

ПРИМЕНЕНИЕ ТЕОРЕТИКО-ВЕРОЯТНОСТНОГО ОПИСАНИЯ БАЗОВЫХ ОПЕРАТОРОВ ГЕНЕТИЧЕСКИХ АЛГОРИТМОВ ДЛЯ ОЦЕНКИ ДИНАМИКИ КАЧЕСТВА ПОПУЛЯЦИИ

Петров А.В.

Кафедра информатики

Цель работы состоит в создании теоретико-вероятностного описания базовых операторов ГА для количественной оценки динамики критериев качества популяции в зависимости от вероятностных параметров алгоритма.

This paper is intended to focus on a probability-theoretical description of basic GA operators and aims at providing a stochastic parameter-dependent, quantitative estimate of criteria dynamics of the generation quality and fitness.

Введение. Практическая реализация генетических алгоритмов (ГА) и их применение для решения прикладных задач в большинстве случаев базируется на описательном материале, не содержащем математических оснований этого класса эволюционных моделей. Однако вопрос сходимости ГА представляет большой интерес не только с теоретической точки зрения. Формализация процесса эволюции может привести к вытеснению случайного поиска и замене его какими-либо детерминированными вычислительными процедурами при решении проблемы выбора оптимальных параметров ГА.

На сегодняшний день разработано несколько подходов к объяснению сходимости процесса эволюции, основанных, в том числе, на утверждениях, которые высказали Холланд (Holland, 1975) и Прайс (Price, 1970). Сформулируем и докажем первую теорему.

1. Теорема Холланда. *Теорема о схеме (Schema Theorem)* [3, 4], предложенная Дж. Холландом, рассматривается сегодня как основная теорема о ГА и является одной из самых известных попыток формального доказательства их сходимости.

О п р е д е л е н и е 1. (Схема, определенный разряд.) Пусть пространство двоичных хромосом Ξ составлено из строк длины L . Тогда схемой \mathcal{H} , по Холланду, называется строка длины L над алфавитом $\{0, 1, *\}$, где "неопределенный" символ $*$ имеет смысл "либо 0, либо 1".

Разряды, содержащие 0 и 1, носят название определенных.

О п р е д е л е н и е 2. (Порядок, определенная длина схемы.) Пусть, по Холланду, \mathcal{H} – схема. Тогда ее порядок $o(\mathcal{H})$ – это количество содержащихся в \mathcal{H} определенных разрядов.

Определенная длина схемы $\delta(\mathcal{H})$ есть расстояние между крайними определенными разрядами \mathcal{H} .

О п р е д е л е н и е 3. (Представитель схемы.) Пусть ξ – двоичная хромосома длины L из множества Ξ , \mathcal{H} – схема той же длины. Тогда ξ принадлежит \mathcal{H} (представляет ее, или является ее представителем), если значения определенных разрядов схемы \mathcal{H} и хромосомы ξ совпадают.

З а м е ч а н и е 1. Количество представителей схемы \mathfrak{H} в множестве Ξ составляет $M_{\Xi}(\mathfrak{H}) = |\{\xi \mid \xi \in \Xi, \xi \in \mathfrak{H}\}| = 2^{L-\alpha(\mathfrak{H})}$.

О п р е д е л е н и е 4. (Качество схемы.) Качество, или приспособленность, схемы \mathfrak{H} в популяции \mathfrak{G}_t определяется средней приспособленностью ее представителей в этой популяции: $\bar{F}(\mathfrak{H}) = \sum_{\mathfrak{g}_{tj} \in \mathfrak{H}} F(\mathfrak{g}_{tj}) / M_t(\mathfrak{H})$, где \mathfrak{g}_{tj} – j -я хромосома популяции \mathfrak{G}_t , $M_t(\mathfrak{H})$ – количество представителей \mathfrak{H} в \mathfrak{G}_t .

Теорема Холланда устанавливает соотношение между количеством представителей схемы в популяциях \mathfrak{G}_t и \mathfrak{G}_{t+1} , $t \geq 1$ и формулируется следующим образом [3]:

Т е о р е м а 1. (О схеме.) В генетическом алгоритме, использующем механизм пропорционального отбора, а также одноточечное скрещивание, происходящее с вероятностью P_c , при заданной кратности хромосом в популяции \mathfrak{G}_t и среднем значении функции качества хромосом \bar{F} , для каждой схемы \mathfrak{H} имеет место следующее неравенство:

$$M_{t+1}(\mathfrak{H}) \geq M_t(\mathfrak{H}) \cdot \frac{\bar{F}(\mathfrak{H})}{\bar{F}} \left(1 - P_c \cdot \frac{\delta(\mathfrak{H})}{L-1}\right). \quad (1)$$

Д о к а з а т е л ь с т в о. При пропорциональном отборе каждая хромосома \mathfrak{g}_{tj} из \mathfrak{G}_t , $1 \leq j \leq N$, выбирается для одного акта скрещивания с вероятностью

$$\Pr(S(\mathfrak{G}_t) = j) = F(\mathfrak{g}_{tj}) / \sum_{q=1}^N F(\mathfrak{g}_{tq}),$$

где $S : \mathfrak{G}_t \rightarrow \{1, 2, \dots, N\}$ – вероятностная функция отбора, N – объем популяции, что соответствует

$$NF(\mathfrak{g}_{tj}) / \sum_{q=1}^N F(\mathfrak{g}_{tq}) = F(\mathfrak{g}_{tj}) / \bar{F}$$

актам отбора \mathfrak{g}_{tj} при формировании поколения \mathfrak{G}_{t+1} . Поскольку интерес представляют только $\mathfrak{g}_{tj} \in \mathfrak{H}$, то в следующее поколение будет отобрано

$$\sum_{\mathfrak{g}_{tj} \in \mathfrak{H}} \frac{F(\mathfrak{g}_{tj})}{\bar{F}} = \frac{M_t(\mathfrak{H}) \cdot \sum_{\mathfrak{g}_{tj} \in \mathfrak{H}} F(\mathfrak{g}_{tj})}{M_t(\mathfrak{H}) \cdot \bar{F}} = M_t(\mathfrak{H}) \cdot \frac{\bar{F}(\mathfrak{H})}{\bar{F}}$$

представителей \mathfrak{H} .

При одноточечном скрещивании схема \mathfrak{H} будет разрушена, если позиция скрещивания, случайно выбранная по равномерному закону распределения, окажется между определенными разрядами. При вероятности скрещивания P_c и определенной длине $\delta(\mathfrak{H})$ вероятность разрушения схемы равна

$$P_1 = P_c \cdot \frac{\delta(\mathfrak{H})}{L-1}.$$

Отсюда следует, что количество представителей \mathfrak{H} , сохранившихся после одноточечного скрещивания, составляет

$$M_{t+1}(\mathfrak{H}) = M_t(\mathfrak{H}) \cdot \frac{\bar{F}(\mathfrak{H})}{\bar{F}} \cdot \left(1 - P_c \cdot \frac{\delta(\mathfrak{H})}{L-1}\right).$$

Нестрогое неравенство (1) указывает на то, что представители схемы \mathfrak{H} могут быть образованы при скрещивании каких-либо хромосом из \mathfrak{G}_t . ■

З а м е ч а н и е 2. Введение мутации хромосом с вероятностью единичной мутации P_m приводит к разрушению схемы порядка $o(\mathfrak{H})$ с вероятностью

$$P_2 = 1 - (1 - P_m)^{o(\mathfrak{H})} \approx P_m \cdot o(\mathfrak{H})$$

для малых P_m и равномерного закона распределения позиций мутации, так что неравенство (1) принимает вид

$$M_{t+1}(\mathfrak{H}) \geq M_t(\mathfrak{H}) \cdot \frac{\bar{F}(\mathfrak{H})}{\bar{F}} \left(1 - P_c \cdot \frac{\delta(\mathfrak{H})}{L-1} - P_m \cdot o(\mathfrak{H}) \right). \quad (2)$$

Теорема Холланда в форме (2) позволяет утверждать, что выживают схемы с низким порядком и малой определенной длиной, характеризующиеся высоким качеством. Гольдберг называет такие схемы *строительными блоками*.

Тем не менее, теорема о схеме подвергается критике, в частности, на том основании, что она обосновывает экспоненциальный рост числа строительных блоков, однако не является доказательством более быстрого, по сравнению со случайным поиском, создания новых решений с высокой приспособленностью. Более того, как отмечает Н. Радклифф, теорема Холланда справедлива даже при случайном поиске, несмотря на то, что она используется для демонстрации превосходства над ним ГА [3]. Большой доказательной силой обладает вторая из названных в п. 1 теорем, однако ее рассмотрение выходит за рамки настоящей статьи.

2. Канонический генетический алгоритм.

О п р е д е л е н и е 5. (Канонический ГА [3].) *Канонический ГА есть динамическая система вида*

$$\Pr(\xi \in \mathfrak{G}_{t+1}) = \sum_{x, x' \in X} \Pr(\xi \leftarrow \varrho(x), \varrho(x')) \times \\ \times \frac{F(\varrho(x))F(\varrho(x'))}{\bar{F}^2} \cdot \Pr(\varrho(x) \in \mathfrak{G}_t) \cdot \Pr(\varrho(x') \in \mathfrak{G}_t), \quad (3)$$

где $\Pr(\xi \leftarrow \varrho(x), \varrho(x'))$ – вероятность того, что результатом применения операторов ГА к хромосомам $\varrho(x)$ и $\varrho(x')$, кодирующим соответственно x и x' , станет хромосома ξ ; X – область определения целевой функции f .

Особенностью описания (3) является вероятностный подход к оценке степени достоверности наличия конкретной хромосомы в популяции очередного поколения. Независимо от [3] и более ранних работ, автором данной статьи получено аналогичное теоретико-вероятностное описание, позволяющее оценить вероятность наличия отдельной хромосомы в популяции следующего поколения после применения базовых операторов ГА и вероятность тождественности существующей популяции одной из возможных. С целью разумного упрощения и в точном соответствии с теоремой 1 набор базовых операторов ограничен пропорциональным отбором, одноточечным скрещиванием и единичной мутацией.

3. Создание начальной популяции. Пусть хромосома, или *генотип*, представляет собой двоичную строку длины L , и все хромосомы являются разрешенными: $N_0 = 2^L$. Пусть также ГА оперирует популяциями объемом N хромосом. Тогда существует N_0^N различных популяций, образующих универсум $\mathfrak{U} = \{u_p | u_p = \{\xi_{pj_1}, \xi_{pj_2}, \dots, \xi_{pj_N}\}\}$, $\xi_{pj_1}, \xi_{pj_2}, \dots, \xi_{pj_N} \in \Xi$, $1 \leq j_1, j_2, \dots, j_N \leq N$. Начальный отбор равновероятен, поэтому при однократном акте начального отбора каждая хромосома $\xi \in \Xi$ может быть включена в комплект u_p с вероятностью $\text{Pr}_1(\xi) = 1/N_0$, а значит, первое поколение \mathfrak{G}_1 тождественно какому-либо $u_p \in \mathfrak{U}$ с вероятностью

$$\text{Pr}(\mathfrak{G}_1 = u_p) = \prod_{j=1}^N \text{Pr}_1(\xi_{pj}) = (1/N_0)^N = 1/N_0^N.$$

Таким образом, существует два многоугольника распределения вероятностей, один из которых задает вероятность включения отдельной хромосомы в очередное поколение, второй – вероятность тождественности очередного поколения конкретной популяции.

4. Пропорциональный отбор. В условиях отсутствия скрещивания популяции \mathfrak{G}_t , $t > 1$ формируются прямым копированием в них хромосом предыдущего поколения. Потеря хромосомы носит безвозвратный характер, так как новый генетический материал не образуется. Если хромосома ξ принадлежит комплекту u_p , то вероятность ее попадания в следующее поколение полностью определяется ее качеством:

$$\begin{aligned} \text{Pr}_{t+1}^s(\xi) &= \sum_{u_p: \xi \in u_p} \frac{F(\xi) \cdot \#(\xi, u_p)}{\sum_{j=1}^N F(\xi_{pj})} \text{Pr}(\mathfrak{G}_t = u_p) = \\ &= \sum_{u_p: \xi \in u_p} \left(F(\xi) \cdot \#(\xi, u_p) \cdot \prod_{j=1}^N \text{Pr}_t(\xi_{pj}) / \sum_{q=1}^N F(\xi_{pq}) \right). \end{aligned}$$

5. Одноточечное скрещивание. Одноточечное скрещивание пары хромосом, происходящее с вероятностью P_c , существенно изменяет динамику системы, разрешая создание нового генетического материала путем перестановки фрагментов генотипов-родителей. Пусть $P_{t+1}^{(I)}(\xi)$ – вероятность переноса хромосомы $\xi \in \Xi$ в следующее поколение путем безусловного копирования, а $P_{t+1}^{(II)}(\xi)$ – вероятность образования $\xi \in \Xi$ в результате применения одноточечного скрещивания. Тогда

$$P_{t+1}^{(I)}(\xi) = (1 - P_c) \text{Pr}_{t+1}(\xi) = \bar{P}_c \cdot \sum_{u_p: \xi \in u_p} \frac{F(\xi) \cdot \#(\xi, u_p)}{\sum_{j=1}^N F(\xi_{pj})} \text{Pr}(\mathfrak{G}_t = u_p) =$$

$$= \bar{P}_c \cdot \sum_{u_p: \xi \in u_p} \left(F(\xi) \cdot \#(\xi, u_p) \cdot \prod_{j=1}^N \text{Pr}_t(\xi_{pj}) / \sum_{q=1}^N F(\xi_{pq}) \right).$$

Для определения $P_{t+1}^{(II)}(\xi)$ допустим, что хромосома $\xi \in \Xi$ может быть получена из любой пары ξ_{i_1}, ξ_{i_2} , такой что $i_1 \neq i_2$, в том случае, если выбранная позиция скрещивания λ является "благоприятной", то есть скрещивание в точке λ приведет к

получению ξ . В этом случае,

$$P_{t+1}^{(II)}(\xi) = \sum_{p=1}^{N_0^N} \left(\sum_{i_1, i_2, i_1 \neq i_2} \frac{F(\xi_{pi_1})}{\sum_{j=1}^N F(\xi_{pj})} \cdot \frac{F(\xi_{pi_2})}{\sum_{q=1}^N F(\xi_{pq}) - F(\xi_{pi_1})} \cdot \frac{n_\lambda(\xi \leftarrow \xi_{pi_1}, \xi_{pi_2})}{L-1} \cdot P_c \right) \times \\ \times \Pr(\mathfrak{G}_t = u_p) = \frac{P_c}{L-1} \sum_{p=1}^{N_0^N} \left(\sum_{i_1, i_2, i_1 \neq i_2} \frac{F(\xi_{pi_1}) F(\xi_{pi_2}) n_\lambda(\xi \leftarrow \xi_{pi_1}, \xi_{pi_2})}{\sum_{j=1}^N F(\xi_{pj}) (\sum_{q=1}^N F(\xi_{pq}) - F(\xi_{pi_1}))} \right) \prod_{r=1}^N \Pr_t(\xi_{pr}),$$

где $n_\lambda(\xi \leftarrow \xi_{pi_1}, \xi_{pi_2})$ – количество позиций скрещивания $\Lambda(\xi \leftarrow \xi_{pi_1}, \xi_{pi_2})$, "благоприятных" для получения хромосомы ξ из пары хромосом ξ_{pi_1}, ξ_{pi_2} . Формально,

$$\Lambda(\xi \leftarrow \xi_{pi_1}, \xi_{pi_2}) = \{ \lambda | \xi_k = \xi_{pi_1 k}, \xi_l = \xi_{pi_2 l}, \\ 1 \leq \lambda < L; 1 \leq k \leq \lambda; \lambda + 1 \leq l < L. \}$$

Таким образом, при одноточечном скрещивании

$$\Pr_{t+1}^c(\xi) = P_{t+1}^{(I)}(\xi) + P_{t+1}^{(II)}(\xi). \tag{4}$$

6. Единичная мутация. Поскольку оператор мутации последовательно применяется к каждому генотипу в популяции, будем считать, что безусловное копирование и скрещивание генотипов ведет к образованию популяции $\mathfrak{G}_{t+1/2}$, которая, в свою очередь, преобразуется в \mathfrak{G}_t в результате действия оператора мутации. Тогда

$$\Pr_{t+1}^m(\xi) = \begin{cases} \Pr_{t+1}^c(\xi), & \text{если } P_m = 0, \\ \Pr_{t+1}^c(\xi) + P_{t+1/2}^{(III)}(\xi) - P_{t+1/2}^{(IV)}(\xi), & \text{если } 0 < P_m \leq 1. \end{cases}$$

где $\Pr_{t+1}^c(\xi)$ определяется из равенства (4), $P_{t+1/2}^{(III)}(\xi)$ и $P_{t+1/2}^{(IV)}(\xi)$ – соответственно приращение и убыль вероятности включения хромосомы $\xi \in \Xi$ в следующее поколение в результате действия оператора мутации.

Единичная мутация¹ хромосомы $\xi \in \Xi$ ведет к ее разрушению, поэтому

$$P_{t+1/2}^{(III)}(\xi) = \sum_{u_p: \xi \in u_p} (P_m(1 - P_m)^{L-1} \cdot \#(\xi, u_p) \Pr(\mathfrak{G}_{t+1/2} = u_p)) = \\ = P_m(1 - P_m)^{L-1} \cdot \sum_{u_p: \xi \in u_p} \left(\#(\xi, u_p) \prod_{j=1}^N \Pr_t(\xi_{pj}) \right).$$

В то же время, хромосома $\xi \in \Xi$ может быть создана вновь:

$$P_{t+1/2}^{(IV)}(\xi) = \sum_{p=1}^{N_0^N} \left(\sum_{j=1}^N n_\mu(\xi \leftarrow \xi_{pj}) P_m(1 - P_m)^{L-1} / L \right) \Pr(\mathfrak{G}_{t+1/2} = u_p) = \\ = \frac{P_m(1 - P_m)^{L-1}}{L} \cdot \sum_{p=1}^{N_0^N} \left(\sum_{j=1}^N n_\mu(\xi \leftarrow \xi_{pj}) \right) \prod_{q=1}^N \Pr_t(\xi_{pq}),$$

¹Случай многократных мутаций рассматривается аналогично.

где $n_\mu(\xi \leftarrow \xi_{pj})$ – количество позиций единичной мутации $M(\xi \leftarrow \xi_{pj})$, "благоприятных" для получения хромосомы ξ из хромосомы ξ_{pj} . Формально,

$$M(\xi \leftarrow \xi_{pj}) = \{\mu \mid \xi_k = \xi_{pj k}, k \neq \mu \text{ и } \xi_k = \neg \xi_{pj k}, k = \mu\}, 1 \leq k, \mu \leq L.$$

7. Вероятностные критерии качества. С учетом сказанного, вычислим значения представляющих интерес вероятностных критериев качества популяции:

1. *средней приспособленности хромосом:*

$$\bar{F}_t(\Xi) = \sum_{i=1}^{N_0} F(\xi_i) \text{Pr}_t(\xi_i);$$

2. *средней приспособленности популяций:²*

$$\bar{F}_t(\mathcal{U}) = \sum_{p=1}^{N_0^N} \bar{F}_t(u_p) \text{Pr}(\mathfrak{G}_t = u_p) = \sum_{p=1}^{N_0^N} \left(\sum_{j=1}^N F(\xi_{pj}) \prod_{q=1}^N \text{Pr}_t(\xi_{pq}) \right) / N;$$

3. *средней меры давления отбора* – отношения приспособленности наилучшего генотипа популяции к средней приспособленности всех генотипов той же популяции:

$$\begin{aligned} \bar{\alpha}_t(\mathcal{U}) &= \sum_{p=1}^{N_0^N} \frac{\max_{j=1}^N F(\xi_{pj})}{\bar{F}_t(u_p)} \text{Pr}(\mathfrak{G}_t = u_p) = \\ &= \sum_{p=1}^{N_0^N} \left(N \max_{j=1}^N F(\xi_{pj}) \prod_{q=1}^N \text{Pr}_t(\xi_{pq}) / \sum_{r=1}^N F(\xi_{pr}) \right). \end{aligned}$$

З а м е ч а н и е 3. Можно показать, что

$$\bar{F}_t(\Xi) = \bar{F}_t(\mathcal{U}). \quad (5)$$

При $N = 1$ доказательство равенства (5) не составляет труда:

$$\bar{F}_t(\mathcal{U}) = \sum_{p=1}^{N_0^N} \bar{F}_t(u_p) \text{Pr}(\mathfrak{G}_t = u_p) = \sum_{i=1}^{N_0} F(\xi_i) \text{Pr}_t(\xi_i) = \bar{F}_t(\Xi).$$

В случае $N > 1$ необходимо учитывать очевидное ограничение, накладываемое на значения вероятностей включения хромосом в популяцию очередного поколения: $\sum_{i=1}^{N_0} \text{Pr}_t(\xi_i) \equiv 1$.

²Средняя приспособленность каждой популяции определяется средним арифметическим приспособленности хромосом, входящих в ее состав, то есть так же, как в теореме Холланда.

Список литературы

- [1] Петров А.В. К вопросу о формализации базовых операторов генетических алгоритмов // Вестник Твер. гос. техн. ун-та. Тверь: ТГТУ, 2003. Вып. 2. С. 67-69.
- [2] Петров А.В. Формальное введение в генетические алгоритмы // Сложные системы: обработка информации, моделирование и оптимизация: Сб. науч. тр. Вып. 2. Тверь: ТвГУ, 2004. С. 78-93.
- [3] Altenberg, L. The Schema Theorem and Price's Theorem. Institute of Statistics and Decision Sciences, Duke University. April 1994.
- [4] Whitley D. A Genetic Algorithm Tutorial. Technical Report CS-93-103, Colorado State University, March 1993.

(1)

$$\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = L, \forall \epsilon > 0 \exists \delta > 0$$

Работа поддержана грантом РФФИ, проект № 02-01-01080