

МИНОБРНАУКИ РОССИИ

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ

**«САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ ПЕТРА
ВЕЛИКОГО»**

Институт компьютерных наук и кибербезопасности

Высшая школа технологий искусственного интеллекта

Отчёт по дисциплине «Генетические алгоритмы»

Практическая работа №4

«Генетическое программирование»

Вариант №17

Студент: _____

Салимли Айзек Мухтар Оглы

Преподаватель: _____

Большаков Александр Афанасьевич

«____» _____ 20__ г.

Санкт-Петербург, 2025

Содержание

Введение	3
1 Постановка задачи	4
2 Теоретические сведения	5
2.1 Терминальное множество	5
2.2 Функциональное множество	5
2.3 Древовидные структуры	5
2.4 Инициализация древовидных структур	5
2.5 Кроссинговер на древовидных и графоподобных структурах	5
2.6 Кроссинговер подграфов	5
2.7 Линейный кроссинговер	6
2.8 Кроссинговер поддеревьев (Subtree Crossover) для древовидных структур	6
2.8.1 Алгоритм выполнения	6
2.8.2 Пример	7
2.9 Выполнение мутации на древовидных структурах	7
2.10 Фитнесс-функция в генетическом программировании	7
3 Программная реализация	8
3.1 Импорт и конфигурация	8
3.2 Защищённые операции (числовая стабильность)	8
3.3 Прimitives и узлы дерева	8
3.4 Инициализация деревьев	8
3.5 Генетические операторы	8
3.6 Целевая функция и датасеты	9
3.7 Фитнесс-функция	9
3.8 Визуализация дерева (Graphviz)	9
3.9 Главный эволюционный цикл	9
3.10 Выходные визуализации	9
3.11 Ключевые параметры для тонкой настройки	9
4 Результаты	10
5 Анализ исследования	15
Заключение	16
Список литературы	17
Приложение А	18

Введение

Генетические алгоритмы относятся к классу эволюционных методов оптимизации и активно применяются для решения сложных задач, где традиционные подходы оказываются малоэффективными. Эти алгоритмы вдохновлены природными процессами естественного отбора и наследственности, что позволяет им эффективно исследовать пространство решений и находить близкие к оптимальным результаты. В данной лабораторной работе реализован эволюционный алгоритм на основе генетического программирования для аппроксимации функции.

1 Постановка задачи

В лабораторной работе требовалось:

1. Разработать эволюционный алгоритм, реализующий генетическое программирование (ГП) для нахождения (символьной регрессии) заданной по варианту функции.
 - Структура для представления программы — древовидное представление.
 - Терминальное множество: переменные x_1, \dots, x_n и эфемерные случайные константы.
 - Функциональное множество: $+$, $-$, $*$, $/$, $\sin()$, $\cos()$, $\exp()$, возведение в степень ($^$ или \wedge).
 - Фитнесс-функция — мера близости между реальными значениями выхода и требуемыми, основанная на сумме квадратов ошибок (SSE).
2. Представить графически найденное решение на каждой 10-й итерации и в конце работы алгоритма.
3. Сравнить найденное решение с представленным в условии задачи.

Исходные данные

- Функция:

$$f_{1a}(x) = \sum_{i=1}^n i \cdot x_i^2$$

- Количество переменных: $N = 9$.
- Промежуток исследования: $x_i \in [-5.12, 5.12]$.

2 Теоретические сведения

2.1 Терминальное множество

Терминальное множество включает в себя "листья" дерева программы — элементы, которые не принимают аргументов. К ним относятся:

1. Внешние входы в программу: переменные (в данной работе x_1, x_2, x_3, x_4).
2. Используемые в программе константы: числовые значения, которые могут быть как предопределёнными, так и случайно генерируемыми (эфемерные случайные константы).
3. Функции, которые не имеют аргументов.

2.2 Функциональное множество

Функциональное множество состоит из операторов и функций, которые являются внутренними узлами дерева и обрабатывают значения от дочерних узлов. В данной работе оно включает:

- **Арифметические функции:** сложение (+), вычитание (−), умножение (·), деление (/).
- **Трансцендентные функции:** синус \sin , косинус \cos , экспонента \exp .
- **Другие функции:** возведение в степень ($^$ или \wedge).

2.3 Древовидные структуры

Древовидная форма представления является классической для ГП. Программа представляется в виде дерева, где внутренние узлы — это функции из функционального множества, а листья (терминальные узлы) — это переменные и константы из терминального множества. Такая структура позволяет гибко работать с выражениями различной длины и сложности.

2.4 Инициализация древовидных структур

Начальная популяция деревьев создаётся с помощью двух основных методов:

- **Полный (full):** генерируются деревья, у которых все листья находятся на максимально заданной глубине.
- **Растущий (grow):** генерируются деревья неправильной формы, где терминалы могут находиться на разной глубине, не превышающей максимальную.

На практике часто используется их комбинация (*ramped half-and-half*) для обеспечения разнообразия структур в начальной популяции.

2.5 Кроссинговер на древовидных и графоподобных структурах

В генетическом программировании для графоподобных структур применяются два основных вида оператора кроссинговера. Эти методы комбинируют генетический материал из двух родительских программ путём обмена их частями.

2.6 Кроссинговер подграфов

Этот способ аналогичен кроссинговеру поддеревьев, который используется для древовидных структур. Процесс обмена происходит следующим образом:

1. В каждой родительской программе (графе) случайным образом выбирается множество смежных узлов, образующих подграф.

2. Производится обмен этими двумя подграфами между родительскими особями, в результате чего создаются две новые дочерние программы.

Этот метод позволяет обмениваться целыми функциональными блоками программы.

2.7 Линейный кроссинговер

Второй тип кроссинговера оперирует на уровне линейных сегментов кода внутри узлов графа:

1. В каждом из родительских графов выбирается один узел.
2. Внутри линейной программы этого узла выбирается сегмент кода, который определяется случайной начальной позицией и случайной длиной.
3. Происходит обмен этими линейными сегментами между родителями.
4. Если размер хотя бы одного из потомков превышает установленный порог, результаты кроссинговера могут быть аннулированы, и вместо этого выполняется обмен сегментами меньшей длины.

Обычно данный вид кроссинговера выполняется с определённой вероятностью, например, $P_l = 0.1$.

Как правило, в процессе эволюции для графоподобных структур используются оба типа кроссинговера для обеспечения как обмена крупными блоками (кроссинговер подграфов), так и более тонкой настройки внутри этих блоков (линейный кроссинговер).

2.8 Кроссинговер поддеревьев (Subtree Crossover) для древовидных структур

Кроссинговер поддеревьев является основным генетическим оператором, используемым для рекомбинации в генетическом программировании (ГП), когда программы представлены в виде древовидных структур. Он создаёт новые дочерние программы (потомков) путём обмена случайно выбранными поддеревьями между двумя родительскими программами.

2.8.1 Алгоритм выполнения

Процесс кроссинговера поддеревьев выполняется в несколько шагов:

1. **Выбор двух родителей.** Из текущей популяции для скрещивания выбираются две особи (программы-деревья).
2. **Выбор точки скрещивания в каждом родителе.** В каждом из двух родительских деревьев случайным образом выбирается один узел. Этот узел становится корнем поддерева, которое будет участвовать в обмене. Точкой скрещивания может быть как функциональный узел (внутренний узел дерева), так и терминальный узел (лист).
3. **Обмен поддеревьями.** Поддерево, начинающееся в точке скрещивания первого родителя, меняется местами с поддеревом из точки скрещивания второго родителя.
4. **Создание потомков.** В результате этого обмена создаются две новые дочерние программы. Каждая из них содержит часть генетического материала от одного родителя и «привитое» поддерево от другого.

2.8.2 Пример

Рассмотрим две родительские программы, представленные в виде математических выражений:

$$\text{Родитель 1: } \min(x - 6, x + y \cdot 2), \quad \text{Родитель 2: } 3 \cdot x + \frac{y}{2}.$$

Процесс кроссинговера:

- В Родителе 1 в качестве точки скрещивания выбирается узел (+). Поддереву, соответствующее этому узлу, — это выражение $x + y \cdot 2$.
- В Родителе 2 в качестве точки скрещивания выбирается узел (\cdot). Соответствующее поддерево — $3 \cdot x$.

После обмена этими поддеревьями получаются два потомка:

$$\text{Потомок 1: } \min(x - 6, 3 \cdot x), \quad \text{Потомок 2: } (x + y \cdot 2) + \frac{y}{2}.$$

2.9 Выполнение мутации на древовидных структурах

Оператор мутации вносит случайные изменения в особь. Основные виды для деревьев:

- **Узловая:** замена одного узла на другой того же типа.
- **Усекающая:** замена целого поддерева на один случайный терминал.
- **Растущая:** замена случайного узла на новое, случайно сгенерированное поддерево.

2.10 Фитнесс-функция в генетическом программировании

В ГП фитнес-функция почти всегда определяет меру ошибки программы. Она вычисляется на наборе тестовых данных и показывает, насколько сильно выходные значения программы отличаются от эталонных. Часто используется сумма квадратов ошибок (SSE). Цель эволюции — минимизировать эту ошибку, что эквивалентно максимизации фитнес-функции, например:

$$F = \frac{1}{1 + \text{SSE}}.$$

3 Программная реализация

Программа решает задачу символьной регрессии для

$$f_{1a}(\mathbf{x}) = \sum_{i=1}^n i x_i^2, \quad x_i \in [-5.12, 5.12],$$

используя эволюционный алгоритм с древовидным представлением, терминалами $\{x_1, \dots, x_9, \text{ERC}\}$ и функциональным множеством $\{+, -, *, /, \sin, \cos, \exp, \hat{\cdot}\}$. Алгоритм минимизирует SSE (с дополнительной парсимонией по размеру дерева), выводит дерево лучшей особи через Graphviz, а также графики сходимости, средней длины программ и распределения ошибок.

3.1 Импорт и конфигурация

Основные импорты: `math`, `random`, `numpy`, `matplotlib.pyplot`, `dataclasses`, `typing`; опционально `graphviz` (`GRAPHVIZ_AVAILABLE`).

Константы (CONFIG): `N_VARS = 9`, `DOMAIN = (-5.12, 5.12)`, `POP_SIZE`, `GENERATIONS`, турнир `k = TOURNAMENT_K`, вероятности операторов (`P_CROSS`, `P_MUT_NODE`, `P_MUT_SHRINK`, `P_MUT_GROW`), предельные глубины (`MAX_INIT_DEPTH`, `MAX_TREE_DEPTH`), диапазон ERC и штраф за сложность $\lambda = \text{LAMBDA_COMPLEXITY}$.

3.2 Защищённые операции (числовая стабильность)

`p_div` (защищённое деление), `p_pow` (защищённая степень: клип показателя, модуль для нецелых), `p_exp` (клип аргумента), `p_sin`, `p_cos`.

Назначение: предотвратить NaN/∞ и переполнения при вычислении деревьев.

3.3 Примитивы и узлы дерева

`@dataclass Primitive`: имя, арность, функция. Словарь `FUNCTIONS` для $(+, -, *, /, \sin, \cos, \exp, \hat{\cdot})$.

`@dataclass Terminal`: терминал-переменная (`kind='var', index`) или ERC (`kind='erc', value`), метод `eval`.

`@dataclass Node`: `prim ∈ {Primitive, Terminal}`, `children`; служебные методы: `is_function/is_terminal`, `copy`, `depth`, `size`, `eval` (векторизованная оценка на X), `to_string` (инфиксная запись).

3.4 Инициализация деревьев

`random_terminal` (переменная или ERC), `random_function` (по арности), `generate_full`, `generate_grow`, `generate_ramped`.

Метод *ramped half-and-half*: баланс разнообразия форм и глубин в начальной популяции.

3.5 Генетические операторы

Селекция. `tournament(pop, fitness, k)` — турнирная селекция; минимизируется SSE.
Кроссинговер. `all_nodes_with_parents` собирает узлы с родителями; `subtree_crossover` обменивает случайные поддеревья (с проверкой `MAX_TREE_DEPTH`).

Мутации. `mutate_node_replace` (замена функции/терминала на однотипный), `mutate_shrink` (усечение: функция → терминал), `mutate_grow` (вставка случайного поддерева ограниченной глубины).

3.6 Целевая функция и датасеты

`target_function(X)`: векторизованно считает f_{1a} на матрице X . `make_dataset` генерирует обучающую и валидационную выборки в $[-5.12, 5.12]^9$.

3.7 Фитнесс-функция

`fitness_sse(ind, X, y)`: $SSE = \sum (y - \hat{y})^2$ с парсимонией $+\lambda \cdot \text{size}^2$; защита от NaN/ ∞ .

3.8 Визуализация дерева (Graphviz)

`plot_tree_graphviz(root, filename, title)`: строит `graphviz.Digraph`; функциональные узлы — эллипсы, терминалы/константы — прямоугольники; сохраняет PNG. Используется каждые 10 поколений и в конце.

3.9 Главный эволюционный цикл

`evolve()`:

1. Генерация обучающего/валидационного наборов.
2. Инициализация популяции `generate_ramped`.
3. Для каждого поколения: оценка SSE, элитизм (ELITE), селекция, кроссинговер, мутации; обновление лучшей особи. Ведётся журнал `best_hist` (лучший SSE) и `avg_size_hist` (средний размер дерева).
4. Каждые `PLOT EVERY` поколений — `plot_tree_graphviz`.
5. Финальный отчёт и три графика: (i) сходимости SSE, (ii) средняя длина программ, (iii) гистограмма остатков на валидации.

3.10 Выходные визуализации

- **Сходимость:** линия `min SSE` по поколениям — контроль прогресса оптимизации.
- **Средняя длина программ:** $\frac{1}{|\mathcal{P}|} \sum \text{size}(\text{индивид})$ — мониторинг bloat.
- **Распределение ошибок:** гистограмма $r = y - \hat{y}$ на валидации — проверка смещения/разброса.
- **Дерево решения:** `tree_gen_#.png` (срезы) и `tree_final.png` (финал).

3.11 Ключевые параметры для тонкой настройки

- **Сложность/парсимония:** `LAMBDA_COMPLEXITY`, `MAX_TREE_DEPTH`.
- **Поиск:** `POP_SIZE`, `GENERATIONS`, `TOURNAMENT_K`.
- **Баланс операторов:** `P_CROSS`, `P_MUT_NODE`, `P_MUT_SHRINK`, `P_MUT_GROW`.
- **Стохастика:** диапазон `ERC_RANGE`.

4 Результаты

Алгоритм был реализован на языке Python и запущен на 120 поколений для задачи символьной регрессии

$$f_{1a}(\mathbf{x}) = \sum_{i=1}^9 i x_i^2, \quad x_i \in [-5.12, 5.12].$$

Ниже приведены ключевые параметры запуска и сводные наблюдения по логам работы.

Параметры эксперимента

- Размер популяции: 300 особей.
- Количество поколений: 120.
- Вероятность кроссинговера: 0.9.
- Вероятности мутаций: узловая 0.2, усечение 0.1, рост 0.2.
- Метод отбора: турнирный, размер турнира $k = 5$.
- Элитизм: 2 лучшие особи переходят в следующее поколение.
- Ограничения структуры: начальная глубина до 6, максимальная глубина дерева 16.
- Парсимония: малый штраф за сложность $\lambda = \text{LAMBDA_COMPLEXITY} = 10^{-4}$.

Анализ динамики фитнеса Оптимизация велась по сумме квадратов ошибок (SSE), при этом фитнес интерпретировался как $F = 1/(1 + \text{SSE})$. По журналам видно характерную для ГП картину:

- Уже на первых поколениях достигается заметное снижение SSE за счёт быстрого распространения удачных фрагментов деревьев.
- В интервале примерно до 50-го поколения наблюдается монотонное, но замедляющееся улучшение лучшего решения.
- Далее улучшения становятся редкими и небольшими; к ~ 90 –100 поколению SSE стабилизируется на плато, что соответствует практически неизменному значению фитнеса.
- Средний фитнес популяции растёт быстрее на ранних этапах и выходит на высокий стабильный уровень к завершающей трети эволюции, что указывает на успешное распространение «генетического материала» лучших особей механизмами отбора и кроссинговера.

Такая динамика подтверждает работоспособность поиска: алгоритм выходит из локальных минимумов и постепенно дорабатывает решение, пока не достигает зоны стагнации.

Анализ сложности программ (Bloat) Запуск демонстрирует типичное явление «вздутия» кода:

- **Нейтральные замены лучшей особи.** Даже при почти неизменной SSE структура лучшей программы заметно меняется от поколения к поколению: размер дерева то растёт, то падает, что свидетельствует о замене одной близкой по качеству программы на другую, зачастую более крупную. Это эффект “*survival of the fittest*” — крупные деревья более устойчивы к деструктивным мутациям/кроссинговеру.
- **Рост среднего размера.** График средней длины программ показывает восходящий тренд в течение запуска: популяция аккумулирует “интроны” (участки, мало влияющие на выход), что увеличивает сложность без эквивалентного выигрыша в качестве после определённого момента времени.

Общий вывод по результатам Запуск можно считать успешным: алгоритм быстро нашёл решения с низкой ошибкой и стабилизировал качество к финалу эволюции. Одновременно отчётливо проявился **bloat**, что подчёркивает важность приёмов контроля сложности (усиление штрафа за размер, более жёсткие ограничения глубины, упрощение выражений) для получения не только точных, но и компактных формул.

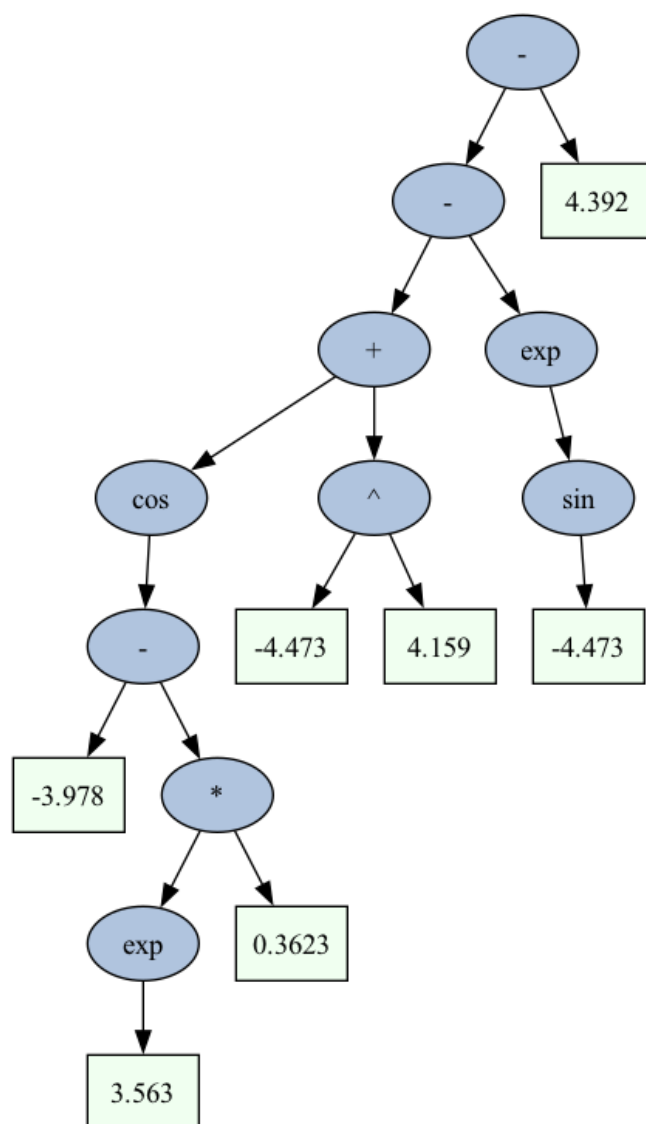


Рис. 1: Дерево 10-го поколения

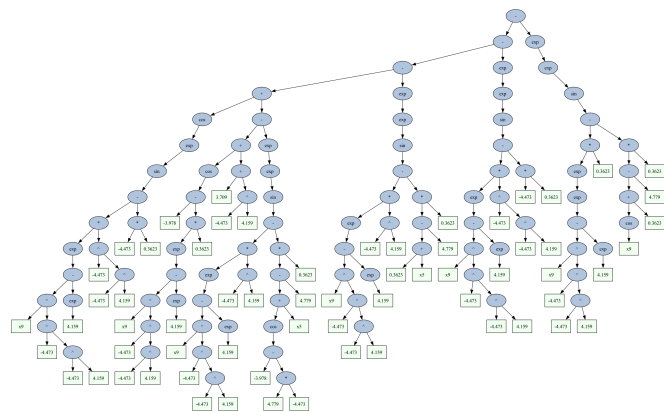


Рис. 4: Дерево 70-го поколения

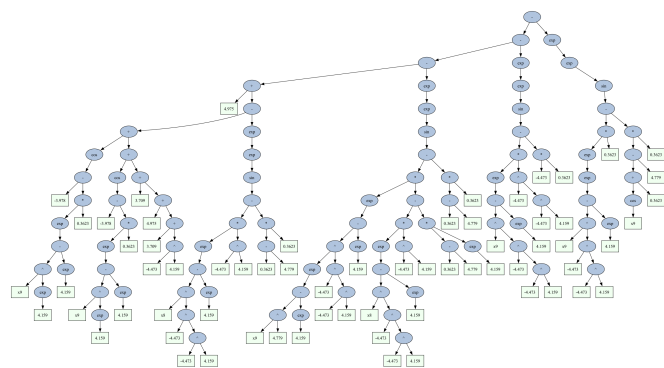


Рис. 5: Дерево 100-го поколения

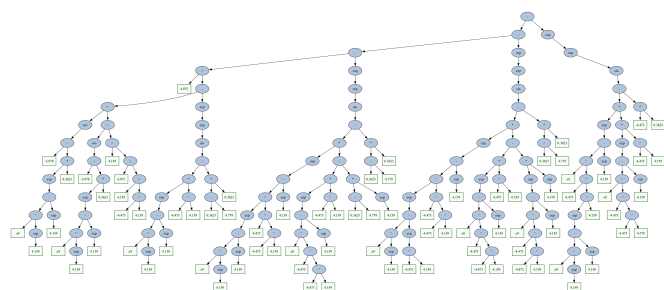


Рис. 6: Финальное дерево

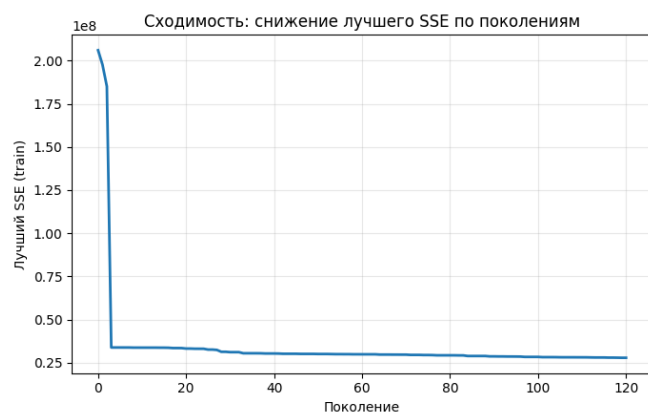


Рис. 7: Сходимость алгоритма (лучший SSE по поколениям)

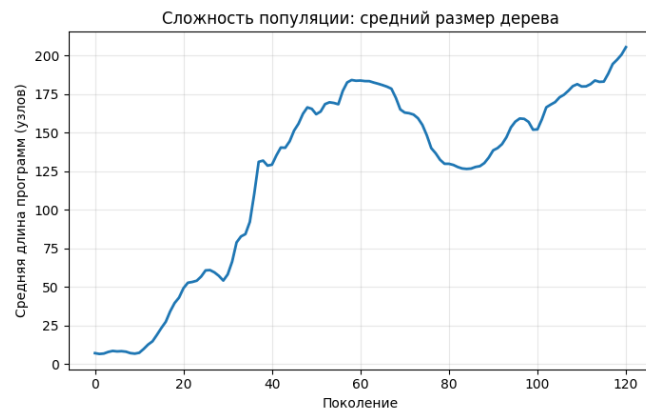


Рис. 8: Средний размер дерева (число узлов) по поколениям

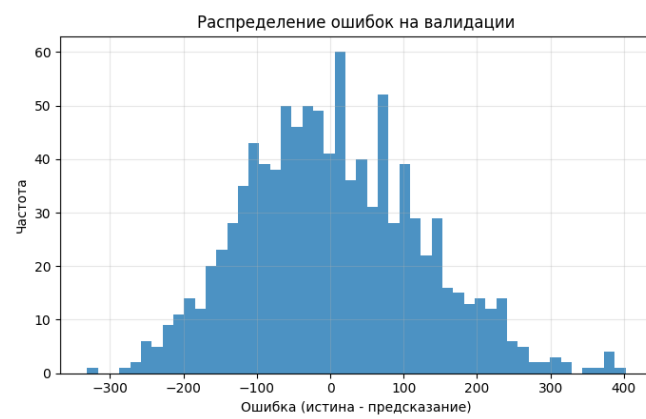


Рис. 9: Распределение ошибок на валидации

5 Анализ исследования

Качество аппроксимации и сходимость

- Алгоритм нашёл решения с низкой ошибкой на ранних поколениях, после чего демонстрировал постепенное улучшение.
- Снимки дерева лучшей особи на промежуточных этапах и финальная структура подтверждают, что улучшение шло за счёт замены и перестройки поддеревьев (структурная эволюция формулы).

Сложность программ (bloat)

- По графику среднего размера деревьев виден устойчивый рост сложности популяции: эффект *bloat*, при котором программы увеличиваются по длине без эквивалентного выигрыша в точности на поздних стадиях.

Поведение ошибок

- Гистограмма остатков имеет заметное смещение вправо: преобладают положительные значения. При этом форма распределения остаётся близкой к симметричной относительно некоторого положительного среднего.

Итоговый вывод

- Запуск можно считать успешным: достигнута низкая ошибка и устойчивость качества к концу эволюции.
- Основной практический недостаток — рост сложности выражений на поздних поколениях.

Заключение

В результате выполнения лабораторной работы №4 были достигнуты следующие результаты:

- освоен теоретический материал;
- создана программа на языке `Python` с использованием среды `Jupyter Notebook`;
- реализован эволюционный алгоритм ГП для функции $f_{1a}(\mathbf{x}) = \sum_{i=1}^n i x_i^2$ и получены приближения на диапазоне $[-5.12, 5.12]^9$;
- выполнена визуализация: срезы деревьев (10, 20, 50, 70, 100 поколения), финальное дерево, графики сходимости, среднего размера программ и распределения ошибок;

Список литературы

1. Методические указания по выполнению лабораторных работ к курсу «Генетические алгоритмы», стр. 119.

```

1  import math
2  import random
3  import numpy as np
4  import matplotlib.pyplot as plt
5  from dataclasses import dataclass, field
6  from typing import Callable, List, Tuple, Union, Optional, Dict
7
8  try:
9      import graphviz
10     GRAPHVIZ_AVAILABLE = True
11 except Exception:
12     GRAPHVIZ_AVAILABLE = False
13
14 SEED = 42
15 N_VARS = 9
16 DOMAIN = (-5.12, 5.12)
17 N_SAMPLES = 2000
18 POP_SIZE = 300
19 GENERATIONS = 120
20 TOURNAMENT_K = 5
21 ELITE = 2
22 P_CROSS = 0.9
23 P_MUT_NODE = 0.2
24 P_MUT_SHRINK = 0.1
25 P_MUT_GROW = 0.2
26 MAX_INIT_DEPTH = 6
27 MAX_TREE_DEPTH = 16
28 ERC_RANGE = (-5.0, 5.0)
29 PLOT_EVERY = 10
30 LAMBDA_COMPLEXITY = 1e-4
31 random.seed(SEED)
32 np.random.seed(SEED)
33
34 def p_div(x, y):
35     denom = np.where(np.abs(y) < 1e-12, 1e-12, y)
36     return x / denom
37
38 def p_pow(x, y):
39     y_clip = np.clip(y, -4.0, 4.0)
40     y_round = np.round(y_clip)
41     is_int = np.isclose(y_clip, y_round)
42     base = np.where(is_int, x, np.abs(x))
43     try:
44         out = np.power(base, np.where(is_int, y_round, y_clip))
45     except FloatingPointError:
46         out = np.power(np.clip(base, -1e6, 1e6), np.where(is_int, y_round,
47                                                         y_clip))
48     return np.clip(out, -1e6, 1e6)
49
50 def p_exp(x):
51     return np.exp(np.clip(x, -50, 50))
52
53 def p_sin(x):
54     return np.sin(x)
55
56 def p_cos(x):
57     return np.cos(x)

```

```

58 @dataclass(frozen=True)
59 class Primitive:
60     name: str
61     arity: int
62     func: Callable
63
64 FUNCTIONS: List[Primitive] = [
65     Primitive('+', 2, lambda a, b: a + b),
66     Primitive('-', 2, lambda a, b: a - b),
67     Primitive('*', 2, lambda a, b: a * b),
68     Primitive('/', 2, p_div),
69     Primitive('sin', 1, p_sin),
70     Primitive('cos', 1, p_cos),
71     Primitive('exp', 1, p_exp),
72     Primitive('^', 2, p_pow),
73 ]
74
75 @dataclass
76 class Terminal:
77     kind: str
78     index: Optional[int] = None
79     value: Optional[float] = None
80
81     def eval(self, X: np.ndarray) -> np.ndarray:
82         if self.kind == 'var':
83             return X[:, self.index]
84         else:
85             return np.full(X.shape[0], self.value, dtype=float)
86
87     def __str__(self):
88         if self.kind == 'var':
89             return f"x{self.index+1}"
90         else:
91             return f"{self.value:.4g}"
92
93 @dataclass
94 class Node:
95     prim: Union[Primitive, Terminal]
96     children: List['Node'] = field(default_factory=list)
97
98     def is_function(self) -> bool:
99         return isinstance(self.prim, Primitive)
100
101     def is_terminal(self) -> bool:
102         return isinstance(self.prim, Terminal)
103
104     def copy(self) -> 'Node':
105         return Node(self.prim, [c.copy() for c in self.children])
106
107     def depth(self) -> int:
108         if self.is_terminal():
109             return 1
110         return 1 + max((ch.depth() for ch in self.children), default=0)
111
112     def size(self) -> int:
113         return 1 + sum(ch.size() for ch in self.children)
114
115     def eval(self, X: np.ndarray) -> np.ndarray:
116         if self.is_terminal():
117             return self.prim.eval(X)
118         args = [c.eval(X) for c in self.children]

```

```

119     try:
120         if self.prim.arity == 1:
121             return self.prim.func(args[0])
122         elif self.prim.arity == 2:
123             return self.prim.func(args[0], args[1])
124     except FloatingPointError:
125         pass
126     return np.full(X.shape[0], np.nan)
127
128 def to_string(self) -> str:
129     if self.is_terminal():
130         return str(self.prim)
131     name = self.prim.name
132     if self.prim.arity == 1:
133         return f"{name}({self.children[0].to_string()})"
134     left = self.children[0].to_string()
135     right = self.children[1].to_string()
136     if name in ['+', '-', '*', '/', '^']:
137         return f"({left} {name} {right})"
138     return f"{name}({left}, {right})"
139
140 def random_terminal() -> Node:
141     if random.random() < 0.5:
142         idx = random.randrange(N_VARS)
143         return Node(Terminal('var', index=idx))
144     else:
145         val = random.uniform(*ERC_RANGE)
146         return Node(Terminal('erc', value=val))
147
148 def random_function(arity: Optional[int] = None) -> Primitive:
149     if arity is None:
150         return random.choice(FUNCTIONS)
151     cands = [f for f in FUNCTIONS if f.arity == arity]
152     return random.choice(cands)
153
154 def generate_full(depth: int) -> Node:
155     if depth <= 1:
156         return random_terminal()
157     prim = random_function(arity=random.choice([1, 2]))
158     return Node(prim, [generate_full(depth - 1) for _ in range(prim.arity)])
159
160 def generate_grow(depth: int) -> Node:
161     if depth <= 1 or (depth > 1 and random.random() < 0.3):
162         return random_terminal()
163     prim = random_function(arity=random.choice([1, 2]))
164     return Node(prim, [generate_grow(depth - 1) for _ in range(prim.arity)])
165
166 def generate_ramped(max_depth: int) -> Node:
167     d = random.randint(2, max_depth)
168     return generate_full(d) if random.random() < 0.5 else generate_grow(d)
169
170 def tournament(pop: List[Node], fitness: List[float], k: int) -> Node:
171     idxs = random.sample(range(len(pop)), k)
172     best = min(idxs, key=lambda i: fitness[i])
173     return pop[best].copy()
174
175 def all_nodes_with_parents(root: Node) -> List[Tuple[Optional[Node], int,
Node]]:
176     out = []
177     def dfs(parent, idx, node):
178         out.append((parent, idx, node))

```

```

179         for i, ch in enumerate(node.children):
180             dfs(node, i, ch)
181     dfs(None, -1, root)
182     return out
183
184 def replace_child(parent: Optional[Node], child_index: int, new_child: Node,
185                  root: Node) -> Node:
186     if parent is None:
187         return new_child
188     parent.children[child_index] = new_child
189     return root
190
191 def subtree_crossover(a: Node, b: Node, max_depth: int) -> Tuple[Node, Node]:
192     a = a.copy()
193     b = b.copy()
194     a_nodes = all_nodes_with_parents(a)
195     b_nodes = all_nodes_with_parents(b)
196     _, _, a_sub = random.choice(a_nodes)
197     _, _, b_sub = random.choice(b_nodes)
198
199     a2 = a.copy()
200     b2 = b.copy()
201     astr = a_sub.to_string()
202     bstr = b_sub.to_string()
203     a2_nodes = all_nodes_with_parents(a2)
204     b2_nodes = all_nodes_with_parents(b2)
205     cand_a = [(p, i, n) for (p, i, n) in a2_nodes if n.to_string() == astr]
206     cand_b = [(p, i, n) for (p, i, n) in b2_nodes if n.to_string() == bstr]
207     if not cand_a or not cand_b:
208         return a, b
209     pa2, ia2, a_sub2 = random.choice(cand_a)
210     pb2, ib2, b_sub2 = random.choice(cand_b)
211
212     new_a = replace_child(pa2, ia2, b_sub2.copy(), a2)
213     new_b = replace_child(pb2, ib2, a_sub2.copy(), b2)
214     if new_a.depth() > max_depth or new_b.depth() > max_depth:
215         return a, b
216     return new_a, new_b
217
218 def mutate_node_replace(tree: Node) -> Node:
219     t = tree.copy()
220     parent, idx, node = random.choice(all_nodes_with_parents(t))
221     if node.is_function():
222         ar = node.prim.arity
223         node.prim = random_function(arity=ar)
224     else:
225         node.prim = Terminal('var', index=random.randrange(N_VARS)) if
226             random.random() < 0.5 \
227             else Terminal('erc', value=random.uniform(*ERC_RANGE))
228     return t
229
230 def mutate_shrink(tree: Node) -> Node:
231     t = tree.copy()
232     funcs = [triple for triple in all_nodes_with_parents(t) if triple[2].
233               is_function()]
234     if not funcs:
235         return t
236     parent, idx, node = random.choice(funcs)
237     return replace_child(parent, idx, random_terminal(), t)

```

```

236 def random_subtree(max_depth: int) -> Node:
237     return generate_ramped(max_depth)
238
239 def mutate_grow(tree: Node, max_depth: int) -> Node:
240     t = tree.copy()
241     parent, idx, node = random.choice(all_nodes_with_parents(t))
242     new_st = random_subtree(max_depth=min(5, max_depth))
243     candidate = replace_child(parent, idx, new_st, t)
244     if candidate.depth() > max_depth:
245         return t
246     return candidate
247
248 def target_function(X: np.ndarray) -> np.ndarray:
249     idxs = np.arange(1, N_VARS + 1, dtype=float)
250     return (X**2 @ idxs)
251
252 def make_dataset(n_samples: int, low: float, high: float) -> Tuple[np.
    ndarray, np.ndarray]:
253     X = np.random.uniform(low, high, size=(n_samples, N_VARS))
254     y = target_function(X)
255     return X, y
256
257 def fitness_sse(ind: Node, X: np.ndarray, y: np.ndarray) -> float:
258     pred = ind.eval(X)
259     if np.any(np.isnan(pred)) or np.any(np.isinf(pred)):
260         return 1e50
261     err = y - pred
262     sse = float(np.sum(err * err))
263     if LAMBDA_COMPLEXITY > 0:
264         sse += LAMBDA_COMPLEXITY * (ind.size() ** 2)
265     return sse
266
267 def plot_tree_graphviz(root: Node, filename: str, title: Optional[str] =
    None):
268     if not GRAPHVIZ_AVAILABLE:
269         return
270     dot = graphviz.Digraph(comment=title or 'tree', format='png')
271     def add(node: Node) -> str:
272         node_id = str(id(node))
273         if node.is_terminal():
274             label = str(node.prim)
275             shape = 'box'
276             fill = 'honeydew'
277         else:
278             label = node.prim.name
279             shape = 'ellipse'
280             fill = 'lightsteelblue'
281         dot.node(node_id, label, shape=shape, style='filled', fillcolor=fill
            )
282         for ch in node.children:
283             ch_id = add(ch)
284             dot.edge(node_id, ch_id)
285         return node_id
286     add(root)
287     out = dot.render(filename, view=False)
288     print(f"[Graphviz] {out}")
289
290 def evolve():
291     X, y = make_dataset(N_SAMPLES, *DOMAIN)
292     X_val, y_val = make_dataset(1000, *DOMAIN)
293     pop = [generate_ramped(MAX_INIT_DEPTH) for _ in range(POP_SIZE)]

```

```

294     fitness = [fitness_sse(ind, X, y) for ind in pop]
295     sizes = [ind.size() for ind in pop]
296
297     best_idx = int(np.argmin(fitness))
298     best = pop[best_idx].copy()
299     best_fit = fitness[best_idx]
300
301     best_hist = [best_fit]
302     avg_size_hist = [float(np.mean(sizes))]
303
304     print(f"Init best SSE: {best_fit:.6g} | depth={best.depth()} | size={
305           best.size()}")
306     print("Best expr:", best.to_string())
307
308     for gen in range(1, GENERATIONS + 1):
309         new_pop: List[Node] = []
310
311         elite_idx = np.argsort(fitness)[:ELITE]
312         for ei in elite_idx:
313             new_pop.append(pop[ei].copy())
314
315         while len(new_pop) < POP_SIZE:
316             if random.random() < P_CROSS and len(new_pop) + 1 < POP_SIZE:
317                 p1 = tournament(pop, fitness, TOURNAMENT_K)
318                 p2 = tournament(pop, fitness, TOURNAMENT_K)
319                 c1, c2 = subtree_crossover(p1, p2, MAX_TREE_DEPTH)
320                 new_pop.extend([c1, c2])
321             else:
322                 p = tournament(pop, fitness, TOURNAMENT_K)
323                 c = p
324                 if random.random() < P_MUT_NODE:
325                     c = mutate_node_replace(c)
326                 if random.random() < P_MUT_SHRINK:
327                     c = mutate_shrink(c)
328                 if random.random() < P_MUT_GROW:
329                     c = mutate_grow(c, MAX_TREE_DEPTH)
330                 new_pop.append(c)
331
332         pop = new_pop[:POP_SIZE]
333         fitness = [fitness_sse(ind, X, y) for ind in pop]
334         sizes = [ind.size() for ind in pop]
335
336         gen_best_idx = int(np.argmin(fitness))
337         gen_best = pop[gen_best_idx]
338         gen_best_fit = fitness[gen_best_idx]
339         if gen_best_fit < best_fit:
340             best_fit = gen_best_fit
341             best = gen_best.copy()
342
343         best_hist.append(best_fit)
344         avg_size_hist.append(float(np.mean(sizes)))
345
346         if gen % 1 == 0:
347             print(f"Gen {gen:4d} | best SSE={best_fit:.6g} | depth={best.
348                   depth():2d} | size={best.size():4d} | avg_size={avg_size_hist
349                           [-1]:.1f}")
350
351         if gen % PLOT_EVERY == 0:
352             title = f"Gen {gen} | SSE(val)={fitness_sse(best, X_val, y_val)
353                     :.4g}"
354             plot_tree_graphviz(best, filename=f"tree_gen_{gen}", title=title
355                                )

```

```

350     final_sse_train = fitness_sse(best, X, y)
351     final_sse_val = fitness_sse(best, X_val, y_val)
352
353
354     plot_tree_graphviz(best, filename='tree_final', title=f"FINAL | SSE(val)
        = {final_sse_val:.4g}")
355
356     plt.figure(figsize=(7,4.5))
357     plt.plot(best_hist, linewidth=2)
358     plt.grid(True, alpha=0.3)
359     plt.tight_layout()
360     plt.show()
361
362     plt.figure(figsize=(7,4.5))
363     plt.plot(avg_size_hist, linewidth=2)
364     plt.xlabel("Evo")
365     plt.ylabel("Mean")
366     plt.title("Tree cap")
367     plt.grid(True, alpha=0.3)
368     plt.tight_layout()
369     plt.show()
370
371     yhat_val = best.eval(X_val)
372     residuals = y_val - yhat_val
373     plt.figure(figsize=(7,4.5))
374     plt.hist(residuals, bins=50, alpha=0.8)
375     plt.xlabel("Err")
376     plt.ylabel("Val")
377     plt.title("Valid")
378     plt.grid(True, alpha=0.3)
379     plt.tight_layout()
380     plt.show()
381
382 if __name__ == "__main__":
383     evolve()

```