

УДК 004.891.021

UDC 004.891.021

ИССЛЕДОВАНИЕ ОСНОВНЫХ ПАРАМЕТРОВ ГЕНЕТИЧЕСКОГО АЛГОРИТМА МЕТОДА ГЕНЕТИЧЕСКИХ СХЕМ В ИНТЕЛЛЕКТУАЛЬНЫХ СИСТЕМАХ, ОСНОВАННЫХ НА ЗНАНИЯХ

RESEARCH OF KEY GENETIC ALGORITHM PARAMETERS OF THE GENETIC SCHEMES METHOD IN INTELLIGENT SYSTEMS, BASED ON KNOWLEDGE

Частикова Вера Аркадьевна
к.т.н., доцент

Chastikova Vera Arkadyevna
Cand.Tech.Sci., associate professor

Кубанский государственный технологический университет, Краснодар, Россия

Kuban State Technological University, Krasnodar, Russia

Данная статья посвящена исследованию влияния основных параметров генетического алгоритма метода генетических схем на эффективность поиска решений в экспертных системах. Рассматриваются следующие параметры генетического алгоритма: численность популяции, длина бинарных кодировок, механизм отбора родительских пар, выбор схемы размножения

This article is devoted to the research of influence of genetic algorithm key parameters of genetic schemes method on efficiency of optimum decisions search in expert systems. The following parameters of genetic algorithm are considered: number of population, length of binary codes, the mechanism of parental pairs selection, the choice of the reproduction scheme

Ключевые слова: ЭКСПЕРНАЯ СИСТЕМА, БАЗА ЗНАНИЙ, ГИПОТЕЗА, ГЕНЕТИЧЕСКИЙ АЛГОРИТМ, БИНАРНАЯ КОДИРОВКА, ЧИСЛЕННОСТЬ ПОПУЛЯЦИИ

Keywords: EXPERT SYSTEM, KNOWLEDGE BASE, HYPOTHESIS, GENETIC ALGORITHM, BINARY CODING, NUMBER OF POPULATION

Для исследования эффективности использования и оценки качественных характеристик теоретически обоснованного и разработанного автором метода генетических схем (ГС) для оптимального поиска решений в продукционных экспертных системах (ЭС) с использованием генетических алгоритмов и специальным образом организованных метазнаний и сравнения его с другими методами и алгоритмами оптимизации был создан специальный программный исследовательский комплекс (ПИК) «Поиск», ядром которого является формальная модель экспертной системы (ФМЭС) [1, 4]. ПИК «Поиск» решает следующие задачи:

- создание базы знаний (БЗ) в соответствии с требуемой сложностью (число уровней, количество терминальных фактов, количество гипотез, максимальное число условий в правиле);
- сохранение базы знаний;
- выбор для исследования любой из хранимых баз знаний;

- создание и поддержка функционирования машины логического вывода [5];
- реализация ряда методов, в том числе метода ГС, поиска оптимальных решений в активной базе знаний;
- поддержка пользовательского интерфейса режима генерации метауровней базы знаний методом ГС с возможностью вариации ряда ключевых параметров как генетических алгоритмов, так и ряда других алгоритмов, используемых методом ГС;
- поддержка процесса поиска оптимальных решений в активной базе знаний любым реализованным методом, либо любой совокупностью выбранных из их числа методов;
- поддержка сравнительного анализа качества поиска выбранных для исследования методов.

Данная статья посвящена изучению эффективности разработанного автором метода ГС для применения его в экспертных системах продукционного типа. Исследованию подлежит влияние основных параметров генетического алгоритма (ГА) на эффективность поиска и их оптимизация, причем, чтобы не ограничивать область применения данного метода в определенной для него сфере, автор намеренно исключает из рассмотрения какие-либо особенности построения деревьев вывода ЭС, делая ставку на универсальность этого метода.

Рассматриваются следующие основные параметры ГА:

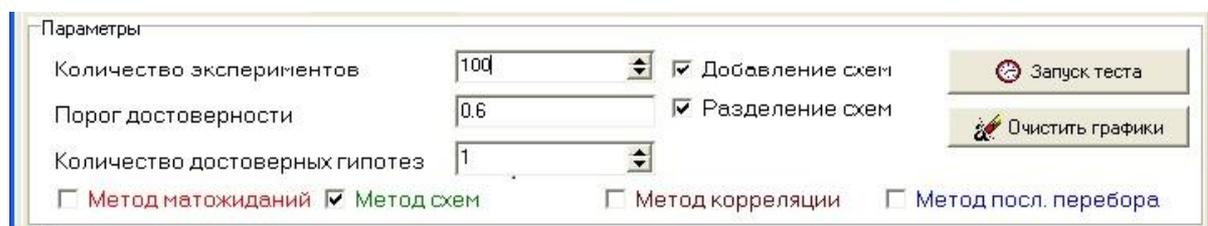
- численность популяции;
- длина бинарных кодировок;
- механизм отбора родительских пар;
- выбор схемы размножения.

Оценка метода ГС, а также сравнение его с другими методами должна базироваться на некотором «алгоритмонеиндепендентном» признаке, в качестве которого используется число оценок целевой функции (ЦФ), или иначе – число гипотез, доказанных вхолостую [5].

Численность популяции

Влияние численности популяции на точность выхода на достоверную гипотезу. Исследование влияния численности популяции на эффективность поиска оптимального решения в ЭС проводилось при различных фиксированных значениях прочих параметров [2]. Характер влияния практически везде одинаков и совпадает с приведенным ниже случаем.

Для значений параметров, установленных в окне «Консультации» ПИК «Поиск» и представленных на рисунках 1 и 2, и различной численности популяций получены результаты, приведенные в таблице 1 и на диаграмме рисунка 3.



Параметры	
Количество экспериментов	100
Порог достоверности	0.6
Количество достоверных гипотез	1
<input type="checkbox"/> Метод матожиданий	<input checked="" type="checkbox"/> Метод схем
<input type="checkbox"/> Метод корреляции	<input type="checkbox"/> Метод посл. перебора
<input checked="" type="checkbox"/> Добавление схем	Запуск теста
<input checked="" type="checkbox"/> Разделение схем	Очистить графики

Рисунок 1 – Значения параметров метода ГС окна «Консультации»

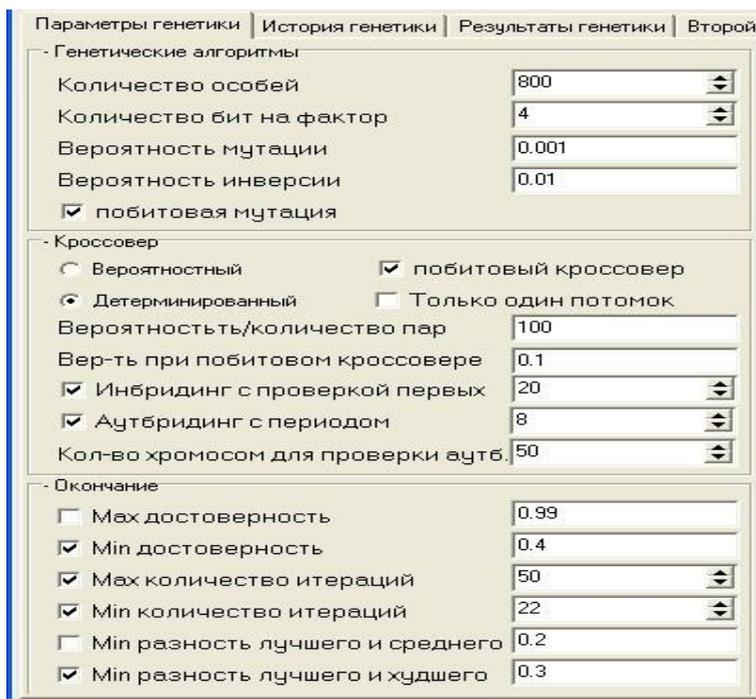


Рисунок 2 – Метод ГС: окно «Параметры генетики»

Таблица 1 – Влияние численности популяции на количество холостых гипотез

Численность популяции	30	50	100	200	300	400	500	600	700	800
Количество холостых гипотез	223	209	147	85	131	88	105	83	130	81

На диаграмме рисунка 3 приведено суммарное число холостых гипотез, полученное при выполнении 100 экспериментов при каждом запуске теста для одного и того же дерева БЗ.



Рисунок 3 – Влияние численности популяции на число гипотез, доказанных вхолостую

Линейный тренд диаграммы влияния численности популяции на число гипотез, доказанных методом ГС вхолостую, приведенный на рисунке 4, свидетельствует о постепенном убывании числа холостых гипотез с ростом численности популяции ГА, и уже при числе особей в популяции большем 200 данный метод дает стабильно удовлетворительные результаты.

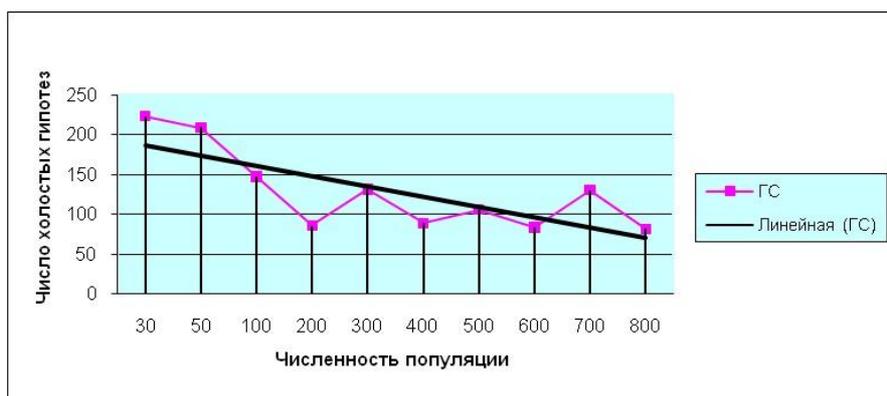


Рисунок 4 – Линейный тренд результатов поиска

Влияние численности популяции на время работы ГА. От того, какова численность популяции, безусловно, зависит время работы ГА. Так, например, для дерева БЗ с 15-ю гипотезами и 4-мя уровнями при 200 особях в популяции максимальное время выполнения ГА $\max(t_{ГА}) = 0.016с.$, а при 800 особей при прочих равных условиях – $\max(t_{ГА})=0.109 с.$ Помимо этого возрастает также время на упаковку хромосом в схемы. Тем не менее, стоит отметить тот факт, что, хотя с ростом численности популяции общее время на формирование метауровней БЗ возрастает, но это возрастание не столь масштабно. К тому же работа ГА осуществляется в рамках подсистемы «Подготовка» и на режим «Консультации» напрямую не влияет. Косвенное же влияние численности популяции на режим «Консультаций» сказывается в том, что увеличение числа особей в поколении приводит к более плотному покрытию хромосомами, а значит и схемами, ландшафта целевой функции и, следовательно, более эффективному использованию метауровней БЗ экспертной системы.

Определив порог численности популяции, преодолев который метод ГС позволяет выйти на стабильные результаты поиска решений в ЭС, в дальнейшем все внимание сосредоточено на исследовании влияния ряда параметров метода ГС, которые позволяют еще более повысить его эффективность.

Длина бинарных кодировок

Для поиска оптимального разбиения пространства параметров на гиперкубы [3], кодируемые хромосомными наборами соответствующей длины $N \cdot L$ (N – размерность задачи, L – длина кодировки одного гена), проводились испытания для $L_1=4$, $L_2=8$, $L_3=12$. Эффективность поиска (при численности популяции 300 особей) определялась по отношению к двум показателям:

- суммарному количеству холостых гипотез по 20 запускам;
- времени работы ГА.

Результаты эксперимента для нескольких деревьев БЗ с 15-ю гипотезами и 4-мя уровнями при 300 особях в популяции приведены в табл. 2 и 3 и на рис. 5, 6.

Таблица 2 – Влияние длины бинарных кодировок на количество холостых гипотез

Число бит на ген	4	8	12
Дерево 1	32	68	56
Дерево 2	5	48	41
Дерево 3	9	37	28

Как видно из представленных диаграмм, для рассматриваемых случаев при увеличении длины бинарного кода гена с 4 до 12 не происходит улучшения эффективности поиска. Вообще, многочисленные эксперименты свидетельствуют, что необходимая длина кодировки в значительной мере зависит от ландшафта целевой функции, ее увеличение

положительно сказывается на функциях, ландшафту которых присущи скачкообразные изменения, и наоборот, непрерывная функция слабо чувствительна к длине кодировки. Иногда даже, как в рассматриваемом случае, увеличение длины кодировки отрицательно влияет на показатели метода.

Очевидно, такая особенность объясняется разбросом значений приспособленности особей, имеющих одинаковые генотипы, этот разброс тем больше, чем крупнее гиперкубы разбиения пространства [3, 4]. Поэтому в процессе поиска для таких функций сложнее определить хромосомные наборы, которые бы соответствовали оптимальному решению.

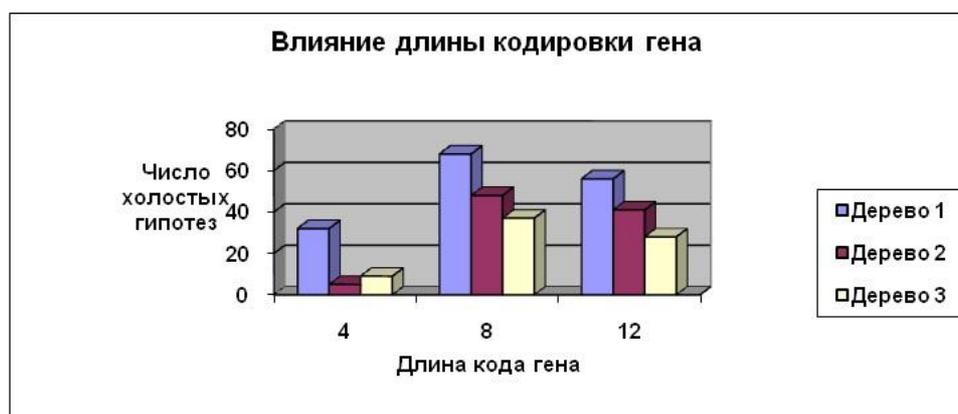


Рисунок 5 – Влияние длины кода на количество холостых гипотез

Осторожное отношение к увеличению длины кодировки вызывает и тот факт, что оно ускоряет процесс сходимости всех членов популяции к лучшему найденному решению. Такой эффект очень часто нежелателен из-за того, что большая часть пространства поиска остается неисследованной, а преждевременная сходимость может не привести к оптимальному решению, кроме того, быстрая сходимость к одной области не гарантирует обнаружения нескольких равных экстремумов. Поэтому в вопросе выбора оптимальной длины кодировки нужно достичь некоторого

компромиссного решения: с одной стороны, L должно быть достаточно большим, чтобы все-таки обеспечить быстрый поиск, с другой стороны, – по возможности малым, чтобы не допускать преждевременной сходимости и оставить алгоритму шанс отыскать несколько оптимальных значений.

Таблица 3 – Влияние длины кода гена на время выполнения ГА

	Число бит на ген	4	8	12
1000* Время ГА	Дерево 1	32	47	63
	Дерево 2	32	63	125
	Дерево 3	30	57	98

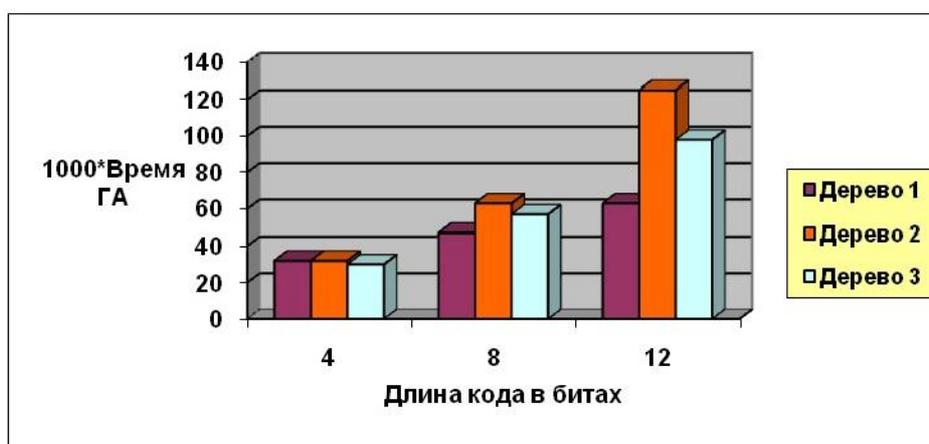


Рисунок 6 – Влияние длины кода гена на время выполнения ГА

Здесь уместно еще раз отметить, что ландшафт ЦФ при работе с ЭС продукционного типа, как правило, неизвестен, и, более того, – у каждой гипотезы он индивидуален [2,5]. Как раз выяснению характерных деталей ландшафтов во многом способствуют экспериментальные исследования с использованием ПИК «Поиск». И в отношении длины бинарной кодировки гена единственный выход – ее экспериментальное определение – такое, которое, будучи использовано для всех ГА, обеспечит в целом лучшую эффективность метода ГС.

Таким образом, экспериментально опробовав ряд режимов для сформированного выше дерева БЗ, необходимо сделать вывод, что по

обоим анализируемым параметрам предпочтительнее установить бинарную длину гена $L=4$.

Механизм отбора родительских пар

В ГС с целью предотвращения преждевременной сходимости процесса поиска к квазиоптимальному решению, заложен специальный механизм – механизм отбора.

В ГС использованы два механизма отбора, дающие в совместном использовании наилучшие результаты: элитный отбор и отбор с вытеснением.

Элитный отбор, основан на построении новой популяции только из лучших особей репродукционной группы, объединяющей в себе родителей, их потомков и мутантов. Он используется для ускорения процесса поиска и поэтому обладает потенциальной опасностью преждевременной сходимости. Поэтому в ГС применен очень осторожный подход в реализации этого отбора: в новое поколение отбирается только одна особь с наилучшим значением функции приспособленности [1].

Для компенсации быстрой сходимости, обеспечиваемой элитным отбором, необходимо предусмотреть некоторый противовес. С этой целью применен отбор родительских пар с вытеснением. В данном отборе перенос особи из репродукционной группы в популяцию нового поколения определяется не только величиной ее приспособленности, но и тем, есть ли уже в формируемой популяции следующего поколения особь с аналогичным хромосомным набором.

Из всех особей с одинаковыми генотипами предпочтение сначала, конечно же, отдается тем, чья приспособленность выше. Таким образом, достигаются две цели: во-первых, не теряются лучшие найденные решения, обладающие различными хромосомными наборами, а во-вторых, в популяции постоянно поддерживается достаточное генетическое

разнообразии. Вытеснение в данном случае формирует новую популяцию скорее из далеко расположенных особей, вместо особей, группирующихся около текущего найденного решения.

Реализованный в методе ГС механизм особенно хорошо себя показал при решении многоэкстремальных задач [1], при этом помимо определения глобальных экстремумов появляется возможность выделить и те локальные экстремумы, значения которых близки к глобальным.

Помимо описанного выше и заложенного алгоритмически в метод ГС механизма отбора пар, обеспечивающего необходимый баланс между поддержанием на удовлетворительном уровне скорости поиска и предотвращением преждевременной сходимости процесса поиска к квазиоптимальному решению, в методе ГС имеется еще и ряд дополнительных рычагов регулирования сходимости поискового процесса. С помощью этих рычагов еще в большей степени можно влиять на расширение области поиска для ее детального исследования и тем самым использовать их для доводки режима работы метода ГС до оптимального уровня.

Выбор схемы размножения

На рисунке 7 представлена панель установки значений параметров кроссовера ГА, с помощью которой перед запуском формирования метауровней БЗ можно выбрать ряд установок, в частности, схему размножения особей в ГА.

- Кроссовер	
<input type="radio"/> Вероятностный	<input checked="" type="checkbox"/> побитовый кроссовер
<input checked="" type="radio"/> Детерминированный	<input type="checkbox"/> Только один потомок
Вероятность/количество пар	100
Вер-ть при побитовом кроссовере	0.1
<input checked="" type="checkbox"/> Инбридинг с проверкой первых	20
<input checked="" type="checkbox"/> Аутбридинг с периодом	8
Кол-во хромосом для проверки аутб.	50

Рисунок 7 – Панель установки параметров кроссовера ГА

Одной из особенностей реализации ГА в ПИК «Поиск» состоит в том, что в нем предусмотрена возможность использования двух схем размножения особей:

- традиционной (вероятностной) схемы;
- альтернативной (детерминированной) схемы.

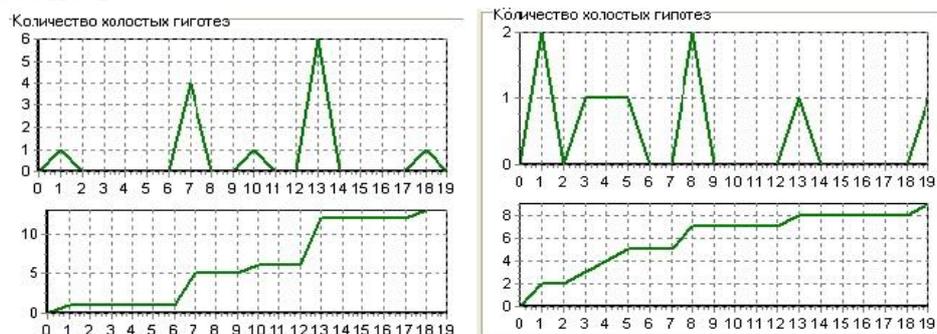
Использование традиционной схемы предполагает ограничение численности потомков посредством вероятности кроссовера. Алгоритмически вероятностный кроссовер организован следующим образом: в неупорядоченном списке хромосом в соответствии с установкой величины вероятности на панели рисунка 7 находится первый родитель, который скрещивается со следующей за ним хромосомой. Для выбора этой схемы репродукции необходимо выбрать режим «Вероятностный» и в поле «Вероятность/количество пар» установить подходящую величину вероятности кроссовера.

Альтернативная схема позволяет использовать фиксированное число брачных пар в каждом поколении. Для формирования потомков в соответствии с этой схемой на панели рисунка 7 необходимо выбрать режим «Детерминированный» и в расположенном рядом поле «Вероятность/количество пар» задать число брачных пар. В рассматриваемой ситуации алгоритм кроссовера для репродукции берет заданное число первых расположенных по порядку пар хромосом.

Вычислительные эксперименты показали, что при определенных обстоятельствах даже для простых случаев поиска оптимальных решений в продукционной экспертной системе нельзя говорить о преимуществе той или иной схемы размножения [1,4]. Так, например, для одного и того же дерева БЗ с 15-ю гипотезами и 4-мя уровнями при 300 особях в популяции в различных сеансах формирования метауровней с последующими консультациями получены результаты, приведенные на рисунках 8 и 9.

Во всех трех сеансах с различным числом экспериментов для двух схем репродукции получены схожие результаты, и небольшие преимущества, имеющиеся у одной схемы размножения в одном сеансе, могут быть потеряны ею в следующем сеансе. Таким образом, можно сделать вывод о том, что при решении задач поиска оптимальных решений в данной предметной области оказываются приемлемыми обе схемы репродукции, а выбор конкретной схемы для определенной базы знаний может быть возложен на ее проектировщика. Оправданность этого выбора может быть проверена или подкреплена проведением экспериментов с использованием ПИК «Поиск» [1].

Сеанс 1



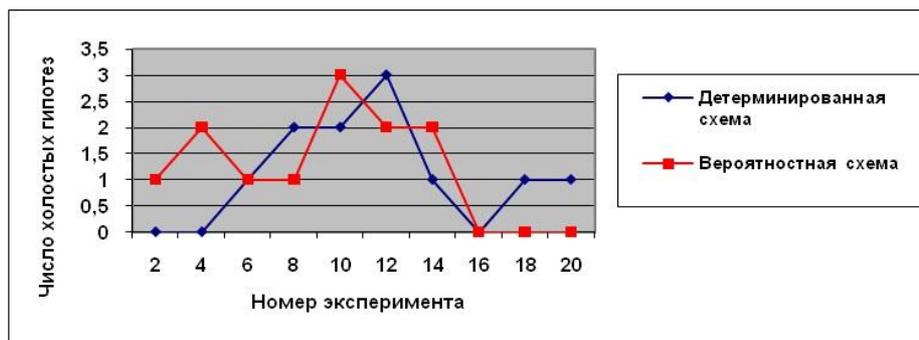
Сеанс 2



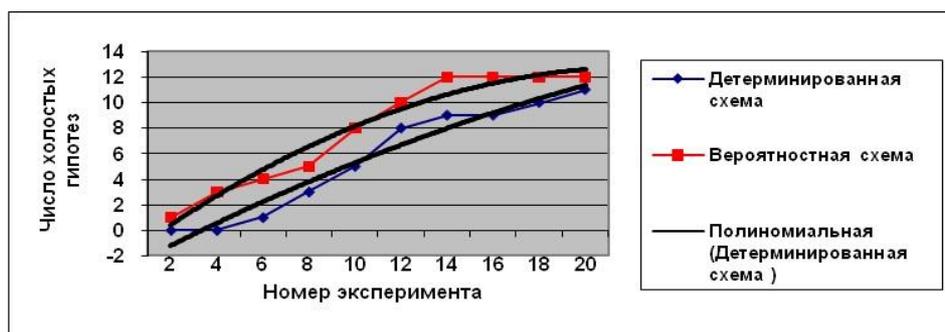
а) детерминированная схема

б) вероятностная схема

Рисунок 8 – Сравнительная оценка двух схем размножения ГА



а) текущие результаты экспериментов



б) суммарный нарастающий итог от эксперимента к эксперименту

Рисунок 9 – Сравнительная оценка двух схем размножения ГА

Литература

1. Частикова В.А. Оптимизация процессов поиска решений в интеллектуальных системах обработки экспертной информации на основе генетических алгоритмов. Автореферат диссертации на соискание ученой степени кандидата технических наук. – Краснодар: Изд-во КубГТУ, 2005.
2. Симанков В.С., Частикова В.А. Генетический поиск решений в экспертных системах. Монография. - Краснодар: Просвещение-Юг, 2008.
3. Частикова В.А. Структура и способы формирования метауровней базы знаний для эффективного поиска решений в продукционной экспертной системе. - Материалы V Международной научно-практической конференции «Интеллектуальные технологии в образовании, экономике и управлении», Воронеж, 2008.
4. Коломиец Т.В., Малыгина М.П. Формирование базы знаний экспертной системы диагностики СУБД. - Известия высших учебных заведений. Северо-Кавказский регион. Серия: Технические науки. 2007. № 3. С. 5-6.
5. Занин Д.Е., Частиков А.П. Эффективность решения задач ранжировки в информационно-поисковых системах на основе динамических нейронных сетей Хопфилда. - Известия высших учебных заведений. Северо-Кавказский регион. Серия: Технические науки. 2008. № 6. С. 62-65.