

Шаг нулевой. Постановка задачи.

Рассмотрим набор A , представляющий множество различных генов для хромосомы в заданном конкретном локусе. Во время генерации особи ген из набора выбирается случайно и с равной вероятностью. Поставим вопрос: какой должна быть численность популяции, чтобы поддерживалось разнообразие генов в популяции? Полагаем, что мутации отсутствуют.

Для простоты положим, что хромосома состоит из одного гена. Набор доступных генов, из которых может состоять хромосома обозначим за $A = (a_1, a_2, \dots, a_M)$. Саму популяцию будем представлять в виде вектора $X = (x_1, x_2, \dots, x_N)$. Особь будем представлять в виде последовательности генов. Например $x_1 = a_1$

Шаг первый. Определения.

Введем определение полной популяции.

Полная популяция – это такая популяция, особей которых представляют ВСЕ гены из набора доступных. Например, пусть $A = (a_1, a_2, \dots, a_M)$ и $X = (a_1, a_2, \dots, a_M)$, тогда X – полная популяция.

Неполная популяция – это такая популяция, которая не является полной. Например пусть $A = (a_1, a_2, \dots, a_M)$ и $X = (a_1, a_2, a_1, \dots, a_2)$, тогда X – неполная популяция.

Замечу, что N заведомо больше или равен M . В противном случае из популяции невозможно создать полную.

Шаг первый. Формулировка идеи.

Сколько различных популяций из N особей можно составить при наборе из M генов? Очевидно, что M^N . Далеко не все эти популяции будут содержать такой набор хромосом, который обеспечивает полноту популяции. Например популяция $X = (a_1, a_2, a_1, \dots, a_2)$. Вероятность сгенерировать такую популяцию, конечно, мала, но существует популяция, и вместе они дают немалую вероятность генерации неполной популяции. Как выше было указано, мутации в данной модели исключены. В таком случае мы никогда не получим популяцию, в которой бы была особь например с геном a_3 .

Обозначим за Y множество всех неполных популяций. Тогда вероятность сгенерировать полную популяцию есть $P = 1 - \frac{CardY}{M^N}$. Очевидно, что $CardY$ зависит от размера популяции N . Поэтому, определяя желаемую вероятность получить полную популяцию P , мы сможем регулировать параметр N .

Шаг второй. Реализация идеи.

Итак, пусть имеется набор $A = (a_1, a_2, \dots, a_M)$ и популяция $X = (x_1, x_2, \dots, x_N)$. Число различных популяций, которые мы можем составить равняется M^N . Однако это количество включает в себя такие популяции, которые не являются полными. Одна из таких популяций есть такая популяция, в которой нет особи с геном a_M $X' = (a_1, a_2, \dots, a_{M-1}, a_1, \dots)$. Однако заметим, что эта популяция является полной, если $A' = (a_1, a_2, \dots, a_{M-1})$. Другими словами, она является полной при некотором наборе A' , являющимся подмножеством множества A !

Исходя из того, что особей в популяции всё также N штук, а набор теперь состоит из генов на единицу меньше, то количество возможных популяций, которые мы можем составить на данном наборе генов есть $(M - 1)^N$. Заметим, что мы точно также могли бы рассматривать популяцию, в которой отсутствует особь с геном a_1 :

$X' = (a_2, a_3, \dots, a_{M-1}, a_2, \dots)$ и набор $A' = (a_2, a_3, \dots, a_M)$. Поэтому на самом деле число различных популяций, которые мы можем составить есть $C_M^{M-1}(M-1)^N$. Однако это число включает лишние популяции, о которых я скажу позже.

Итак, теперь $X' = (a_1, a_2, \dots, a_{M-1}, a_1, \dots)$ – полный набор при $A' = (a_1, a_2, \dots, a_{M-1})$. Число различных популяций, которые мы можем составить равняется $(M-1)^N$. Однако это количество включает в себя такие популяции, которые не являются полными, например $X'' = (a_1, a_2, \dots, a_{M-2}, a_1, \dots)$. Однако популяция $X'' = (a_1, a_2, \dots, a_{M-2}, a_1, \dots)$ является полной при наборе $A'' = (a_1, a_2, \dots, a_{M-2})$. Число различных популяций, которые могут быть составлены из этого набора $C_{M-1}^{M-2}(M-2)^N$

Дальше следуют повторяющиеся рассуждения до тех пор, пока мы не упрямся на набор из одного гена, например $A''' = (a_1)$. Из данного набора можно образовать единственную популяцию, которая будет являться полной $X''' = (a_1, \dots, a_1)$. Таким образом искомая вероятность:

$$P = \frac{M^N - C_M^{M-1}[(M-1)^N - C_{M-1}^{M-2}[(M-2)^N - \dots]]}{M^N} =$$

$$1 - \frac{C_M^{M-1}[(M-1)^N - C_{M-1}^{M-2}[(M-2)^N - \dots]]}{M^N}$$

Но это не итоговый вид формулы. Каждый раз, когда мы образуем подмножество **A** возникают возможные повторяющиеся друг друга популяции. Например, рассмотрим набор генов из 4-ёх генов: **{a, b, c, d}**.

1. На первом шаге мы образуем подмножества
{a, b, c}, {a, b, d}, {a, c, d}, {b, c, d}
Видим, что получается по $1! = 1$ одинаковому подмножества (тождественному самому себе)

2. На втором шаге мы образуем подмножества
{a, b}, {a, c}, {b, c}, {a, b}, {a, d}, {b, d}, {a, c}, {a, d}, {c, d}, {b, c}, {b, d}, {c, d}
Видим, что получается по $2! = 2$ одинаковых подмножества

3. На третьем шаге мы образуем подмножества
{a}, {b}, {a}, {c}, {b}, {c}, {a}, {b}, {a}, {d}, {b}, {d}, {a}, {c}, {a}, {d}, {c}, {d}, {b}, {c}, {b}, {d}, {c}, {d}
Видим, что получается по $3! = 6$ одинаковых подмножеств.

Если положить первоначально набор из 5-ти генов **{a, b, c, d, e}**, то в последней итерации получится $4! = 24$ одинаковых подмножеств. Интуиция подсказывает, что на 6-ти одинаковых генах будет уже $5! = 120$ одинаковых подмножеств. Этот факт вносит корректировку в формулу:

$$P = 1 - \frac{\frac{C_M^{M-1}}{1!}[(M-1)^N - \frac{C_{M-1}^{M-2}}{2!}[(M-2)^N - \dots[\frac{C_{M-(M-1)+1}^{M-(M-1)}}{(M-1)!}(M-(M-1))^N]]}{M^N}$$