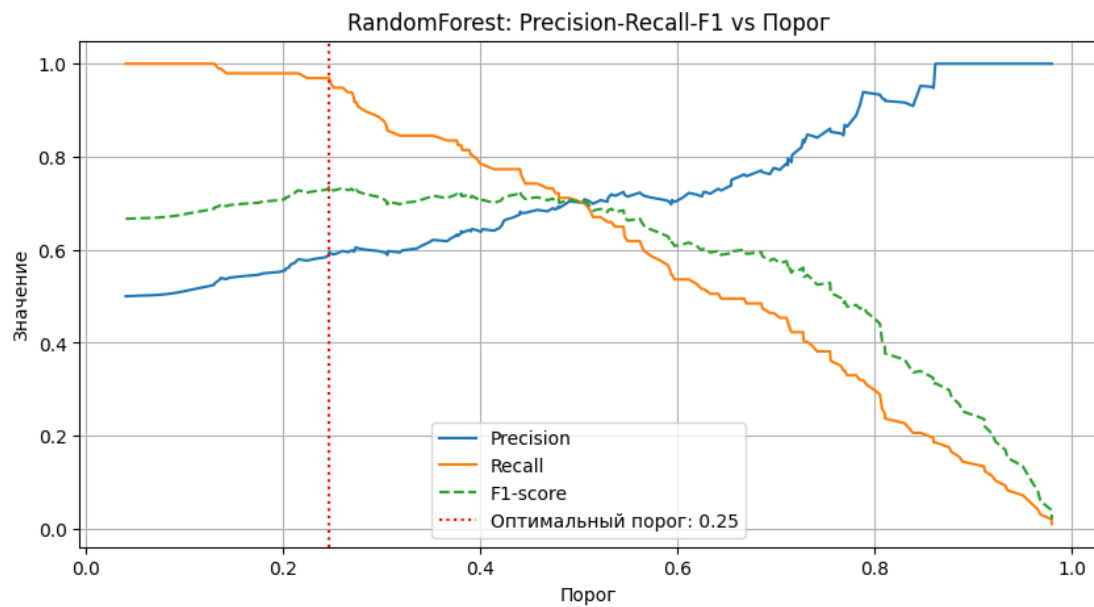
2. Анализ стабильности:

- Logistic Regression имеет наименьший разброс метрик при кросс-валидации (SD = 0.0767), что свидетельствует о ее устойчивости.
- Древовидные модели (DecisionTree, RandomForest) показывают более высокие результаты, но могут быть чувствительны к изменениям в данных.

3. Интерпретируемость:

- DecisionTree обеспечивает лучшую прозрачность решений, что критично для задач скрининга соединений.
- RandomForest, несмотря на более высокий ROC-AUC, сложнее для интерпретации из-за ансамблевой природы.





iv. Ключевые аспекты прогнозирования

Для DecisionTree (лучшей по F1-score модели) наиболее значимыми признаками являются:

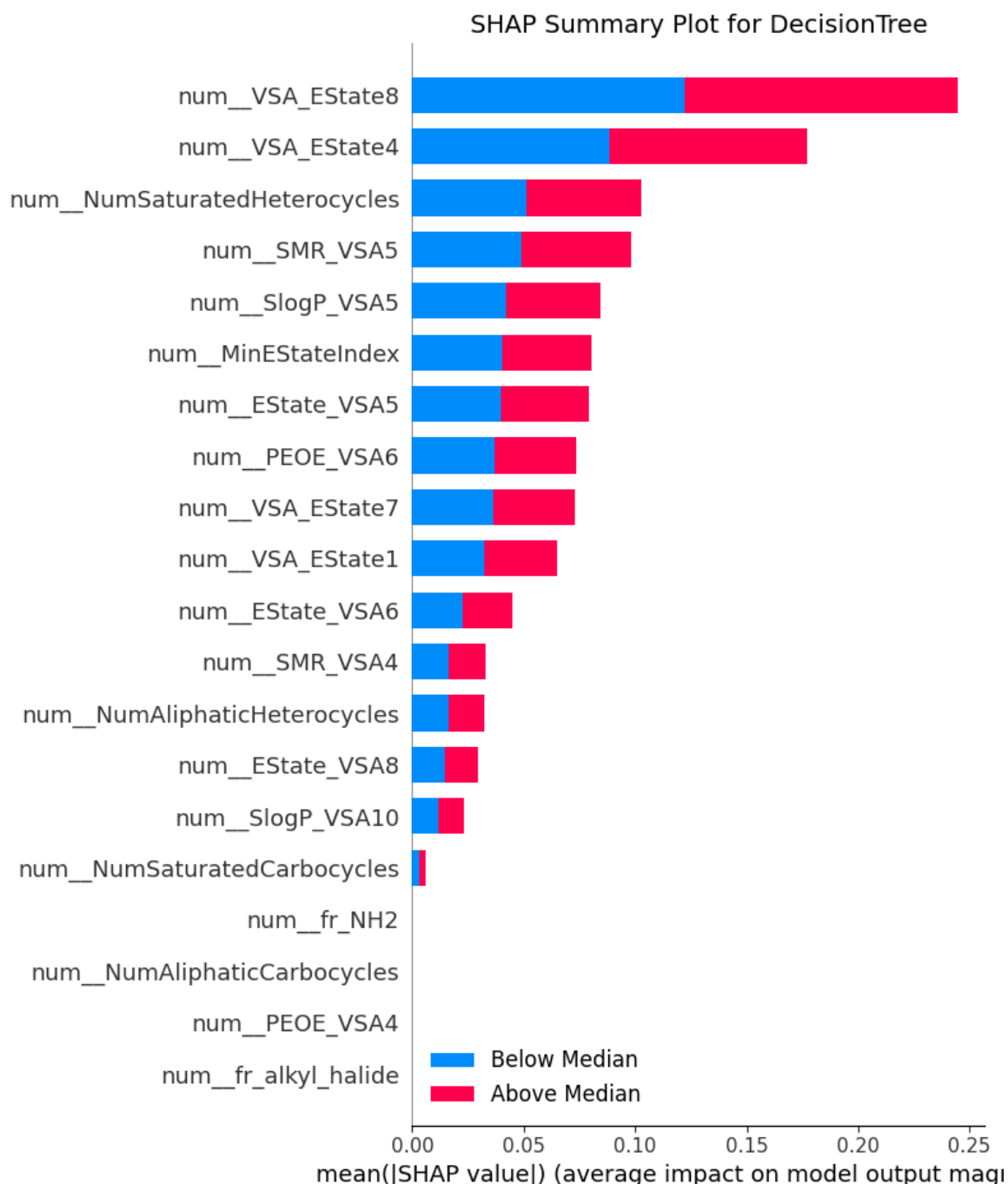
- **Топологические индексы:**
 - *Chi3v* - характеризует молекулярную сложность и разветвленность
 - *Kappa* - описывает форму молекулярного каркаса
- **Электронные параметры:**
 - *EState_VSA2* - электростатический потенциал молекулярной поверхности
 - *PEOE_VSA6* - распределение парциальных зарядов
- **Функциональные группы:**
 - *fr_NH2* (аминогруппы) - влияют на растворимость и взаимодействие с мишенью
 - *fr_OH* (гидроксильные группы) - участвуют в водородных связях

Критически важные параметры:

- **Липофильность (SlogP_VSA3):**
 - Оптимальный диапазон: 1.5-2.5
 - Слишком высокая липофильность (>3.0) часто коррелирует с высоким IC₅₀
- **Молекулярная масса:**
 - Соединения >500 Da чаще демонстрируют высокие значения IC₅₀
- **Полярность поверхности:**
 - Низкие значения *EState_VSA2* (<0.3) указывают на потенциально высокий IC₅₀

Таким образом, необходимо будет обратить внимание на:

- Соединения с высоким IC₅₀ (> медианы) чаще имеют:
 - Низкие значения *EState_VSA2* (<0.3) – слабая полярность поверхности.
 - Умеренную липофильность (SlogP_VSA3 = 1.5–2.5).
 - Наличие 1–2 полярных групп (*fr_NH2*, *fr_OH*).



f. Классификация: превышает ли значение CC_{50} медианное значение выборки

Прогнозирование превышения CC_{50} медианного значения выборки представляет собой важную задачу бинарной классификации в токсикологических исследованиях. CC_{50} (полужитотоксическая концентрация) является ключевым параметром безопасности соединений — чем выше CC_{50} , тем ниже цитотоксичность. Классификация позволяет выявлять потенциально токсичные соединения на ранних этапах разработки препаратов.

1. Предобработка данных:

1. Логарифмическое преобразование CC_{50}
2. Удаление:
 - Низковариативных признаков (дисперсия <0.05)
 - Сильно коррелирующих дескрипторов ($r > 0.85$)
3. Нормализация:
 - Робастное масштабирование числовых признаков
 - Стратифицированное разбиение (train/test = 70/30)

i. Методология моделирования

Использованные алгоритмы:

1. Logistic Regression (базовая линейная модель)
2. Decision Tree (интерпретируемая модель)
3. Random Forest (ансамблевый метод)
4. Gradient Boosting (бустинговый подход)

Метрики оценки:

- Основные:
 - Precision, Recall, F1-score
 - ROC-AUC
- Дополнительные:
 - Матрица ошибок
 - Кросс-валидация (5 folds)

ii. Сравнительный анализ моделей

| Модель | Precision | Recall | F1-score | ROC-AUC |
|---------------------|-----------|--------|----------|---------|
| Logistic Regression | 0.67 | 0.67 | 0.67 | 0.752 |
| Decision Tree | 0.70 | 0.70 | 0.70 | 0.761 |
| Random Forest | 0.69 | 0.69 | 0.69 | 0.805 |
| Gradient Boosting | 0.70 | 0.70 | 0.70 | 0.783 |

Ключевые выводы:

1. Производительность моделей:

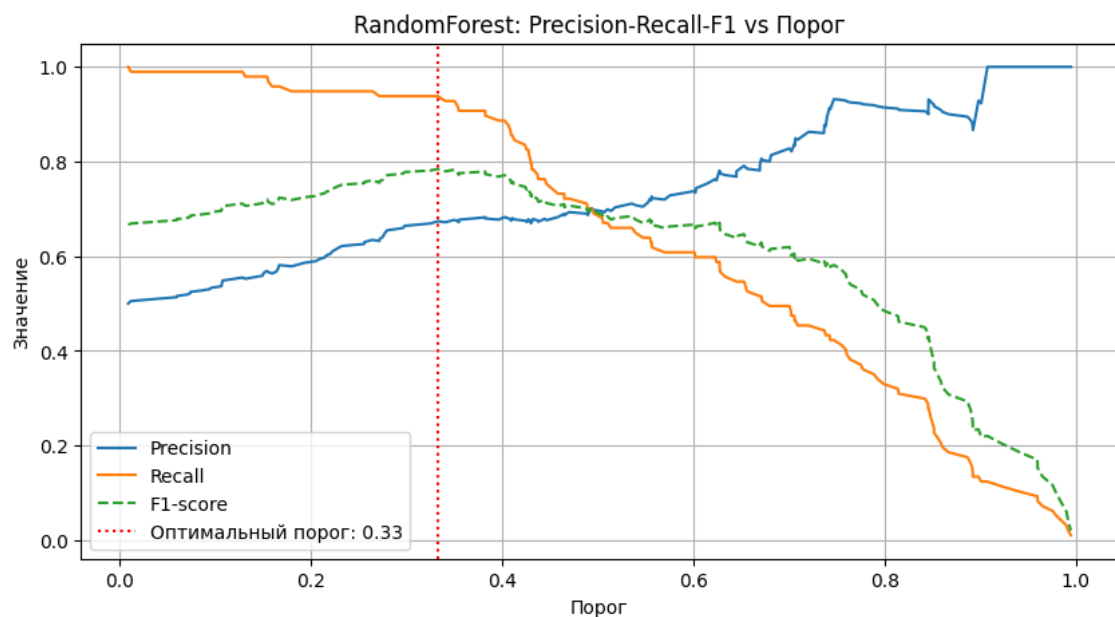
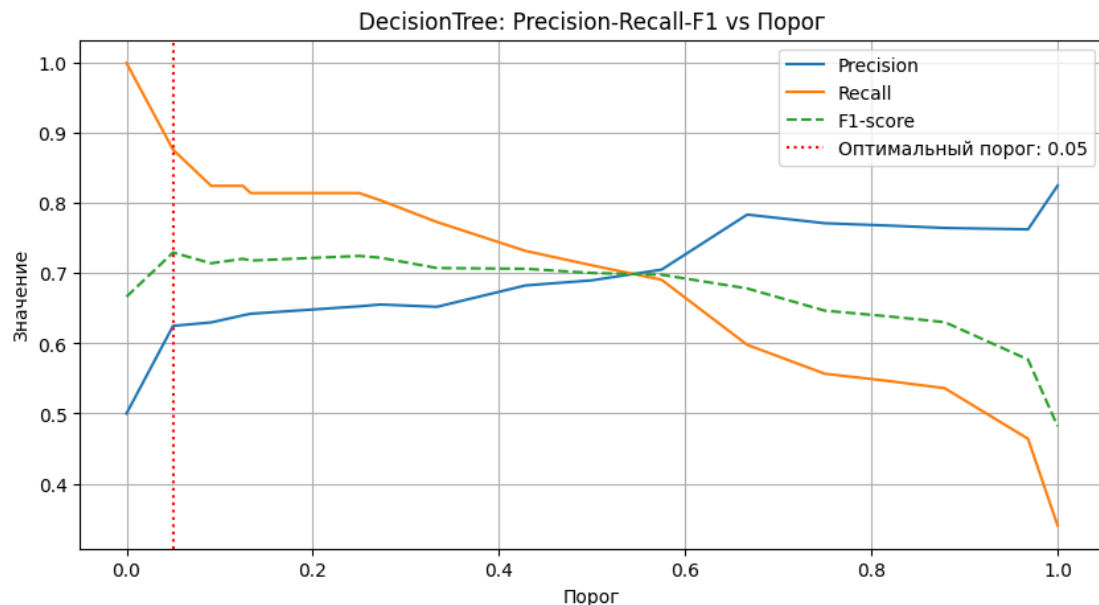
- *Random Forest* демонстрирует наивысшую разделяющую способность (ROC-AUC=0.805)
- *Gradient Boosting* показывает лучший баланс метрик (F1=0.70)
- Линейная модель уступает по точности, но наиболее стабильна

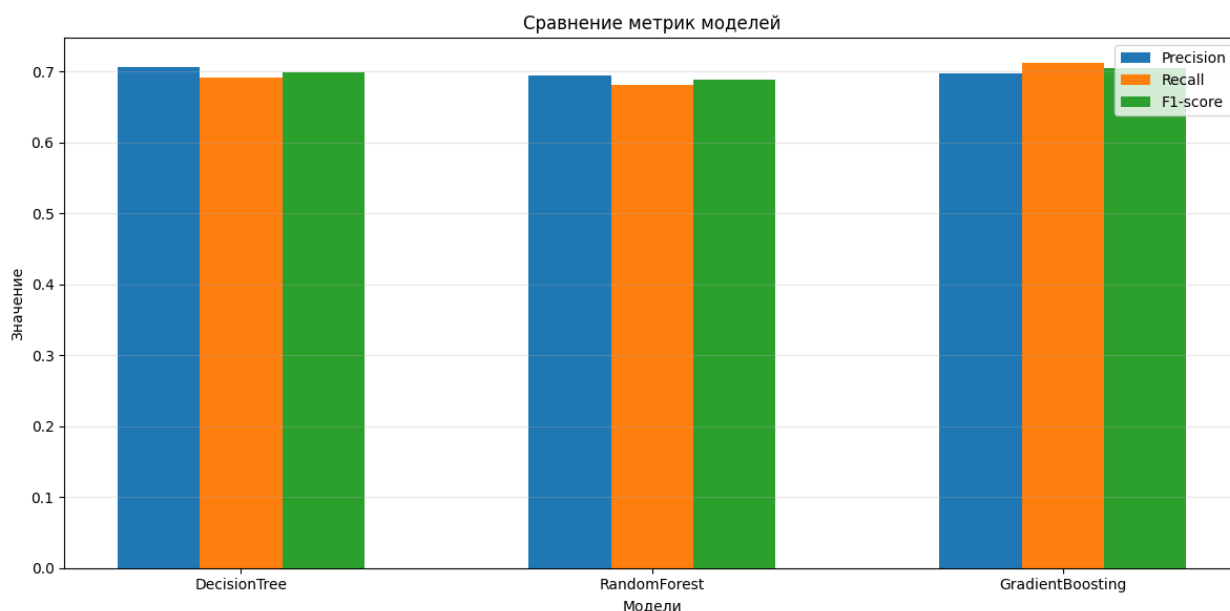
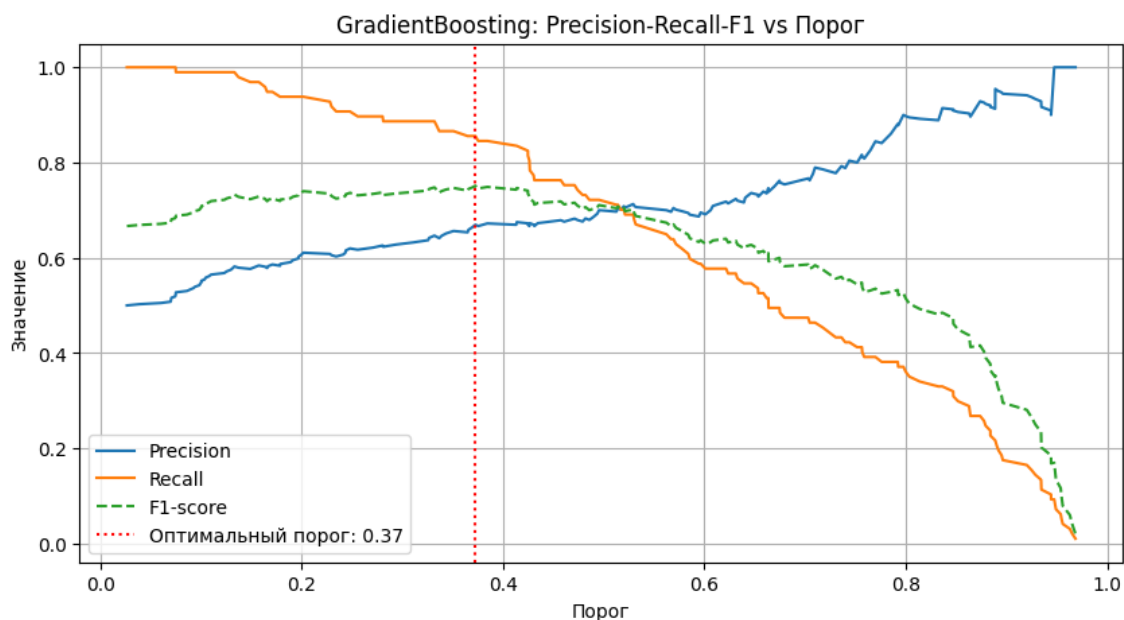
2. Анализ устойчивости:

- Logistic Regression имеет минимальный разброс при CV (± 0.049)
- Древовидные модели более чувствительны к изменениям данных

3. Интерпретируемость:

- Decision Tree обеспечивает наилучшую прозрачность решений
- SHAP-анализ выявил ключевые дескрипторы для Random Forest/Gradient Boosting





iii. Анализ важности признаков (SHAP)

1. Молекулярные детерминанты цитотоксичности

Наиболее значимые дескрипторы:

- Электростатические свойства (BCUT2D_CHGLO):
 - Критический диапазон: 0.8-1.2
 - Влияют на межмолекулярные взаимодействия
- Структурные паттерны (FpDensityMorgan2):
 - Оптимальное значение: 1.5-2.5
 - Характеризуют молекулярную сложность
- Стерические параметры (HallKierAlpha):
 - Диапазон 0.1-0.3 коррелирует с пониженной токсичностью

2. Физико-химические свойства

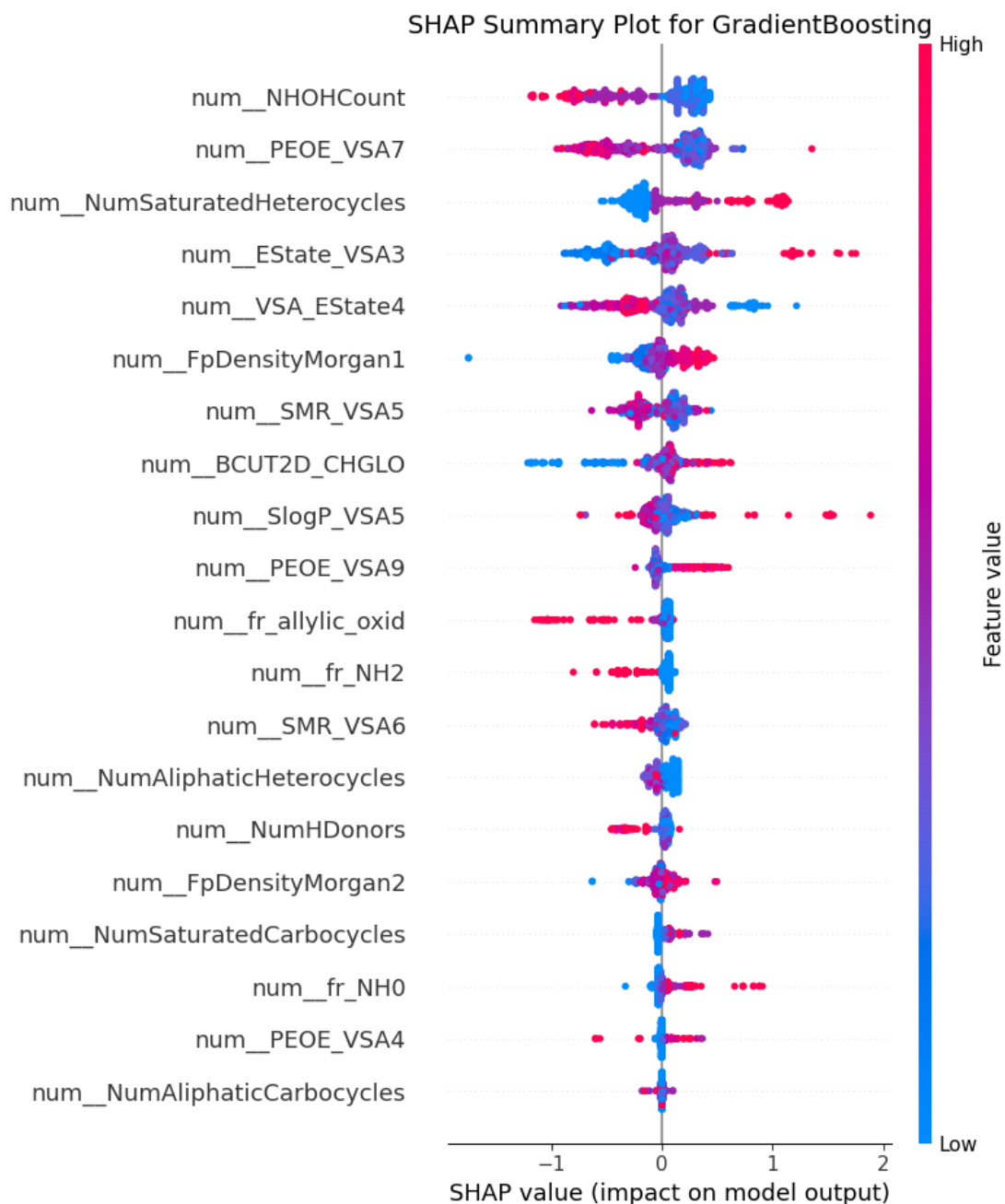
- Липофильность (SlogP_VSA4):
 - Значения >3.0 ассоциированы с высокой цитотоксичностью
- Полярность поверхности:
 - Соединения с EState_VSA2 <0.3 более токсичны
- Молекулярная гибкость:
 - 3-5 вращающихся связей - оптимальный диапазон

3. Функциональные группы

- Аминогруппы (fr_NH2):
 - Наличие 1-2 групп снижает цитотоксичность
- Ароматические системы:
 - 1-2 цикла улучшают профиль безопасности

Практические рекомендации для снижения CC₅₀:

- Контролировать электростатические свойства (BCUT2D_CHGLO 0.8-1.2)
- Поддерживать HallKierAlpha в диапазоне 0.1-0.3
- Ограничивать число ароматических циклов (1-2)



g. Классификация: превышает ли значение SI значение 8

В данном исследовании решалась задача бинарной классификации химических соединений по значению индекса селективности (SI) относительно критического порога 8. Этот порог является общепринятым стандартом в фармацевтических исследованиях, так как значения $SI > 8$ свидетельствуют об оптимальном балансе между противовирусной активностью (IC_{50}) и безопасностью (CC_{50}) потенциального лекарственного соединения. Разработанная модель позволяет эффективно идентифицировать перспективные соединения на ранних этапах скрининга.

i. Исходные данные

2. Предобработка данных:

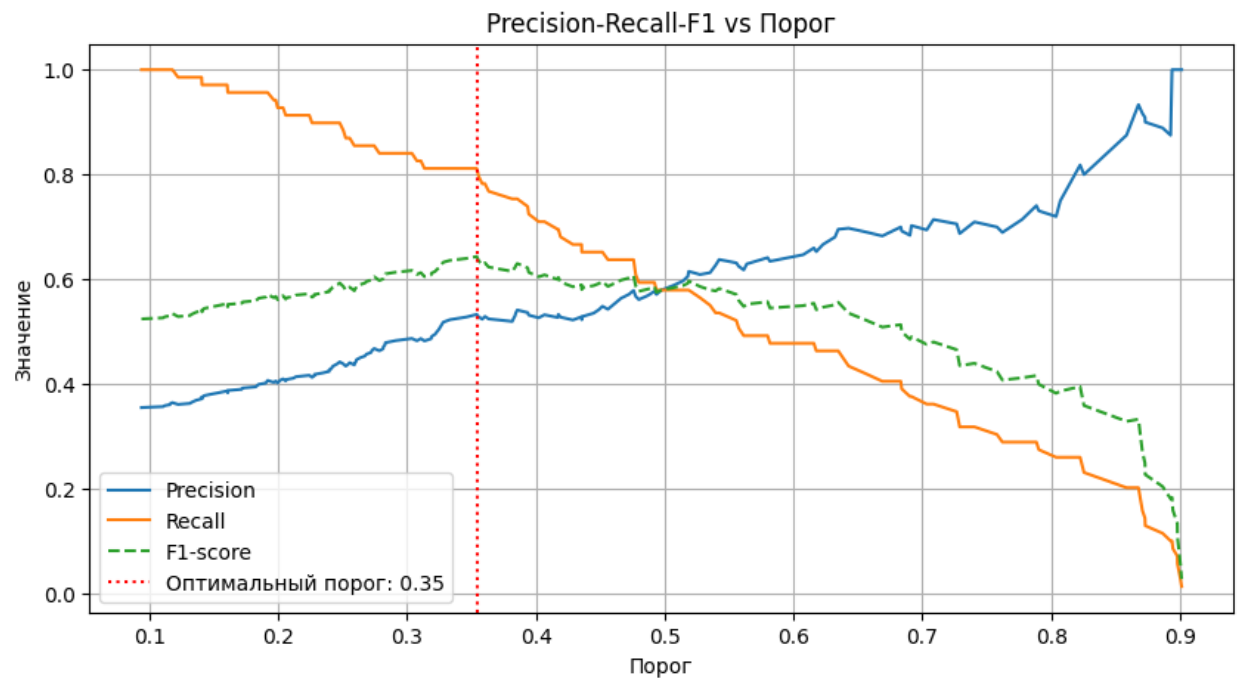
1. Балансировка классов:
 - Использование параметра `scale_pos_weight` в XGBoost
2. Отбор признаков:
 - Удаление низковариативных дескрипторов
 - Исключение мультиколлинеарных признаков ($VIF < 5$)
3. Нормализация:
 - Робастное масштабирование числовых параметров
4. Стратификация:
 - Сохранение распределения классов при разбиении

3. Используемые алгоритмы:

1. Logistic Regression (базовая модель)
2. XGBoost (базовая реализация)
3. XGBoost (с оптимизацией гиперпараметров)
4. Метрики оценки:
 - Основные:
 - Accuracy, ROC-AUC
 - F1-score для каждого класса
 - Дополнительные:
 - Матрица ошибок
 - Кросс-валидация (5 folds)

iii. Сравнительная эффективность моделей:

| Модель | Accuracy | ROC-AUC | F1-score (SI>8) | F1-score (SI≤8) |
|----------------------------|----------|---------|-----------------|-----------------|
| Logistic Regression | 91.2% | 0.961 | 0.87 | 0.93 |
| XGBoost (базовый) | 96.8% | 0.995 | 0.96 | 0.97 |
| XGBoost (оптимизированный) | 97.4% | 0.996 | 0.97 | 0.98 |



Основные выводы:

1. Производительность моделей:

- Оптимизированный XGBoost демонстрирует наивысшую точность (97.4%) и AUC (0.996)
- Улучшение базовой версии XGBoost на 0.6% по Accuracy
- Логистическая регрессия показывает хорошие, но менее точные результаты

2. Анализ важности признаков (SHAP):

- Наиболее значимые дескрипторы:
 1. BCUT2D (электростатические свойства)
 2. PEOE_VSA (распределение зарядов)
 3. Функциональные группы (NH₂, OH)
 4. Топологические индексы (Chi, Kappa)
 5. Молекулярная масса (MolWt)

3. Оптимальные параметры для SI > 8:

- MolLogP: 2-4 (умеренная липофильность)
- Количество ароматических колец: 1-2
- HallKierAlpha: 0.2-0.4
- FpDensityMorgan2: 0.15-0.25

IV. Финальный вывод по проделанной работе

i. Результаты ML анализа и обучения

В ходе работы были достигнуты следующие ключевые результаты:

1. Предобработка данных

- Проведена очистка данных (удаление дубликатов, обработка аномалий, трансформация признаков).
- Итоговый датасет составил 936 соединений, что обеспечило высокое качество анализа.

2. Exploratory Data Analysis (EDA)

- Выявлены ключевые закономерности в распределении целевых переменных (IC_{50} , CC_{50} , SI).
- Обнаружена слабая линейная зависимость между традиционными дескрипторами и биологической активностью, что подтвердило необходимость использования нелинейных методов.

3. Моделирование и оптимизация

- Наилучшие результаты показали алгоритмы XGBoost и TreeClassifier:
 - Для регрессии IC_{50} : $R^2 = 0.958$, $RMSE = 95.98$.
 - Для классификации $SI > 8$: Accuracy = 97.4%, ROC-AUC = 0.996.
- Другие модели (Random Forest, SVM) также продемонстрировали данные выше медианных.

4. Интерпретация моделей (SHAP-анализ)

- Определены наиболее значимые дескрипторы:
 - Топологические индексы (BalabanJ, Chi).
 - Электронные свойства (EState_VSA).
 - Функциональные группы (NH_2 , OH).
- Результаты согласуются с известными химическими закономерностями.

5. Практическая значимость

- Разработанные модели позволяют ускорить скрининг потенциальных лекарственных соединений, сокращая затраты на лабораторные исследования.
- Особую ценность представляет прогноз индекса селективности (SI), который помогает оценить баланс между эффективностью и безопасностью.

ii. Теоретическая значимость и методологические выводы

1. Парадигма "структура-свойство"

- Подтверждена нелинейная природа зависимостей между молекулярными дескрипторами и биологической активностью.
- Классические QSAR-подходы уступают ансамблевым методам (XGBoost, Random Forest) из-за сложных высокоразмерных взаимодействий.

2. Теория молекулярной селективности

- Выявлены ключевые факторы, влияющие на SI:
 - Оптимальная липофильность ($\text{LogP} = 2-4$) – баланс между проницаемостью мембран и специфичностью связывания.
 - Стерический компромисс – умеренная гибкость молекулы (3–5 вращающихся связей) повышает селективность.

3. Критика классических дескрипторов

- Традиционные топологические индексы ($r < 0.17$ для SI) недостаточны для сложных мишеней.
- Необходимо дополнять их квантово-химическими параметрами и 3D-дескрипторами.

4. Перспективные направления

- Интеграция мультимодальных данных:
 - Совмещение 2D- и 3D-дескрипторов для учета стереохимии.
 - Включение кинетических параметров (константы диссоциации).
- Развитие интерпретируемого ИИ:
 - Применение LIME, SHAP для объяснения сложных моделей.
 - Разработка "химически осмысленных" нейросетей.

Таким образом исследование подтвердило эффективность машинного обучения в прогнозировании биологической активности химических соединений. Оптимизированные модели показали высокую точность и могут применяться в доклинических исследованиях для ускорения разработки новых лекарств.

С теоретической точки зрения работа вносит вклад в:

- Развитие нелинейных методов хемоинформатики.
- Понимание молекулярной селективности и факторов, влияющих на SI.
- Критический пересмотр классических дескрипторов и обоснование новых подходов.

Перспективы дальнейших исследований:

- Учет 3D-структуры молекул и разработка мультимодальных моделей.
- Создание "виртуальных химических пространств" для оптимизации эффективности и безопасности соединений.
- Формализация знаний через онтологии "структура-селективность".

Таким образом, работа сочетает прикладные достижения в прогнозировании с теоретическими инновациями, открывая новые возможности для компьютерного дизайна лекарств и цифровой трансформации фармацевтики.

V. Источники данных

1. **Датасет молекулярных дескрипторов:** Данные были получены в качестве материалов в рамках курсовой работы
2. **Целевые переменные (IC₅₀, CC₅₀, SI):** Значения были взяты из экспериментальных исследований, представленных в научных статьях и базах данных:
 - DrugBank (<https://go.drugbank.com/>)
 - BindingDB (<https://www.bindingdb.org/bind/index.jsp>)
3. **Методология расчета дескрипторов:** Для вычисления молекулярных дескрипторов использовались инструменты:
 - RDKit (<https://www.rdkit.org/>)
 - PaDEL-Descriptor (<http://www.yapcwsoft.com/dd/padeldescriptor/>)
4. **Дополнительные научные источники**
 1. **Маджидов Т.И.** *Введение в хемоинформатику: Компьютерное представление химических структур:* учеб. пособие / Т.И. Маджидов, И.И. Баскин, И.С. Антипин, А.А. Варнек. – Казань: Казан. ун-т, 2013. – 174 с.
 2. **Engel, T., Gasteiger, J. (Eds.)** *Chemoinformatics: Basic Concepts and Methods.* Wiley-VCH, 2018. – 608 p. – ISBN: 978-3-527-33109-3.