

## ОПИСАНИЕ ДИВЕРГЕНЦИИ СУБПОПУЛЯЦИЙ В ИЕРАРХИЧЕСКОЙ СИСТЕМЕ ПРИ АНАЛИЗЕ ИЗОНИМИИ. II. ВЕРОЯТНОСТИ НЕИЗОНИМНЫХ ВСТРЕЧ

© 2023 г. В. П. Пасеков<sup>1</sup>, \*

<sup>1</sup>Федеральный исследовательский центр “Информатика и управление” Российской академии наук, Москва, 119991 Россия

\*e-mail: pass40@mail.ru

Поступила в редакцию 01.05.2023 г.

После доработки 22.05.2023 г.

Принята к публикации 25.05.2023 г.

Рассматриваются типичные метапопуляции человека с иерархической подразделенностью на части (субпопуляции), соответствующие классификации на основе административно-территориального деления (скажем, село, сельсовет, район, область и так далее) или генеалогического подхода, базирующегося на этногенезе, а также на других принципах биологической классификации. Целью настоящей работы является анализ общих свойств распределения концентрации фамилии по субпопуляциям при их иерархической структуре. Внимание концентрируется на описании фамильной дивергенции субпопуляций, в качестве показателя которой рассматривается общая вероятность встреч ( $H_s$ ) лиц с разными фамилиями, понимаемая как вероятность встречи в выбранной наугад субпопуляции, единицы наблюдения, из рассматриваемой метапопуляции с иерархической структурой. Данная вероятность является фамильным аналогом концентрации гетерозигот в метапопуляции со случайным скрещиванием в ее субпопуляциях, единицах наблюдения. Получено разложение  $H_s$  по уровням иерархии, обобщающее эффект Валунда в популяционной генетике. Общая вероятность неизонимных встреч в иерархически подразделенной метапопуляции меньше вероятности случайных встреч в ней на сумму средних внутригрупповых дисперсий концентрации фамилии, соответствующих отдельным уровням. Такие свойства являются чисто статистическими характеристиками иерархической структуры, а не особенностью конкретной популяционной системы и не выводятся из закономерностей той или иной модели микроэволюции. Они вычислительно формулируются одинаково для любой иерархической системы, хотя в общем случае не совпадают количественно. Полученные результаты относятся к сельским и городским иерархическим метапопуляциям как отдельным компонентам всего населения.

**Ключевые слова:** иерархическая структура популяций, метапопуляции, концентрации фамилии в субпопуляциях человека, характеристики неоднородности субпопуляций, разложение вероятности встреч носителей разных фамилий по уровням иерархии.

**DOI:** 10.31857/S0016675823110097, **EDN:** TNBXXK

Случайный генный дрейф [1] приводит к генетической дивергенции популяций с общим происхождением. Аналогично в результате случайного дрейфа фамилий [2, 3] происходит фамильная дивергенция таких популяций. Оба этих процесса протекают синхронно в одной и той же популяции с ограниченной численностью. Фамильный состав следующего поколения с численностью  $N$  мужской составляющей можно упрощенно рассматривать как результат случайной выборки с возвращением размера  $N$  из совокупности фамилий мужчин родительского поколения. Генетический состав нового поколения диплоидной популяции при отсутствии систематических давлений является результатом выборки  $4N$  гамет роди-

тельского поколения ( $2N$  размер популяции с учетом женщин и  $4N$  размер выборки гамет, формирующих  $2N$  диплоидных потомков — по две гаметы на потомка). В результате имеется полная аналогия модели фамильного дрейфа с генетической моделью диплоидной популяции со случайным скрещиванием, рассматриваемой в отношении одного аутосомного локуса с множественными аллелями. Роль фамилий играют аллели. Аналогом случайной встречи пары индивидуумов с разными фамилиями (следовательно, с разными родоначальниками) является появление гетерозиготного генотипа при случайном скрещивании, а аналогом пар однофамильцев будут гомозиготы.

Очевидно, что закономерности динамики фамильного и генетического состава популяции имеют один и тот же характер выборочного дрейфа. Когда закономерности данных процессов сходны, то по наблюдению за результатами одного из них можно делать выводы о результатах протекания другого. Однако оба процесса не идентичны. Для одной и той же диплоидной популяции случайный процесс выборочных колебаний концентраций фамилий в ряду поколений интенсивнее генного дрейфа (соответствует вчетверо меньшей выборке фамилий, чем выборка гамет). Подробнее об этом говорится в [2, 3].

Популяционно-генетический анализ оперирует с данными по генетической структуре популяций. Получение данных по фамильной структуре менее трудоемко и менее затратно. Поэтому желательнее проанализировать пути экстраполяции результатов дивергенции популяций по фамильным данным в генетические выводы. При этом важно уметь описывать свойства дивергенции в типичных для населения ситуациях, и мы начнем с такого описания.

Для метапопуляций человека типична иерархическая подразделенность на части (субпопуляции), соответствующие классификации на основе административно-территориального деления (скажем, село, сельсовет, район, область и так далее) или генеалогического подхода, базирующегося на этногенезе, а также на других принципах биологической классификации. В предыдущей части [4] данной работы в качестве характеристики фамильной дивергенции в иерархически подразделенной метапопуляции служила дисперсия распределения фамилии по субпопуляциям. Теперь обратимся к другой характеристике — общей (полной) вероятности случайной встречи (столкновения) индивидуума с фиксированной фамилией с носителем какой-либо другой, которую понимаем как вероятность встречи в выбранной наугад субпопуляции, единице наблюдения, из рассматриваемой метапопуляции с иерархической структурой.

Напомним, что под *иерархической системой* понимается система с многоуровневой структурой, элементы которой связаны отношениями подчинения, причем отдельный элемент, в свою очередь, является иерархической системой (более низкого ранга). У нас иерархической системой является метапопуляция, подразделенная на субпопуляции различного уровня иерархии. Повторим, что вероятность рассматриваемой встречи аналогична доле (концентрации) гетерозигот в популяционной генетике. Это позволяет при соответствующих предположениях оценивать такой важный генетический показатель как коэффициент инбридинга популяции [5].

Итак, объектом нашего анализа служит реальная или теоретическая метапопуляция  $s$  с иерархической подразделенностью на субпопуляции  $\{s_i\}$  по уровням иерархии. При подразделенности системы *отношение подчинения означает вхождение подсистемы низкого уровня в качестве составной части в соответствующую подсистему (группу) более высокого ранга*. Например, субпопуляции (скажем, сёла) группируются в сельсоветы, районы и т.д. с уровнями иерархии 2, 3 ... соответственно, и множество субпопуляций на каждом из этих уровней представляет собой разбиение рассматриваемой метапопуляции. Здесь выполняется иерархическое подчинение одной субпопуляции другой, т.е. вхождение первой в качестве части во вторую с более высоким уровнем иерархии.

Дадим несколько слов относительно обозначений и терминологии. Под концентрациями фамилий в популяции подразумеваются концентрации однофамильцев (носителей данной фамилии). Векторы набраны полужирным шрифтом, к обозначениям фамильных аналогов популяционно-генетических характеристик (дисперсии, коэффициента фамильного инбридинга) добавлено окончание  $s$  ( $Vs$  и  $Fs$  соответственно). Символ тождества ( $\equiv$ ) у нас используется в смысле равенства по определению, а символ  $\blacktriangleleft$  обозначает конец доказательства.

Напомним, что идентификацию и положение субпопуляции в иерархической структуре можно задавать с помощью мультиномеров, как это было сделано нами ранее в предыдущей части [4]. Пусть номер конкретного села (первый уровень) обозначается как  $s_1$ ; номер сельсовета (второй уровень), куда входит это село, как  $s_2$ ; номер района (третий уровень), включающего указанные сельсовет и село, как  $s_3$ ; и т.д. Тогда мультиномер  $s_1 \equiv (s_1, s_2, s_3, \dots)$  однозначно определяет рассматриваемое село среди прочих сел первого уровня, вектор  $s_2 \equiv (s_2, s_3, s_4, \dots)$  — идентификатор сельсовета, вектор  $s_i \equiv (s_i, s_{i+1}, \dots)$  идентифицирует субпопуляцию  $i$ -го уровня среди прочих таких же субпопуляций внутри соответствующей группы следующего уровня  $i + 1$ . Это аналогично почтовому адресу, в котором административные составляющие заменены числами.

В результате субпопуляция  $s_1$  входит в надлежащую субпопуляцию  $s_2$ , а  $s_i$  входит в  $s_{i+1}$  и т.д., т.е. субпопуляция некоторого уровня иерархии включает в себя в качестве своей части соответствующие субпопуляции более низкого уровня. Между объектами и их идентификаторами имеется взаимно однозначное соответствие, и мы иногда будем писать идентификатор вместо названия объекта (села, сельсовета и т.д.). Кроме того, повторим, что множество субпопуляций на каждом отдельно выбранном уровне иерархии представ-

ляет собой разбиение всей метапопуляции, т.е. составляет ее целиком.

Данный способ нумерации приложим также к концентрации  $x$  интересующей фамилии, которую в селе  $s_1$  будем обозначать как  $x(s_1)$ , а концентрацию фамилии в субпопуляции  $i$ -го уровня как  $x(s_i)$ . Когда некоторые из номеров  $\{s_i\}$  рассматриваются как случайные величины (например, при выборе субпопуляции наугад), а другие как фиксированные, то для наглядности будем писать фиксированные величины после вертикальной черты. В частности, будем рассматривать  $x(s_{i-1}|s_i)$  как *случайную величину*, значениями которой являются концентрации фамилии, получаемые при выборе наугад субпопуляции ( $i - 1$ )-го уровня из фиксированной группы  $s_i$ . *Первый аргумент  $s_{i-1}$  у  $x(s_{i-1}|s_i)$  указывает на случайно выбираемую субпопуляцию уровня  $i - 1$ , а второй ( $s_i$ ) — на содержащую ее фиксированную группу уровня  $i$ .*

Таким образом, концентрация фамилии в фиксированной субпопуляции  $s_i$  обозначается как  $x(s_i)$ , а  $x(s_{i-1}|s_i)$  дает случайную величину, принимающую значения концентраций фамилии в субпопуляциях уровня  $i - 1$ , выбираемых наугад из  $s_i$ . Обозначим математическое ожидание (среднее значение) случайной величины  $x(s_{i-1}|s_i)$  как  $x(s_i)$ , а ее дисперсию как  $Vs(x(s_{i-1}|s_i))$ .

*Цель настоящей работы — продолжение анализа в [4] общих свойств распределения концентрации фамилии по субпопуляциям при их иерархической подразделенности. Они являются чисто статистическими характеристиками иерархической структуры, а не особенностью конкретной популяционной системы.* Анализируемые свойства присущи любой иерархически подразделенной метапопуляции и не выводятся из закономерностей той или иной модели микроэволюции. Они вычислительно формулируются одинаково для любой иерархической системы, хотя в общем случае не совпадают количественно. Теоретически это может позволить выделить специфические особенности исследуемого материала. По-прежнему внимание концентрируется на изучении дивергенции субпопуляций, в качестве показателя которой теперь рассматривается не дисперсия концентрации фамилии в субпопуляциях [4], а вероятности встреч пары лиц с неизонимными фамилиями. Обратим внимание на то, что *в рассматриваемых вероятностях учитывается порядок фамилий в паре*, что не вполне традиционно.

Структура изложения материала следующая. После напоминания системы идентификации субпопуляций в метапопуляции с иерархической структурой рассматриваются вероятности встреч носителей разных фамилий в качестве показателей дивергенции субпопуляций на разных уровнях иерархии. Далее подробно рассматривается из-

менчивость субпопуляций в случае многоуровневой иерархической структуры метапопуляции. Выведено разложение по уровням иерархии общей (полной) вероятности неизонимных встреч в выбранной наугад субпопуляции, единице наблюдения. В найденном разложении каждому отдельному уровню иерархии соответствует неотрицательный компонент, равный соответствующей средней внутригрупповой дисперсии концентрации фамилии.

## ПОДРАЗДЕЛЕННЫЕ МЕТАПОПУЛЯЦИИ

В случае изучения популяций человека классификация и группировка данных часто производятся на основе административно-территориального деления, имеющего иерархический характер (скажем, село, сельсовет, район, область и др.), генеалогического подхода на основе этногенеза, лингвистических данных и пр. Получаемая группировка субпопуляций приближенно является иерархической. Иерархическая структура метапопуляции отражается на ее свойствах, в частности на распределении фамилий в популяциях человека, где типична опора на официальные данные иерархического характера, сбор и обработку материалов в соответствии с ними. Настоящая статья мотивирована анализом фамильных данных с ориентацией на популяционную генетику. Использование фамилий для получения выводов о генетической структуре популяций основывается на существующих параллелях в передаче на популяционном уровне потомкам фамилий и аутосомных аллелей (см., например, [2, 3]). Плодотворность такого использования продемонстрирована в ряде работ [6–8] (изонимные браки) [9] (фундаментальная монография), в том числе в исследованиях популяций России [10] (медико-генетические аспекты).

Повторим, что у нас объектами являются субпопуляции, их совокупность образует метапопуляцию (иерархическую систему), подчинение означает принадлежность одной субпопуляции другой в качестве ее части. При этом у самой метапопуляции наивысший уровень иерархии, и все субпопуляции являются ее частями. В роли числового признака субпопуляции может рассматриваться любая числовая характеристика, но мы ограничимся преимущественно дискретными признаками.

На популяционном уровне в качестве признака субпопуляции рассматривается среднее значение соответствующей характеристики для множества входящих в нее субпопуляций более низкого уровня. Вычисление средних значений основано на суммировании, и желательно выяснить результат разбиений метапопуляции или субпопуляции на более мелкие группы для соответствующего среднего значения признака. При его вычислении суммирование происходит по входящим в объект базовым единицам самого низкого первого уровня,

которые у нас считаются заданными (служат начальными данными). Варьирование характера разбиения приводит лишь к перегруппировке слагаемых, а результат остается неизменным. Таким образом, среднее значение признака рассматриваемого объекта не зависит от характера разбиения последнего, что подробнее рассматривается далее.

Важным примером характеристики субпопуляции, состоящей из индивидуумов с дискретным признаком  $T$ , является средняя величина  $T$ , принимающего значения 1 и 0 в зависимости от его категории (1 для интересующей категории и 0 в противном случае; скажем, интересующей категорией может быть фамилия, аллель, гетерозигота и т.п.). Покажем, что среднее значение  $T$  в произвольной субпопуляции  $s_i$  представляет собой не что иное как концентрацию в ней носителей интересующей категории. Эта концентрация также не зависит от характера подразделенности  $s_i$ .

**Замечание 1.** Пусть субпопуляция  $s_i$  состоит из индивидуумов с дискретным признаком  $T$ , кодируемым единицей при интересующей категории и кодируемым нулем в противном случае.

Тогда среднее значение  $T$  в  $s_i$  равно концентрации носителей рассматриваемой категории в данной субпопуляции (доле ее носителей; иначе говоря, вероятности случайно встретить такого носителя в субпопуляции при одинаковых шансах встретить любого индивидуума).

**Доказательство.** Среднее значение  $T$  находится как сумма значений  $T$  (по условию это 1 и 0) с весами, равными вероятностям, с которыми  $T$  принимает данные значения, т.е. равно

$$1 \times Pr\{T = 1\} + 0 \times Pr\{T = 0\} = Pr\{T = 1\},$$

где  $Pr\{T = \dots\}$  обозначает вероятность наблюдения соответствующего значения  $T$ . Вероятность  $Pr\{T = 1\}$  интерпретируется как вероятность при выборе наугад из  $s_i$  индивидуума получить носителя интересующей категории. При одинаковых шансах выбрать любого индивидуума вероятность  $Pr\{T = 1\}$  равна доле (концентрации) носителей этой категории в  $s_i$ . ◀

Далее рассматривается только популяционный уровень организации, т.е. сами носители остаются за кадром, а анализируемыми объектами являются субпопуляции, характеризующие средними значениями (концентрациями) интересующих показателей, например фамилии. Для субпопуляций самого низкого первого уровня иерархии значения концентрации фамилии полагаем заданными исходно (рассматриваем их как начальные данные).

Обратим внимание на то, что даже когда предметом изучения является единственная иерархически подразделенная метапопуляция, выводы относительно ее свойств можно делать в вероятностных терминах в результате трактовки харак-

теристик субпопуляции как результатов ее случайного выбора из соответствующего множества субпопуляций, частей разбиения метапопуляции. При этом поскольку концентрация является средним значением для  $T$ , то к ней приложимы известные свойства средних значений (математических ожиданий).

Рассмотрим метапопуляцию  $s_j$ . Покажем, что концентрация  $T(s_j)$  носителей интересующей категории дискретного признака  $T$  в  $s_j$  равна средней величине для его концентраций в составляющих разбиение  $s_j$  субпопуляциях  $\{s_i\}$  и не зависит от разбиения.

**Замечание 2.** Пусть метапопуляция  $s_j$  с общей численностью  $N(s_j)$  разбита на какие-либо группы  $\{s_i\}$  с численностями  $\{N(s_i)\}$ . Положим численность носителей рассматриваемой категории  $T$  в метапопуляции  $s_j$  обозначена как  $T_N(s_j)$ .

Тогда концентрация  $T(s_j) \equiv T_N(s_j)/N(s_j)$  носителей данной категории  $T$  в метапопуляции  $s_j$ , произвольно разбитой на группы  $\{s_i\}$ , находится как

$$\begin{aligned} T(s_j) &= \sum_{s_i \in s_j} T(s_i) \frac{N(s_i)}{N(s_j)} \equiv \\ &\equiv \sum_{s_i \in s_j} T(s_i) Pr(s_i | s_j) = E_{s_i} \{T(s_i) | s_j\}, \end{aligned}$$

$$0 \leq Pr(s_i | s_j) \equiv \frac{N(s_i)}{N(s_j)} \leq 1, \quad \sum_{s_i} Pr(s_i | s_j) = 1,$$

т.е. как математическое ожидание случайной величины  $T(s_i | s_j)$ , концентрации  $T$  в  $s_i$  при выборе наугад  $s_i$  из  $s_j$  с вероятностью  $Pr(s_i | s_j) \equiv N(s_i)/N(s_j)$ . Здесь  $E$  является символом операции получения математического ожидания, нижний индекс у  $E$  указывает на используемую при усреднении переменную.

Значение  $T(s_j)$  не зависит от разбиения метапопуляции  $s_j$ , и  $T(s_j)$  не меняется при его изменении.

**Доказательство.** Концентрация  $T(s_j)$  рассматриваемых носителей в разбитой на группы  $\{s_i\}$  метапопуляции  $s_j$  определяется как отношение их численности  $T_N(s_j)$  в  $s_j$ , складывающейся из численностей в подгруппах  $\{T_N(s_i)\}$ , к общей численности  $N(s_j)$  метапопуляции  $s_j$ , складывающейся из численностей ее подгрупп  $\{N(s_i)\}$ . Для субпопуляций  $s_i$  и  $s_j = \{s_i\}$  имеем

$$\begin{aligned} T_N(s_i) &\equiv T(s_i)N(s_i), \\ T_N(s_j) &= \sum_{s_i \in s_j} T_N(s_i) = T(s_i)N(s_i), \end{aligned}$$

$$T(s_j) \equiv T_N(s_j)/N(s_j) = \sum_{s_i \in s_j} T(s_i) \frac{N(s_i)}{N(s_j)}.$$

Следовательно, концентрация  $T(s_j)$  в  $s_j$  равняется среднему взвешенному значению для концентраций  $\{T(s_i)\}$  в составляющих разбиение  $s_j$  субпопуля-

циях  $\{s_j\}$  с весами, равными относительным численностям субпопуляций  $\{N(s_i)/N(s_j)\}$ . Очевидно, эти веса в сумме по  $\{s_j\}$  дают единицу и неотрицательны. Их значения  $\{N(s_i)/N(s_j)\}$  можно интерпретировать как вероятности  $\{Pr(s_i|s_j)\}$  попасть на индивидуума из соответствующей субпопуляции  $s_j$  при выборе наугад индивидуума из  $s_j$  (как вероятности выбора субпопуляции  $s_i$  из  $s_j$  при одинаковых шансах каждого быть выбранным).

Таким образом, формула для  $T(s_j)$  показывает, что  $T(s_j)$  является математическим ожиданием для  $T(s_i|s_j)$  при выборе наугад соответствующей субпопуляции  $s_i$  из  $s_j$ :

$$\begin{aligned} T(s_j) &\equiv \sum_{s_i \in s_j} T(s_i) \frac{N(s_i)}{N(s_j)} = \\ &= \sum_{s_i \in s_j} T(s_i) Pr(s_i|s_j) = E_{s_i} \{T(s_i)|s_j\}. \end{aligned}$$

Очевидно, приведенный вывод верен для любого разбиения  $s_j$  на субпопуляции  $\{s_i\}$ . При этом ее общая численность  $N(s_j)$  и общее количество носителей  $T_N(s_j)$  в ней не изменятся в результате выбора другого разбиения. Следовательно, их отношение (концентрация носителей) не зависит от разбиения  $s_j$  и находится по такой же формуле. ◀

Типичным признаком при изучении изонимии служит концентрация  $x$  какой-либо фамилии (точнее, носителей данной фамилии). Как и для признака  $T$ , среднее значение  $x(s_j)$  для концентрации фамилии в некоторой метапопуляции  $s_j$  равно средней величине для концентраций в составляющих ее разбиение субпопуляциях. Скажем, если  $s_j$  разбита на  $\{s_i\}$ , то  $x(s_j) = E_{s_i} \{x(s_i)|s_j\}$ , а если рассматривается разбиение на  $\{s_i\}$ , то  $x(s_j) = E_{s_i} \{x(s_i)|s_j\}$ .

Это свойство средних значений признаков в подразделенных популяциях рассматривается далее с более общих позиций как выражение полного математического ожидания.

### СВОЙСТВА СРЕДНИХ ЗНАЧЕНИЙ В ИЕРАРХИЧЕСКОЙ СИСТЕМЕ

Пусть дана какая-либо метапопуляция  $s$  с иерархической структурой подразделенности. На первом уровне  $s$  разбита на субпопуляции  $\{s_1\}$ , а на  $i$ -м уровне на  $\{s_j\}$ . Понятно, что каждая из субпопуляций  $s_j$  в свою очередь, подразделена на соответствующие подмножества субпопуляций первого уровня. Положим, что рассматривается случайная переменная величина  $T$  (ею может быть концентрация носителей какой-нибудь категории признака  $T$ ). Пусть значения концентрации  $\{T(s_1)\}$  заданы для всех субпопуляций первого уровня иерархии  $\{s_1\}$ , а значения  $T(s_j)$  на более высоких уровнях находятся по определению как

$$\begin{aligned} T(s_i) &\equiv E_{s_1} \{T(s_1)|s_i\} = \sum_{s_1} T(s_1) Pr(s_1|s_i), \\ i &= 2, 3, \dots; \quad T(s) \equiv E_{s_1} \{T(s_1)|s\}. \end{aligned} \quad (1)$$

Таким образом, на каждом из уровней иерархии  $i$  значение  $T(s_i)$  является средним значением для  $T(s_1)$  при выборе наугад  $s_1$  из  $s_i$ .

Рассмотрим некоторые из свойств математического ожидания для иерархических систем, формулируемые для произвольной случайной величины  $T$ . Напомним широко используемую далее формулу полного математического ожидания (частным случаем которой является полученная выше формула для  $T(s)$  в предыдущем замечании):

$$E\{T\} = E_A \{E_T \{T|A\}\} = E_T \{E_T \{T|A_i\} Pr(A_i)\}.$$

Здесь случайное событие  $A$  принадлежит множеству событий  $\{A_i\}$ , составляющих полную систему несовместимых случайных событий, реализующихся с вероятностями  $\{Pr(A_i)\}$  и таких, что обязательно происходит одно из них;  $E_T \{T|A_i\}$  означает условное (условие пишем после вертикальной черты) математическое ожидание для случайной величины  $T$  (нижний индекс у  $E$  указывает на переменную, которая служит для усреднения) при условии реализации соответствующего случайного события  $A_i$ .

У нас при выборе наугад субпопуляции  $s_i$  из  $s_j$  роль полной системы случайных событий играет множество субпопуляций  $\{s_i\}$ , представляющих собой разбиение рассматриваемой популяции  $s_j$ . Аналогично при случайном выборе  $s_1$  из  $s_j$  полную систему образуют субпопуляции  $\{s_1\}$ . Тогда для случайной величины  $T$  формула полного математического ожидания выглядит следующим образом:

$$\begin{aligned} E\{T(s_j)\} &\equiv E_{s_1} \{T(s_1)|s_j\} = \\ &= E_{s_1} \{E_{s_i} \{T(s_1)|s_i\}|s_j\} = E_{s_i} \{T(s_i)|s_j\}, \end{aligned} \quad (2)$$

т.е. среднее значение  $T$  у некоторой группы  $s_j$  является средним для математических ожиданий  $E_{s_1} \{T(s_1)|s_j\} \equiv T(s_i)$  у разбивающих ее подгрупп  $\{s_i\}$ . Это согласуется с полученными ранее формулами для среднего значения случайной величины  $T$  или для концентрации  $x$  в подразделенной популяции.

В частности, при  $j = i + 1$  получаем рекуррентное уравнение

$$T(s_{i+1}) = E_{s_i} \{T(s_i)|s_{i+1}\}. \quad (3)$$

Если в (2) вместо  $s_j$  рассматривать всю метапопуляцию  $s$ , разбитую на  $\{s_j\}$ , по формуле полного математического ожидания и определению (1)

$$\begin{aligned} T(s) &= E_{s_1} \{T(s_1)|s\} = E_{s_1} \{E_{s_i} \{T(s_1)|s_i\}|s\} = \\ &= \sum_{s_i} T(s_i) Pr(s_i|s), \quad i = 1, 2, \dots, \end{aligned} \quad (4)$$

т.е.  $E_{s_1} \{E_{s_i} \{T(s_1)|s_i\}|s\} = E_{s_1} \{T(s_1)|s\}$ .

Мы имеем два выражения для одного и того же значения  $T(\mathbf{s})$  – по определению (1) и по формуле (4), т.е.

$$T(\mathbf{s}) \equiv E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}\} = E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}_i\}, \quad (5)$$

$$i = 1, 2, 3, \dots, \quad T(\mathbf{s}_i) \equiv E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}_i\}.$$

Эта формула показывает в явном виде равенство полного (общего) среднего значения  $T(\mathbf{s})$  средней величине для математических ожиданий отдельных частей  $\{T(\mathbf{s}_i)\}$  одного и того же уровня иерархии  $i$ , какими бы ни были эти части и выбранный уровень ( $i$  принимает значения, соответствующие отдельным уровням иерархии метапопуляции  $\mathbf{s}$ ). Поскольку разбиение  $\mathbf{s}$  на  $\{\mathbf{s}_i\}$  может быть любым, то  $T(\mathbf{s})$  не зависит от характера подразделенности  $\mathbf{s}$ .

Аналогичная формула верна для  $T(\mathbf{s}_j)$ , когда  $\mathbf{s}_j$  разбита на субпопуляции  $\{\mathbf{s}_i\}$ . Если  $\mathbf{s}$  состоит из  $\{\mathbf{s}_j\}$ , а  $\mathbf{s}_j$  из  $\{\mathbf{s}_i\}$ , то

$$E_{s_j} \{E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}_j\} | \mathbf{s}\} = E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}\}. \quad (6)$$

Повторим, что здесь случайная величина  $T$  принимает значения  $\{T(\mathbf{s}_i)\}$  с вероятностями  $\{Pr(\mathbf{s}_i | \mathbf{s}_j)\}$  выбора наугад соответствующей субпопуляции  $\mathbf{s}_i$  на уровне иерархии  $i$  из субпопуляции  $\mathbf{s}_j$  на уровне  $j$ . Среднее значение  $E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}_j\}$  для  $T(\mathbf{s}_i | \mathbf{s}_j)$  обозначается как  $T(\mathbf{s}_j)$  и характеризует составную часть  $\mathbf{s}_j$  метапопуляции  $\mathbf{s}$ . Это среднее значение само является случайной величиной  $T(\mathbf{s}_j | \mathbf{s})$ , если  $\mathbf{s}_j$  наугад выбирается из фиксированной метапопуляции  $\mathbf{s} \equiv \{\mathbf{s}_j\}$ . Согласно формуле полного математического ожидания полное (общее) среднее значение рассматриваемой случайной величины  $T(\mathbf{s}_i | \mathbf{s}_j)$  для всей метапопуляции  $\mathbf{s} \equiv \{\mathbf{s}_j\}$  равно математическому ожиданию для средних значений в ее отдельных частях  $\{\mathbf{s}_j\}$ , что формально отображается формулами (5)–(6).

Признак  $T$  определен для субпопуляций на любом уровне иерархии. При этом его базовыми значениями являются значения  $\{T(\mathbf{s}_1)\}$ , заданные для субпопуляций  $\{\mathbf{s}_1\}$  первого уровня. Через них выражаются величины  $T$  для субпопуляций последующих уровней как средние значения в соответствующих группировках базовых значений. Поэтому обозначение  $T(\mathbf{s}_i)$  подразумевает вычисление при условии значений  $\{T(\mathbf{s}_1)\}$ , а средние величины  $T$  для ряда субпопуляций одного и того же уровня не зависят от выбранного уровня и совпадают с величиной для всей метапопуляции  $\mathbf{s}$ :

$$E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}\} = E_{s_i} \{T(\mathbf{s}_i) | \mathbf{s}_i\} =$$

$$= \dots E_{s_j} \{T(\mathbf{s}_j) | \mathbf{s}\} = \dots = T(\mathbf{s}),$$

что говорит о независимости значения  $T(\mathbf{s})$  от разбиений  $\mathbf{s}$ , в том числе неиерархических.

## ПРИМЕРЫ СЛУЧАЙНЫХ ВЕЛИЧИН ПРИ ИЗУЧЕНИИ ИЗОНИМИИ

Далее в качестве случайных величин будем рассматривать такие признаки субпопуляции как концентрация фамилии  $x$  и вероятность  $Hs$  неизоимных встреч (столкновений) в субпопуляции. Проанализируем также свойства дисперсии  $V_{s_{in}}(x(\mathbf{s}_i | \mathbf{s}_j))$ , т.е. дисперсии стоящего в скобках аргумента (случайной величины  $x(\mathbf{s}_i | \mathbf{s}_j)$ ).

### КОНЦЕНТРАЦИЯ ФАМИЛИИ

Рассмотрим подразделенную субпопуляцию  $\mathbf{s}_2$ . Положим, что в качестве  $T$  рассматривается концентрация  $x$  носителей интересующей фамилии,  $x(\mathbf{s}_2)$  обозначает концентрацию фамилии во всей фиксированной подразделенной субпопуляции  $\mathbf{s}_2$ . Она находится по определению (1) как среднее значение для концентраций  $\{x(\mathbf{s}_1 | \mathbf{s}_2)\}$  в субпопуляциях единицах наблюдения первого уровня иерархии  $\{\mathbf{s}_1\}$  внутри  $\mathbf{s}_2$ :

$$x(\mathbf{s}_2) \equiv E_{s_1} \{x(\mathbf{s}_1) | \mathbf{s}_2\} = \sum_{s_1} x(\mathbf{s}_1) Pr(\mathbf{s}_1 | \mathbf{s}_2).$$

Среднее значение концентрации фамилии  $x(\mathbf{s}_3)$  для распределения  $\{x(\mathbf{s}_1 | \mathbf{s}_3)\}$  в фиксированной субпопуляции  $\mathbf{s}_3$  находится по определению (1) и по формуле полного математического ожидания (4) как

$$x(\mathbf{s}_3) \equiv E_{s_1} \{x(\mathbf{s}_1 | \mathbf{s}_3)\} = E_{s_2} \{x(\mathbf{s}_2) | \mathbf{s}_3\} =$$

$$= E_{s_2} \{E_{s_1} \{x(\mathbf{s}_1) | \mathbf{s}_2\} | \mathbf{s}_3\} = \sum_{s_2} x(\mathbf{s}_2) Pr(\mathbf{s}_2 | \mathbf{s}_3).$$

Согласно (4) для произвольного уровня иерархии  $j$  (при  $j > i$ , когда  $\mathbf{s}_i$  входит в  $\mathbf{s}_j$  как составная часть, элемент разбиения  $\mathbf{s}_j$ )

$$x(\mathbf{s}_j) \equiv E_{s_i} \{x(\mathbf{s}_i) | \mathbf{s}_j\} = \sum_{s_i} x(\mathbf{s}_i) Pr(\mathbf{s}_i | \mathbf{s}_j) =$$

$$= \sum_{s_i} x(\mathbf{s}_i) Pr(\mathbf{s}_i | \mathbf{s}_j), \quad \mathbf{s}_j = \{\mathbf{s}_i\}.$$

Если здесь вместо  $\mathbf{s}_j$  взять всю метапопуляцию  $\mathbf{s}$ , то ясно, что концентрация фамилии  $x(\mathbf{s})$  не зависит от использования при ее подсчетах уровня разбиения  $\mathbf{s}$  (как мы видели для произвольной случайной величины  $T$  в (5)).

Для иерархически подразделенной метапопуляции  $\mathbf{s}$  мы имеем при заданных значениях  $\{x(\mathbf{s}_1)\}$ :

$$x(\mathbf{s}_i) \equiv E_{s_1} \{x(\mathbf{s}_1) | \mathbf{s}_i\} =$$

$$= \sum_{s_1} x(\mathbf{s}_1) Pr(\mathbf{s}_1 | \mathbf{s}_i), \quad \mathbf{s}_i = \{\mathbf{s}_1\};$$

$$\begin{aligned} x(s_j) &\equiv E_{s_1} \{x(s_1)|s_j\} = E_{s_1} \{x(s_i)|s_j\} = \\ &= E_{s_1} \{E_{s_i} \{x(s_1)|s_i\}|s_j\} = \sum_{s_i} x(s_i)Pr(s_i|s_j), \\ s_j &= \{s_j\}, \quad s_i = \{s_i\}; \\ x(s) &\equiv E_{s_1} \{x(s_1)|s\} = \\ &= E_{s_1} \{x(s_i)|s\} = \sum_{s_i} x(s_i)Pr(s_i|s). \end{aligned}$$

При  $i = j - 1$  получаем еще одну запись рекуррентного уравнения  $x(s_j) \equiv E_{s_{j-1}} \{x(s_{j-1})|s_j\}$ , которую можно исследовать с целью получения выражения для  $x(s_j)$  на произвольном уровне иерархии  $j$  в зависимости от начальных значений  $\{x(s_1)\}$ .

### КЛАССИФИКАЦИЯ ДИСПЕРСИИ КОНЦЕНТРАЦИИ

Рассмотрим внутригрупповую дисперсию  $V_{s_{in}}(s_2) \equiv V_s(x(s_1|s_2))$  концентрации  $x(s_1)$  по субпопуляциям  $\{s_1\}$  внутри  $s_2$ , т.е. дисперсию распределения  $\{x(s_1)\}$  значений  $x$  у субпопуляций  $\{s_1\}$ , являющихся единицами наблюдения и содержащихся в  $s_2$ . Она является случайной при выборе наугад  $s_2$  из  $s_3$ . Признак  $V_{s_{in}}$  определен для субпопуляций с уровнем не ниже второго (говорить о внутригрупповой дисперсии для  $s_1$  не имеет смысла, так как в ней нет субпопуляций, различия между которыми характеризуются дисперсией). Формально можно определить значение  $V_s(x(s_n|s_n))$  равным нулю.

Начальные значения для случайной величины  $x$  заданы для субпопуляций первого уровня иерархии как  $\{x(s_1)\}$ . Эти значения  $\{x(s_1)\}$  внутри  $s_2$  определяют начальные значения для  $V_s(x(s_1|s_2)) = V_{s_{in}}(x(s_1|s_2))$ , внутригрупповой дисперсии  $V_{s_{in}}$  для  $s_2$ . Повторим, что данные значения случайны, если  $s_2$  наугад выбирается из фиксированной субпопуляции  $s_3$ .

Напомним классификацию дисперсий в случае метапопуляции  $s_3$  с тремя уровнями иерархии. Здесь выделяют три типа дисперсий. Под *общей (полной) дисперсией*  $V_{s_{tot}}(x(s_1|s_3))$  метапопуляции понимают дисперсию распределения концентрации  $\{x(s_1)\}$  признака по субпопуляциям первого уровня иерархии, содержащимся в  $s_3$ . Рассматриваемая метапопуляция  $s_3$  разбита на группы (субпопуляции)  $\{s_2\}$  и дисперсия  $V_{s_{betw}}(x(s_2|s_3))$  распределения концентрации  $\{x(s_2)\}$  по этим группам называется *межгрупповой дисперсией* в метапопуляции. Наконец, каждая группа состоит из соответствующих субпопуляций  $\{s_1\}$  с концентрациями  $\{x(s_1)\}$ . Дисперсия  $V_s(x(s_1|s_2))$  распределения концентраций  $\{x(s_1)\}$  внутри  $s_2$  называется *внутригрупповой дисперсией*  $V_{s_{in}}(x(s_1|s_2))$  группы  $s_2$ . Метапопуляция  $s_3$  в

целом характеризуется *средней внутригрупповой дисперсией*  $E_{s_2} \{V_{s_{in}}(x(s_2|s_3))\}$ .

**Ремарка 3.** Когда у метапопуляции количество уровней иерархии более трех, приведенная классификация дисперсии становится расплывчатой. Отметим, что дисперсия  $V_s$  полностью определяется ее аргументом, а индексы *tot*, *betw*, *in* служат только для облегчения ориентации в отношении аспекта рассмотрения данной дисперсии. Без них свободно можно было бы обойтись без ущерба для логики изложения.

При многоуровневой иерархии метапопуляцию  $s$  по-прежнему можно охарактеризовать, скажем, в отношении концентрации  $x$  некоторого признака *полной дисперсией*  $V_{s_{tot}}(x(s_1|s))$ , являющейся дисперсией распределения признака по субпопуляциям самого низкого, первого, уровня иерархии, содержащимся в  $s$ .

Еще одной характеристикой субпопуляции  $s$  является *межгрупповая дисперсия*  $s$  на уровне  $i$ , т.е. дисперсия  $V_s(x(s_i|s))$  распределения концентрации фамилии по группам уровня  $i$ , на которые разбита исходная субпопуляция  $s$ . Для каждого из промежуточных уровней  $u$  с  $s$  будет своя межгрупповая дисперсия.

Другой характеристикой  $s$  служит *внутригрупповая дисперсия*  $V_{s_{in}}(s_i|s_j)$  субпопуляции  $s_j$  на уровне  $i$ , т.е. дисперсия распределения  $\{x(s_i|s_j)\}$  значений  $x$  у субпопуляций  $s_i$  в  $s_j$ . Внутригрупповую дисперсию можно интерпретировать когда  $i$  равно единице как полную дисперсию  $s_j$ , т.е. когда в  $s_j$  субпопуляции  $\{s_i\}$  являются единицей наблюдения.

Напомним, что внутригрупповая дисперсия носит случайный характер при выборе наугад  $s_j$  из  $s$ . Тогда вся субпопуляция  $s$  характеризуется математическим ожиданием для  $\{V_{s_{in}}(s_i|s)\}$ , иначе говоря, *средней внутригрупповой дисперсией на уровне  $i$  для субпопуляций  $j$ -го уровня иерархии в  $s$ :*

$$E_{s_j} \{V_{s_{in}}(s_i|s_j)|s\} = \sum_{s_j} V_{s_{in}}(s_i|s_j)Pr(s_j|s).$$

Вспомним, что при анализе концентрации фамилии  $x$  ее среднее значение в субпопуляции (метапопуляции) высокого уровня, скажем, в  $s_k$  выражается согласно (5) через значения  $u$  входящих в нее субпопуляций уровнем ниже как  $x(s_k) = E_{s_j} \{x(s_j)|s_k\}$ . Казалось бы, если в (5) положить  $T(s_j) = V_{s_{in}}(s_j|s_j)$ , то аналогично  $V_{s_{in}}(s_i|s_k) = E_{s_j} \{V_{s_{in}}(s_i|s_j)|s_k\}$ . Однако это равенство неверно, что мы подробно проанализируем далее. Суть в том, что перетасовка данных при группировке субпопуляций  $\{s_j\}$  не приводит к изменению результатов для концентрации фамилии как показано ранее, в отличие от дисперсии концентрации.

### ВЕРОЯТНОСТИ ВСТРЕЧ (СТОЛКНОВЕНИЙ) В ПОДРАЗДЕЛЕННОЙ МЕТАПОПУЛЯЦИИ

Рассмотрим концентрацию (вероятность)  $Hs(\mathbf{s})$  встреч (столкновений) индивидуума с данной фамилией с носителем какой-либо другой в метапопуляции  $\mathbf{s}$ , определяемую как среднее значение  $Hs$  для групп одного уровня иерархии. Согласно свойствам средних значений (5) результат усреднения не зависит от того, какое при этом используется разбиение  $\mathbf{s}$  на группы на данном уровне и от самого уровня. Поскольку концепция встреч  $Hs$  скорее умозрительная, чем дающая легко доступный метод определения  $Hs$  в произвольной популяции (за исключением случая брачных пар), то возникает проблема как задать  $Hs$  на базовом уровне.

Здесь приходится прибегать к косвенным методам, основанным на дополнительных предположениях. Таким предположением является допущение о случайном характере встреч в группах  $\{s_1\}$ . Оно разумно, когда разбиение метапопуляции состоит из территориально малых групп (базовых субпопуляций  $\{s_1\}$  самого низкого уровня иерархии), в каждой из которых существованием каких-либо препятствий для встреч (столкновений) индивидуумов и нарушениями внешних границ при встречах можно пренебречь.

При случайной встрече двух индивидуумов комбинирование ее участников в пары происходит независимо. В результате независимости участников случайных столкновений вероятность встречи равна произведению вероятностей встретить каждого из участников по отдельности, равных их концентрациям. Так, вероятность  $Hs(x(s_1))$  встречи индивидуума с интересующей фамилией (с концентрацией  $x(s_1)$  ее носителей в  $s_1$ ) и индивидуума с какой-либо иной с учетом порядка участников, очевидно, находится как  $Hs_0(x(s_1)) \equiv x(s_1)(1 - x(s_1))$ . Эта формула аналогична выражению для концентрации гетерозигот по закону Харди–Вайнберга в популяционной генетике при учете порядка аллелей в гетерозиготе, поскольку в обоих случаях наблюдаем случайные объединения в парах.

Метапопуляция, состоящая из неподразделенных субпопуляций, характеризуется средней величиной для значений  $Hs$  по субпопуляциям. Например, общая (полная) вероятность  $Hs(x(s_1|s_2))$  встречи индивидуумов с разными фамилиями в *подразделенной* метапопуляции  $s_2$  (скажем, в сельсовете  $s_2$ , подразделенном на села) определяется как среднее значение  $Hs(x(s_1))$  для входящих в нее субпопуляций (сел  $\{s_1\}$ ), под которой понимается вероятность при выборе наугад субпопуляции  $s_1$  из рассматриваемой метапопуляции. Тем самым подразумевается, что *встречи происходят только внутри отдельных сел* по аналогии с рядом моделей популяционной генетики, в которых случай-

ное скрещивание идет внутри отдельной элементарной популяции после формирования ее генетического состава к моменту скрещивания с учетом всех факторов динамики. Чтобы подчеркнуть, что когда речь идет о средней вероятности встреч в метапопуляции, а не в отдельной неподразделенной субпопуляции, то среднее значение вероятности встреч называем *общей (полной) вероятностью*  $Hs(x(s_2))$ .

Очевидно, при произвольном разбиении метапопуляции  $\mathbf{s}$  на  $\{s_1\}$  у нас нет аргументов в пользу случайного характера встреч в каждой из субпопуляций  $\{s_1\}$  и соответственно в пользу вероятности встреч в виде  $Hs(s_1) = Hs_0(s_1) \equiv x(s_1)(1 - x(s_1))$ . При разбиении метапопуляции  $\mathbf{s}$  на  $\{s_1\}$ , когда в  $\{s_1\}$  встречи случайны,

$$\begin{aligned} Hs(\mathbf{s}) &= \sum_{s_1} Hs_0(s_1) Pr(s_1|\mathbf{s}) = \\ &= \sum_{s_1} x(s_1)(1 - x(s_1)) Pr(s_1|\mathbf{s}). \end{aligned}$$

Здесь и далее для краткости мы пишем  $Hs(\mathbf{s})$  вместо  $Hs(x(\mathbf{s}))$ .

Повторим, что роль вероятности  $Hs(s_1) = Hs(x(s_1))$  случайной встречи рассматриваемых индивидуумов с учетом порядка их фамилий в неподразделенной субпопуляции  $s_1$  с концентрацией интересующей фамилии  $x(s_1)$  играет  $x(s_1)(1 - x(s_1))$ , а вероятность встречи  $Hs(s_2)$  в метапопуляции  $s_2 = \{s_1\}$  определяется как среднее значение  $Hs(s_1)$  при выборе наугад  $s_1$  из метапопуляции  $s_2$ . Таким образом, полной вероятностью встречи  $Hs(s_2)$  в метапопуляции  $s_2$  будет

$$Hs(s_2) \equiv E_{s_1} \{Hs(s_1)|s_2\} = E_{s_1} \{x(s_1)(1 - x(s_1))|s_2\}.$$

Напомним, что  $Pr(s_1|s_2)$  обозначает вероятность случайного выбора села  $s_1$  из фиксированного сельсовета  $s_2$ ;  $E_{s_1}$  – символ операции получения математического ожидания (среднего значения) случайной величины, стоящей в фигурных скобках; нижний индекс  $s_1$  у  $E_{s_1}$  указывает на служащую для усреднения переменную  $s_1$ . По определению вероятности  $Hs(s_i)$  рассматриваемой встречи для произвольного уровня иерархии  $i$  задается формулой

$$\begin{aligned} Hs(s_i) &\equiv E_{s_1} \{Hs(s_1)|s_i\} = \sum_{s_1} Hs(s_1) Pr(s_1|s_i) = \\ &= \sum_{s_1} x(s_1)(1 - x(s_1)) Pr(s_1|s_i). \end{aligned}$$

На примере вероятности (концентрации)  $Hs$  мы видим, что в качестве характеристики субпопуляций может служить не только концентрация фамилии  $x$ , но и функция  $Hs$  от  $x$ . В каждой субпопуляции  $s_1$  значение  $Hs$  находится как функция  $x(1 - x)$  от концентрации  $x$  в  $s_1$ . Казалось бы есте-

ственно предположить, что среднее значение  $Hs$  по субпопуляциям  $s_1$ , составляющим  $s_i$ , т.е.  $Hs(x(s_i))$ , находится как такая же функция  $x(1 - x)$  теперь от значения  $x$  в  $s_i$ . Однако  $Hs(x(s_i))$  не равно  $x(s_i)(1 - x(s_i))$ , что на первый взгляд может показаться парадоксальным.

**ВЕРОЯТНОСТЬ ВСТРЕЧ  
В ПОДРАЗДЕЛЕННОЙ МЕТАПОПУЛЯЦИИ  
И КОЭФФИЦИЕНТ ФАМИЛЬНОГО  
ИНБРИДИНГА**

Обратимся к анализу вероятности (доли, концентрации) встреч  $Hs$ , являющейся аналогом концентрации гетерозигот в популяционной генетике. Как уже говорилось, эту концентрацию можно уподобить вероятности соответствующего столкновения в некоторой совокупности частиц. Напомним, что здесь подразумевается отыскание  $Hs(s_i)$  в подразделенной популяции при условии заданных значений вероятности столкновений в субпопуляциях первого уровня.

На первом уровне иерархии вероятность встречи  $Hs(s_1)$  в субпопуляции  $s_1$  двух индивидуумов (с заданной фамилией с концентрацией  $x(s_1)$ ) и какой-либо другой с учетом их порядка) находится в предположении случайных столкновений как  $Hs(s_1) \equiv Hs_0(s_1) = x(s_1)(1 - x(s_1))$  аналогично закону Харди–Вайнберга. Для субпопуляций более высокого уровня  $i$  используется обозначение  $Hs(s_i)$ , в котором зависимость от значений  $Hs(s_1)$  неясна. Концентрация  $Hs(s_i)$  находится как среднее значение для  $\{Hs(s_1)\}$  в субпопуляциях  $\{s_1\}$ , составляющих разбиение  $s_i$ . Когда  $Hs(s_1) = x(s_1)(1 - x(s_1))(1 - F(s_1))$ , для вероятности  $Hs(s_i)$  используется обозначение  $Hs(s_i|Fs_1)$ , явным образом указывающее на вид  $Hs(s_1)$ .

По определению значение  $Hs$  для каждой субпопуляции уровня иерархии  $i$  выше единицы равно средней величине для значений  $Hs$  у всех входящих в  $s_i$  субпопуляций на первом уровне (или на одном из других предшествующих  $i$  уровней) согласно (5), а на первом уровне в предположении случайного характера встреч равно  $x(s_1)(1 - x(s_1))$ :

$$\begin{aligned} Hs(s_1) &\equiv Hs_0(s_1) \equiv x(s_1)(1 - x(s_1)), \\ Hs(s_i) &\equiv E_{s_i} \{Hs_0(s_1)|s_i\} = \\ &= E_{s_i} \{x(s_1)(1 - x(s_1))|s_i\} = \\ &= \sum_{s_1} x(s_1)(1 - x(s_1)) Pr(s_1|s_i), \quad i = 2, 3, \dots \end{aligned} \tag{7}$$

Иначе говоря, полная вероятность случайной встречи (концентрация  $Hs(s_i)$ ) в  $s_i$  находится как среднее значение  $Hs(s_1)$  при выборе наугад субпопуляции  $s_1$  из  $s_i$ . По формуле полного математического ожидания (5)

$$Hs(s_j) \equiv E_{s_2} \{Hs(s_2)|s_j\} = E_{s_i} \{Hs(s_i)|s_j\}, \quad j > i > 1.$$

В частности, значения  $Hs$  на соседних уровнях иерархии связывает рекуррентное соотношение

$$\begin{aligned} Hs(s_{i+1}) &\equiv E_{s_i} \{Hs(s_i)|s_{i+1}\}, \text{ или} \\ Hs(s_i) &\equiv E_{s_{i-1}} \{Hs(s_{i-1})|s_i\}; \\ Hs(s_1) &\equiv x(s_1)(1 - x(s_1)). \end{aligned} \tag{8}$$

Отсюда вытекает согласующаяся с (6) запись формулы полного математического ожидания для  $Hs$

$$E_{s_i} \{Hs(s_i)|s_{i+1}\} = E_{s_i} \{E_{s_{i-1}} \{Hs(s_{i-1})|s_i\}|s_{i+1}\}. \tag{9}$$

Повторим, что (8) означает, что полная вероятность рассматриваемых встреч в субпопуляции  $s_i$  является средней величиной для вероятностей встреч в субпопуляциях  $\{s_{i-1}\}$ , вместе составляющих разбиение  $s_i$ . Выясним зависимость вероятности подобной встречи от фамильной дивергенции между субпопуляциями, рассмотренной нами ранее в первой части [4] в терминах дисперсии концентрации рассматриваемой фамилии.

**Результат 4 (фамильный аналог эффекта Валунда).** Пусть дана метапопуляция  $s$ , разбитая на субпопуляции  $\{s_1\}$  с концентрациями  $\{x(s_1)\}$  рассматриваемой фамилии в них и с вероятностями рассматриваемых неизонимных встреч  $\{x(s_1)(1 - x(s_1))\}$ .

Тогда общая (полная) вероятность  $Hs(s)$  случайной встречи двух индивидуумов (с фиксированной фамилией и с какой-либо иной при учете порядка фамилий) в наугад выбранной из  $s$  субпопуляции  $s_1$  имеет вид

$$\begin{aligned} Hs(s) &\equiv E_{s_1} \{Hs(s_1)|s\} = \\ &= \sum_{s_1} x(s_1)(1 - x(s_1)) Pr(s_1|s) = \\ &= x(s)(1 - x(s)) - V_{S_{betw}}(x(s_1|s)), \\ V_{S_{betw}}(x(s_1|s)) &\equiv V_S(x(s_1|s)) = \\ &= \sum_{s_1} (x(s_1) - x(s))^2 Pr(s_1|s). \end{aligned} \tag{10}$$

Здесь  $V_{S_{betw}}(x(s_1|s))$  обозначает межгрупповую дисперсию концентраций  $\{x(s_1|s)\}$  рассматриваемой фамилии в субпопуляциях (группах)  $\{s_1\}$ , составляющих разбиение метапопуляции  $s$ ;  $Pr(s_1|s)$  дает вероятность случайного выбора субпопуляции  $s_1$  из метапопуляции  $s$ .

Иным образом общую (полную) вероятность  $Hs(s)$  можно представить как

$$\begin{aligned} Hs(s) &= x(s)(1 - x(s)) \left( 1 - \frac{V_{S_{betw}}(x(s_1|s))}{x(s)(1 - x(s))} \right) = \\ &= x(s)(1 - x(s))(1 - Fs(s)), \\ Fs(s) &\equiv \frac{V_{S_{betw}}(x(s_1|s))}{x(s)(1 - x(s))} = \frac{V_{S_{betw}}(x(s_1|s))}{Hs_0(s)} \geq 0, \end{aligned} \tag{11}$$

где  $Fs(\mathbf{s})$  является фамильным аналогом коэффициента инбридинга  $S$ . Райта в популяционной генетике.

**Доказательство.** Пусть неподразделенная субпопуляция  $\mathbf{s}_1$  рассматривается как единица наблюдения. Например, в сельском населении это может быть село, а метапопуляция  $\mathbf{s}$  соответствует, скажем, сельсовету. Дисперсия  $V_{S_{betw}}(x(\mathbf{s}_1|\mathbf{s}))$  относится к распределению концентрации заданной фамилии по селам  $\{\mathbf{s}_1\}$  внутри фиксированного сельсовета  $\mathbf{s}$ . Данная дисперсия  $V_{S_{betw}}(x(\mathbf{s}_1|\mathbf{s}))$  характеризует фамильную дивергенцию сел внутри сельсовета.

Напомним, что для любой случайной величины  $x$  выполняется равенство  $E\{x^2\} = (E\{x\})^2 + V(x)$ , так как  $V(x) \equiv E\{x^2\} - (E\{x\})^2$ . Для дальнейшего эту формулу с учетом равенства  $E_{s_i}\{x(\mathbf{s}_i)|\mathbf{s}_{i+1}\} = x(\mathbf{s}_{i+1})$  удобно представить как

$$E_{s_i}\{x^2(\mathbf{s}_i)|\mathbf{s}_{i+1}\} = (E_{s_i}\{x(\mathbf{s}_i)|\mathbf{s}_{i+1}\})^2 + V_S(x(\mathbf{s}_i|\mathbf{s}_{i+1})) = x^2(\mathbf{s}_{i+1}) + V_S(x(\mathbf{s}_i|\mathbf{s}_{i+1})). \quad (12)$$

Как уже говорилось, для субпопуляции  $\mathbf{s}_1$  первого уровня иерархии вероятность случайной встречи рассматриваемой пары индивидуумов с учетом их порядка равна  $x(\mathbf{s}_1)(1 - x(\mathbf{s}_1))$  подобно закону Харди–Вайнберга в популяционной генетике. Значение  $Hs(\mathbf{s})$  для  $\mathbf{s}$  по определению находится как среднее значение  $Hs(\mathbf{s}_1)$  при случайном выборе  $\mathbf{s}_1$  из  $\mathbf{s}$ , т.е.

$$\begin{aligned} Hs(\mathbf{s}) &= E_{s_1}\{Hs(\mathbf{s}_1)|\mathbf{s}\} = E_{s_1}\{x(\mathbf{s}_1)(1 - x(\mathbf{s}_1))|\mathbf{s}\} = \\ &= E_{s_1}\{x(\mathbf{s}_1)|\mathbf{s}\} - E_{s_1}\{x^2(\mathbf{s}_1)|\mathbf{s}\} = \\ &= x(\mathbf{s}) - (V_S(x(\mathbf{s}_1|\mathbf{s})) + (E_{s_1}\{x(\mathbf{s}_1)|\mathbf{s}\})^2) = \\ &= x(\mathbf{s})(1 - x(\mathbf{s})) - V_S(x(\mathbf{s}_1|\mathbf{s})) \end{aligned}$$

согласно (12). Очевидно, данное выражение  $Hs(\mathbf{s})$  можно переписать с использованием коэффициента фамильного инбридинга как  $x(\mathbf{s})(1 - x(\mathbf{s}))(1 - F_S(\mathbf{s}))$ , т.е. формула (11) верна. ◀

Согласно доказанному результату в метапопуляции  $\mathbf{s}$ , разбитой на несколько неподразделенных субпопуляций  $\{\mathbf{s}_1\}$ , в среднем вероятность указанной встречи в субпопуляции  $\mathbf{s}_1$  из  $\mathbf{s}$  меньше на величину межгрупповой дисперсии  $V_S(x(\mathbf{s}_1|\mathbf{s}))$ , чем вероятность независимых сочетаний фамилий двух индивидуумов во всей метапопуляции  $\mathbf{s}$ . Эта формула является практически полным аналогом эффекта Валунда в популяционной генетике (см., например, [1]). Чем больше фамильная дивергенция субпопуляций  $\{\mathbf{s}_1\}$ , тем меньше общая вероятность  $Hs(\mathbf{s})$ .

Изложенное показывает важность дисперсии распределения концентрации рассматриваемой фамилии по субпопуляциям для анализа их дивергенции. Это вызывает интерес к изучению

компонентов дисперсии концентрации. К данному случаю приложимо правило сложения дисперсий, которое в нашей ситуации выглядит следующим образом.

**Замечание 5 (правило сложения дисперсий).** Пусть метапопуляция  $\mathbf{s}$  разбита на субпопуляции  $\{\mathbf{s}_2\}$ , каждая из которых состоит из соответствующих неподразделенных групп  $\{\mathbf{s}_1\}$ , являющихся единицами наблюдения с концентрациями  $\{x(\mathbf{s}_1)\}$  рассматриваемой фамилии в них.

Тогда общая (полная) дисперсия  $V_{S_{tot}}(x(\mathbf{s}_1|\mathbf{s}))$  распределения концентрации фамилии по единицам наблюдения  $\{\mathbf{s}_1\}$  во всей метапопуляции  $\mathbf{s}$  равна сумме межгрупповой дисперсии  $V_{S_{betw}}(x(\mathbf{s}_2|\mathbf{s}))$ , характеризующей дивергенцию концентраций фамилии  $\{x(\mathbf{s}_2)\}$  между субпопуляциями  $\{\mathbf{s}_2\}$ , и средней внутригрупповой дисперсии  $W(x(\mathbf{s}_1|\mathbf{s}_2)|\mathbf{s}) \equiv E_{s_2}\{V_{S_{in}}(x(\mathbf{s}_1|\mathbf{s}_2))|\mathbf{s}\}$ , характеризующей среднюю фамильную дивергенцию концентраций  $x(\mathbf{s}_1)$  субпопуляций  $\{\mathbf{s}_1\}$  внутри отдельных субпопуляций из множества  $\{\mathbf{s}_2\}$ :

$$V_{S_{tot}}(x(\mathbf{s}_1|\mathbf{s})) = V_{S_{betw}}(x(\mathbf{s}_2|\mathbf{s})) + E_{s_2}\{V_{S_{in}}(x(\mathbf{s}_1|\mathbf{s}_2))|\mathbf{s}\}. \quad (13)$$

**Доказательство.** По определению полная (общая) дисперсия  $V_{S_{tot}}(x(\mathbf{s}_1|\mathbf{s}))$  популяции  $\mathbf{s}$  равна  $E_{s_1}\{(x(\mathbf{s}_1) - E_{s_1}\{x(\mathbf{s}_1)|\mathbf{s}\})^2\}$ . Вспомним, что  $E_{s_2}\{x(\mathbf{s}_2)|\mathbf{s}\} = x(\mathbf{s})$ ,  $E_{s_1}\{x(\mathbf{s}_1)|\mathbf{s}_i\} = x(\mathbf{s}_i)$ , т.е.  $V_{S_{tot}}(x(\mathbf{s}_1|\mathbf{s})) = E_{s_1}\{(x(\mathbf{s}_1) - x(\mathbf{s}))^2|\mathbf{s}\}$ . Представим  $(x(\mathbf{s}_1) - x(\mathbf{s}))^2$  в виде

$$\begin{aligned} &(x(\mathbf{s}_1) - x(\mathbf{s}_2) + x(\mathbf{s}_2) - x(\mathbf{s}))^2 = \\ &= (x(\mathbf{s}_1) - x(\mathbf{s}_2))^2 + (x(\mathbf{s}_2) - x(\mathbf{s}))^2 + \\ &+ 2(x(\mathbf{s}_1) - x(\mathbf{s}_2))(x(\mathbf{s}_2) - x(\mathbf{s})). \end{aligned}$$

Математическое ожидание правой части равно сумме средних значений ее отдельных слагаемых, которые находим с применением формулы полного математического ожидания (2). При этом учтем, что за знак  $E$  можно выносить независимые от усредняемой величины сомножители.

$$\begin{aligned} &E_{s_1}\{(x(\mathbf{s}_1) - x(\mathbf{s}_2))^2|\mathbf{s}\} = \\ &= E_{s_2}\{E_{s_1}\{(x(\mathbf{s}_1) - x(\mathbf{s}_2))^2|\mathbf{s}_2\}|\mathbf{s}\} = \\ &= E_{s_2}\{V_{S_{in}}(x(\mathbf{s}_1|\mathbf{s}_2))|\mathbf{s}\}, \end{aligned}$$

т.е. первое слагаемое равно средней внутригрупповой дисперсии.

Для второго слагаемого находим

$$\begin{aligned} &E_{s_1}\{(x(\mathbf{s}_2) - x(\mathbf{s}))^2|\mathbf{s}\} = \\ &= E_{s_2}\{E_{s_1}\{(x(\mathbf{s}_2) - x(\mathbf{s}))^2|\mathbf{s}_2\}|\mathbf{s}\} = \\ &= E_{s_2}\{(x(\mathbf{s}_2) - x(\mathbf{s}))^2 E_{s_1}\{1|\mathbf{s}_2\}|\mathbf{s}\} = V_{S_{betw}}(x(\mathbf{s}_2|\mathbf{s})), \end{aligned}$$

т.е. получаем межгрупповую дисперсию.

Наконец, покажем, что среднее значение третьего слагаемого равно нулю:

$$E_{s_1} \{2(x(s_1) - x(s_2))(x(s_2) - x(s)) | s\} = E_{s_2} \{2(x(s_2) - x(s)) E_{s_1} \{(x(s_1) - x(s_2)) | s_2\} | s\} = 0,$$

так как среднее отклонение от средней величины  $E_{s_1} \{(x(s_1) - x(s_2)) | s_2\}$  всегда равняется нулю.

В итоге получаем требуемое выражение для полной дисперсии

$$V_{S_{tot}}(x(s_1 | s)) = V_{S_{betw}}(x(s_2 | s)) + E_{s_2} \{V_{S_{in}}(x(s_1 | s_2)) | s\}. \blacktriangleleft$$

Данное правило приложимо к любой совокупности, разбитой на группы ее элементов, характеризующихся числовым признаком. В частном случае метапопуляции с произвольным количеством уровней иерархии общая дисперсия концентрации фамилии в ней равна сумме межгрупповой дисперсии распределения концентрации по субпопуляциям (группам) одного и того же произвольно выбранного уровня  $k$  плюс средняя внутригрупповая дисперсия концентраций в них. Скажем, в разбитой на субпопуляции  $\{s_1\}$  метапопуляции  $s = s_m$  с  $m$  уровнями иерархии правило сложения дисперсий выполняется, когда вместо субпопуляций  $\{s_2\}$  рассматриваются  $\{s_k\}$ :

$$V_{S_{tot}}(x(s_1 | s_m)) = V_{S_{betw}}(x(s_k | s_m)) + E_{s_k} \{V_{S_{in}}(x(s_1 | s_k)) | s_m\}, \quad 2 < k < m. \quad (14)$$

Выше мы считали, что встречи индивидуумов в субпопуляциях  $\{s_1\}$  низшего уровня иерархии (единицах наблюдения) случайны. Благодаря этому  $Hs(s_1) = Hs_0(s_1) = x(s_1)(1 - x(s_1))$ . В реальных исследованиях данное условие может не выполняться (скажем, при существовании скрытой подразделенности в  $s_1$  или из-за иных причин), и желательно рассмотреть такую ситуацию, поскольку для нее приведенный аналог результата Валунда неверен.

В абстрактном случае, скажем, в совокупности сталкивающихся частиц нескольких типов предположение о случайной природе их столкновений не гарантировано. Распространенным выражением для вероятности  $Hs$  неслучайных столкновений частиц рассматриваемого типа с концентрацией  $x$  с иными типами является используемая в популяционной генетике вероятность столкновений вида

$$Hs = x(1 - x)(1 - Fs) = Hs_0(1 - Fs), \quad 0 \leq Hs \leq 1, \quad 0 \leq |Fs| \leq 1.$$

При  $Fs > 0$  происходит своего рода отталкивание, и разноименные частицы сталкиваются реже, чем чисто случайно. Обратим внимание, что здесь неявно предполагается независимость  $Fs$  от типа сталкивающихся частиц (в противном случае рассматриваемая вероятность столкновения

зависит от концентраций всех типов). В случае только двух типов (рассматриваемого и искусственного объединения всех прочих типов) коэффициент  $Fs$  имеет смысл корреляции между сталкивающимися типами.

**Следствие 6.** *Рассмотрим метапопуляцию  $s_2$ , разбитую на субпопуляции  $\{s_1\}$  с концентрациями  $\{x(s_1)\}$  носителей рассматриваемой фамилии в них. Пусть вероятность  $Hs(s_1)$  их встреч в  $s_1$  с носителями иной фамилии неслучайна и формулируется как*

$$Hs(s_1) = x(s_1)(1 - x(s_1))(1 - Fs(s_1)) = Hs_0(s_1)(1 - Fs(s_1)), \quad s_2 = \{s_1\},$$

где  $Hs_0(s_1)$  обозначает вероятность  $x(s_1)(1 - x(s_1))$  чисто случайной встречи в субпопуляции  $s_1$  на первом уровне иерархии.

Тогда общая (полная) вероятность  $Hs(s_2 | Fs_1)$  рассматриваемой встречи в наугад выбранной из  $s_2$  субпопуляции  $s_1$  при независимости концентрации  $x$  рассматриваемой фамилии и коэффициента фамильного инбридинга  $Fs(s_1)$  в субпопуляциях  $\{s_1\}$  имеет вид

$$Hs(s_2 | Fs_1) = (x(s_2)(1 - x(s_2)) - V_{S_{betw}}(x(s_1 | s_2)))(1 - E_{s_1} \{Fs(s_1) | s_2\}) + V_{S_{betw}}(x(s_1 | s_2)) = V_{S_{betw}}(x(s_1 | s_2)) = \sum_{s_1} (x(s_1) - x(s_2))^2 Pr(s_1 | s_2). \quad (15)$$

Здесь  $V_{S_{betw}}(x(s_1 | s_2))$  обозначает межгрупповую дисперсию распределения концентрации  $\{x(s_1 | s_2)\}$  рассматриваемой фамилии по субпопуляциям (группам)  $\{s_1\}$ , составляющим разбиение метапопуляции  $s_2$ ;  $Pr(s_1 | s_2)$  дает вероятность случайного выбора субпопуляции  $s_1$  из метапопуляции  $s_2$ .

Иным образом общую (полную) вероятность  $Hs(s_2 | Fs_1)$  можно представить как

$$Hs(s_2 | Fs_1) = Hs_0(s_2) \left(1 - \frac{V_{S_{betw}}(x(s_1 | s_2))}{x(s_2)(1 - x(s_2))}\right) \times (1 - E_{s_1} \{Fs(s_1) | s_2\}) = Hs_0(s_2)(1 - Fs(s_2)) \times (1 - E_{s_1} \{Fs(s_1) | s_2\}), \quad (16)$$

$$Fs(s_2) \equiv \frac{V_{S_{betw}}(x(s_1 | s_2))}{x(s_2)(1 - x(s_2))} \geq 0,$$

где  $Fs(s_2)$  является фамильным аналогом случайного коэффициента инбридинга  $F_{ST}$ , а средний коэффициент фамильного инбридинга  $E_{s_1} \{Fs(s_1) | s_2\}$  для субпопуляций  $\{s_1\}$  в  $s_2$  служит аналогом коэффициента  $F_{IS}$  неслучайного инбридинга С. Райта в популяционной генетике.

**Доказательство.** Общая вероятность  $Hs(s_2)$  рассматриваемых встреч в метапопуляции  $s_2$  по определению находится как математическое ожидание

для значений  $Hs(s_1)$  в субпопуляциях  $\{s_1\}$ , составляющих разбиение  $s_2$ :

$$Hs(s_2) \equiv E_{s_1} \{Hs(s_1)|s_2\} = \\ = E_{s_1} \{x(s_1)(1-x(s_1))(1-Fs(s_1))|s_2\}.$$

При независимости  $x$  и  $Fs$  данное выражение равно произведению математических ожиданий сомножителей и

$$Hs(s_2) = E_{s_1} \{x(s_1)(1-x(s_1))|s_2\} \times \\ \times E_{s_1} \{1-Fs(s_1)|s_2\} = x(s_2)(1-x(s_2)) \times \\ \times (1-Fs(s_2))(1-E_{s_1} \{Fs(s_1)|s_2\}).$$

Здесь коэффициент фамильного инбридинга  $Fs(s_2)$  характеризует подразделенность метапопуляции  $s_2$  (фамильную дивергенцию субпопуляций  $s_1$  внутри  $s_2$ ) и по определению равен отношению дисперсии распределения концентрации фамилии по субпопуляциям  $\{s_1\}$  к вероятности случайных встреч  $Hs_0(s_2) \equiv x(s_2)(1-x(s_2))$  в  $s_2$ :

$$Fs(s_2) \equiv \frac{Vs_{betw}(x(s_1|s_2))}{x(s_2)(1-x(s_2))} = \frac{Vs_{betw}(x(s_1|s_2))}{Hs_0(s_2)}, \\ Hs_0(s_2) \equiv x(s_2)(1-x(s_2)).$$

Он соответствует коэффициенту случайного инбридинга  $F_{ST}$  в популяционной генетике, отражающему генетическую дивергенцию между субпопуляциями.

Коэффициент  $E_{s_1} \{Fs(s_1)|s_2\}$  является аналогом неслучайного коэффициента инбридинга  $F_{IS}$  в генетике популяций и отражает усредненное нарушение случайного характера встреч (закона Харди–Вайнберга) внутри субпопуляций, служащих единицей наблюдения. ◀

Повторим, что здесь неявно предполагается независимость  $Fs(s_1)$  и концентрации  $x$  рассматриваемой фамилии. Коэффициенты  $\{Fs(s_1)\}$  для субпопуляций первого уровня заданы изначально и по ним находятся средние коэффициенты неслучайного фамильного инбридинга для субпопуляций более высоких уровней. Для единообразия записи вероятности встреч метапопуляции  $s_2$  можно использовать коэффициент общего инбридинга  $F_{IT}$  в  $s_2$ . В популяционной генетике он учитывает одновременно влияние как дивергенции между субпопуляциями, так и нарушения закона Харди–Вайнберга внутри них на гетерозиготность  $H$  метапопуляции  $s_2$  (и на корреляцию между гаминами генотипа):

$$H(s_2) = x(s_2)(1-x(s_2))(1-F_{IT}), \\ 1-F_{IT} \equiv (1-F_{ST})(1-F_{IS}).$$

### РАЗЛОЖЕНИЕ ПО УРОВНЯМ ИЕРАРХИИ ВЕРОЯТНОСТИ ВСТРЕЧ ПРИ ИЕРАРХИЧЕСКОЙ ПОДРАЗДЕЛЕННОСТИ МЕТАПОПУЛЯЦИИ

Рассмотренный случай подразделенности метапопуляции на неподразделенные группы первого уровня можно считать простейшим случаем иерархической структуры. Теперь обратимся к подразделенной метапопуляции  $s_m$  с произвольным количеством уровней иерархии  $m$ . Если  $s = s_m$ , то дисперсию распределения концентрации по субпопуляциям  $\{s_k\}$  на уровне  $k$  в метапопуляции  $s_m$  обозначим как  $Vs_{betw}(x(s_k|s_m))$ ,  $k < m$ . В частности,  $Vs(x(s_1|s_m)) = Vs_{betw}(x(s_1|s_m))$ , причем  $Vs_{betw}(x(s_1|s_m))$  одновременно является общей дисперсией  $Vs_{tot}(x(s_1|s_m))$  для распределения концентрации фамилии по всем субпопуляциям, служащим единицей наблюдения.

Проанализируем выражение для вероятности встреч  $Hs$  в метапопуляции в зависимости от количества ее уровней иерархии  $m$ . Дадим решение этой задачи при независимости  $x$  и  $Fs$  путем последовательного увеличения  $m$  и индукции.

В предшествующем следствии мы нашли, что при фиксированной метапопуляции  $s = s_2$

$$Hs(s_2) \equiv E_{s_1} \{Hs(s_1)|s_2\} = \\ = (x(s_2)(1-x(s_2)) - Vs(x(s_1|s_2)))(1-E_{s_1} \{Fs(s_1)|s_2\}).$$

Теперь зафиксируем  $s_3$ , подставим в рекуррентное уравнение (8), определяющее  $Hs(s_3)$ , найденное значение  $Hs(s_2)$ , учтем эффект Валунда (11) и заменим математическое ожидание произведения произведением математических ожиданий сомножителей в силу предполагаемой независимости  $x$  и  $Fs$ :

$$Hs(s_3|Fs_1) \equiv E_{s_2} \{Hs(s_2)|s_3\} = \\ = E_{s_2} \{(x(s_2)(1-x(s_2)) - Vs(x(s_1|s_2)))|s_3\} \\ (1-E_{s_2} \{E_{s_1} \{Fs(s_1)|s_2\}|s_3\}) = \\ = (x(s_3)(1-x(s_3)) - Vs(x(s_2|s_3)) - \\ - E_{s_2} \{Vs(x(s_1|s_2))|s_3\})(1-E_{s_1} \{Fs(s_1)|s_3\}).$$

**Ремарка 7.** Заметим, что  $Vs(x(s_2|s_3))$  можно записать как  $E_{s_1} \{Vs(x(s_2|s_3))|s_3\}$ , поскольку метапопуляция  $s_3$  единственна. Поэтому

$$E_{s_2} \{Hs(x(s_2|s_3))\} = \left( x(s_3)(1-x(s_3)) - \sum_{i=1}^2 E_{s_{i+1}} \{Vs(x(s_i|s_{i+1}))|s_3\} \right) (1-E_{s_1} \{Fs(x(s_1)|s_3\}).$$

Найденное выражение для  $Hs(s_3)$  и ремарка подсказывают, что для субпопуляций  $j$ -го уровня

$$Hs(s_j) = \left( x(s_j)(1 - x(s_j)) - \sum_{i=1}^{j-1} E_{s_{i+1}} \{Vs(x(s_i|s_{i+1})|s_j)\} \right) (1 - E_{s_1} \{Fs(s_1|s_j)\}).$$

С другой стороны, поскольку  $Hs(s_j) = (x(s_j)(1 - x(s_j)) - Vs(x(s_1|s_j)))(1 - E_{s_1} \{Fs(s_1|s_j)\})$  согласно фамильному аналогу эффекта Валунда (11), где роль метапопуляции  $s$  играет  $s_j$ , то одновременно мы видим разложение дисперсии  $Vs(x(s_1|s_j))$  на сумму членов  $E_{s_{i+1}} \{Vs(x(s_i|s_{i+1})|s_j)\}$ .

**Результат 8.** Рассмотрим иерархически подразделенную метапопуляцию  $s_m$  с уровнями иерархии  $i = 1, 2, \dots, m$ . Положим, что в ее субпопуляциях  $\{s_i\}$  концентрации интересующей фамилии равны  $\{x(s_i)\}$  и встречи индивидуумов происходят случайно с вероятностями  $\{Hs(s_i) = Hs_0(s_i) \equiv x(s_i)(1 - x(s_i))\}$ .

Пусть из  $s_m$  наугад выбрана субпопуляция  $s_1$  первого уровня. Тогда вероятность  $Hs(s_m)$  встречи в  $s_1$  пары индивидуумов с интересующей фамилией и с какой-нибудь иной с учетом их порядка (полная вероятность встречи) допускает следующее разложение:

$$\begin{aligned} Hs(s_m) &\equiv E_{s_1} \{Hs(s_1)|s_m\} = \\ &= x(s_m)(1 - x(s_m)) - Vs_{betw}(x(s_1|s_m)) = \\ &= Hs_0(s_m) - \sum_{i=1}^{m-1} E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\}. \end{aligned} \tag{17}$$

Таким образом, общая (полная) вероятность  $Hs(s_m)$  равна разности между значением вероятности  $Hs_0(s_m) \equiv x(s_m)(1 - x(s_m))$  случайной встречи такой пары во всей метапопуляции  $s_m$  и разложением межгрупповой дисперсии  $Vs_{betw}(x(s_1|s_m))$  на средние внутригрупповые дисперсии  $E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\}$ , соответствующие отдельным уровням иерархии  $i$ .

**Доказательство.** Воспользуемся индукцией. Выше мы видели, что разложение (17) для  $Hs(s_m)$  верно при некотором  $m$ . Покажем, что тогда оно справедливо для  $Hs(s_{m+1})$ . Согласно рекуррентному уравнению (8)

$$\begin{aligned} Hs(s_{m+1}) &\equiv E_{s_m} \{Hs(s_m)|s_{m+1}\} = \\ &= E_{s_m} \{x(s_m)(1 - x(s_m)) - \\ &- \sum_{i=1}^{m-1} E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\}|s_{m+1}\} = \\ &= E_{s_m} \{x(s_m)(1 - x(s_m))\} - \\ &- E_{s_m} \left\{ \sum_{i=1}^{m-1} E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\}|s_{m+1} \right\} = \\ &= x(s_{m+1})(1 - x(s_{m+1})) - Vs(x(s_m|s_{m+1})) - \\ &- \sum_{i=1}^{m-1} E_{s_m} \{E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\}|s_{m+1}\}. \end{aligned}$$

Рассмотрим отдельное слагаемое в сумме по  $i$ . Согласно (6), где роль  $i$  играет  $i + 1$ , а  $T(s_{i+1}) = Vs_{in}(x(s_i|s_{i+1}))$ , имеем

$$\begin{aligned} E_{s_m} \{E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\}|s_{m+1}\} &= \\ &= E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_{m+1})\}. \end{aligned}$$

Поскольку в соответствии со сделанной ремаркой  $Vs_{in}(x(s_m|s_{m+1})) = E_{s_{m+1}} \{Vs_{in}(x(s_m|s_{m+1})|s_{m+1})\}$ , этот член можно включить в рассматриваемую сумму, верхний предел которой увеличится на единицу и станет равен  $(m - 1) + 1 = m$ :

$$\begin{aligned} Hs(s_{m+1}) &\equiv E_{s_m} \{Hs(s_m)|s_{m+1}\} = x(s_{m+1})(1 - x(s_{m+1})) - \\ &- \sum_{i=1}^m E_{s_{i+1}} \{Vs(x(s_i|s_{i+1})|s_{m+1})\}. \end{aligned}$$

Полученное выражение совпадает с (17), если вместо  $m$  поставить  $m + 1$ . ◀

Из (17) отчетливо видно, что вероятность  $Hs$  для уровня выше первого меньше значения, соответствующего случайным встречам, и чем больше сумма средних внутригрупповых дисперсий, тем больше увеличивается указанный разрыв. Полученное разложение является некоторого рода обобщением эффекта Валунда [10]. Заметим, что здесь не требуется независимости концентрации фамилии и коэффициента фамильного инбридинга.

**Следствие 9.** В иерархически подразделенной метапопуляции  $s_m$  с  $m$  уровнями иерархии полная дисперсия концентрации фамилии  $Vs_{tot}(x(s_1|s_m))$  допускает разложение

$$\begin{aligned} Vs_{tot}(x(s_1|s_m)) &= \sum_{i=1}^{m-1} E_{s_{i+1}} \{Vs(x(s_i|s_{i+1})|s_m)\} = \\ &= \sum_{i=1}^{m-1} (Vs(x(s_i|s_m)) - Vs(x(s_{i+1}|s_m))). \end{aligned} \tag{18}$$

**Доказательство** первого равенства содержится в предыдущем результате, т.е. там попутно получено разложение дисперсии, найденное нами ранее в [4] более длинным способом.

Второе равенство вытекает из правила сложения дисперсий (14), согласно которому  $E_{s_k} \{Vs_{in}(x(s_n|s_k)|s_m)\} = Vs_{tot}(x(s_n|s_m)) - Vs_{betw}(x(s_k|s_m)) \geq 0$ . В нашем случае при  $n = i, k = i + 1$

$$\begin{aligned} E_{s_{i+1}} \{Vs_{in}(x(s_i|s_{i+1})|s_m)\} &= Vs(x(s_i|s_m)) - \\ &- Vs(x(s_{i+1}|s_m)), \quad n < k < m. \end{aligned}$$

Таким образом, полная дисперсия  $Vs(x(s_1|s_m))$  равна сумме приращений межгрупповой дисперсии при переходе от уровня  $i + 1$  к предыдущему  $i$ .

Подобно тому, как пройденный путь складывается из приращений по шагам, так и при суммировании приращений дисперсии по  $i$  получим полную дисперсию  $Vs(x(s_i|s_m))$ , поскольку все члены в сумме сокращаются, кроме  $Vs(x(s_1|s_m))$ , а  $Vs(x(s_m|s_m))$  разумно определить как нуль. ◀

Так как левая часть  $E_{s_{i+1}}\{Vs(x(s_i|s_{i+1}))|s_m\}$  последней формулы неотрицательна (как среднее значение неотрицательной дисперсии), то при переходе к более высокому уровню  $i + 1$  межгрупповая дисперсия (дисперсия концентрации фамилии в субпопуляциях одного и того же уровня) не увеличивается. Это интуитивно ожидается, поскольку при группировке данных разброс значений не возрастает, и дисперсия средних значений для сгруппированных данных не может быть больше дисперсии для исходных данных.

### ЗАКЛЮЧЕНИЕ

Пройдемся по основным линиям данной работы — ее мотивации, подходу к анализу, результатам. Интерес к семейной структуре мотивирован вниманием к генеалогии, родственным связям и др. аспектам, так или иначе связанным с генетикой населения. Настоящая статья ориентирована на популяционно-генетическую сторону изучения распределения фамилий, и эта ориентация приводит к использованию и обобщениям традиционных методов анализа генетической структуры по данным о распределении концентраций маркеров в ряде популяций.

Выбор такого подхода обусловлен тем, что популярные модели процессов выборочного генного дрейфа Райта—Фишера и дрейфа фамилий в популяциях ограниченной численности при стандартных предположениях совпадают (отличаются лишь интенсивностью, см. [2, 3]). Следовательно, это должно позволить при соответствующих поправках экстраполировать результаты одного из них на ожидаемые результаты для другого. Природа процессов дрейфа с динамикой генетического и семейного состава популяции обусловлена случайными колебаниями количества потомков у отдельных индивидуумов.

Если у некоторого родителя-мужчины больше сыновей, чем в среднем, то в следующем поколении будет больше носителей его фамилии и его генов. На популяционном уровне это приводит к случайным колебаниям семейного и генетического составов популяции от поколения к поколению, представляющим сущность процессов дрейфа. Оба процесса приводят к дивергенции родственных популяций (генетической и семейной). Очевидно, требуется развивать методы, характеризующие дивергенцию популяций (неоднородность системы популяций в целом) для типичных для человека популяционных структур.

Данные по фамилиям часто доступны в иерархически структурированном виде, например в соответствии с административно-территориальной подразделенностью. В качестве показателей дивергенции служат межгрупповая дисперсия распределения концентрации (кодминантного аллеля, фамилии) в группах (субпопуляциях), нарушения закона Харди—Вайнберга с эффектом Валунда, коэффициент случайного инбридинга. Рассмотрены семейные аналоги этих показателей и даны их обобщения для иерархически подразделенной метапопуляции с произвольным количеством уровней иерархии. Данные обобщения вычислительно универсальны в том смысле, что не зависят от характера разбиений метапопуляции на каждом из уровней, от типа микроэволюционного процесса, приведшего к текущему состоянию метапопуляции, от существования миграций и пр. Неоднородность (дивергенция) в иерархических системах характеризуется алгоритмически одинаковыми показателями, различающимися между системами количественно.

В настоящей статье в качестве основного показателя семейной дивергенции субпопуляций рассматривается общая вероятность встреч (столкновений)  $Hs$  лиц с разными фамилиями (аналог концентрации гетерозигот), понимаемая как вероятность встречи в выбранной наугад субпопуляции, единицы наблюдения, из рассматриваемой метапопуляции с иерархической структурой субпопуляций. Получено разложение  $Hs$  по уровням иерархии, обобщающее эффект Валунда в популяционной генетике. Общая вероятность неизонимных встреч в иерархически подразделенной метапопуляции меньше вероятности случайных встреч в ней на сумму средних внутригрупповых дисперсий концентрации фамилии, соответствующих отдельным уровням.

Настоящая статья не содержит каких-либо исследований с использованием в качестве объекта животных.

Настоящая статья не содержит каких-либо исследований с участием в качестве объекта людей.

### СПИСОК ЛИТЕРАТУРЫ

1. Ли Ч. Введение в популяционную генетику. М.: Мир, 1978. 555 с.
2. Пасеков В.П. К анализу случайных процессов изонимии. I. Структура изонимии // Генетика. 2021. Т. 57. № 10. С. 1194–1204. <https://doi.org/10.31857/S001667582110009X>
3. Пасеков В.П. К анализу случайных процессов изонимии. II. Динамика дивергенции популяций // Генетика. 2021. Т. 57. № 11. С. 1318–1329. <https://doi.org/10.31857/S0016675821101114>
4. Пасеков В.П. Описание дивергенции субпопуляций в иерархической системе при анализе изонимии

- мии. I. Дисперсия как показатель дивергенции // Генетика. 2022. Т. 58. № 6. С. 713–727.
5. *Wright S.* The interpretation of population structure by *F* statistics with special regard to systems of mating // *Evolution*. 1965. V. 19. P. 395–420.
  6. *Гинтер Е.К., Зинченко Р.А., Ельчинова Г.И. и др.* Роль факторов популяционной динамики в распространении наследственной патологии в российских популяциях // *Мед. генетика*. 2004. Т. 3. № 12. С. 548–555.
  7. *Ревазов А.А., Парадеева Г.М., Русакова Г.И.* Пригодность русских фамилий в качестве квазигенетического маркера // *Генетика*. 1986. Т. 22. № 4. С. 699–703.
  8. *Crow J.F., Mange A.P.* Measurement of inbreeding from the frequency of marriages between persons of the same surname // *Social Biology*. 1982. V. 29. № 1/2. P. 101–105.
  9. *Lasker W.G.* Surnames and Genetic Structure. Cambridge: Cambr. Univ. Press, 1985. 148 p.
  10. *Сорокина И.Н., Чурносов М.И., Балтуцкая И.В. и др.* Антропогенетическое изучение населения центральной России. М.: Изд-во РАМН, 2014. 336 с.

## Description of Divergence of Subpopulations in the Hierarchical System under the Analysis of Isonymy. II. Probabilities of Non-Isonymic Encounters

V. P. Passekov<sup>a, \*</sup>

<sup>a</sup>*Federal Research Center “Computer Science and Control” of Russian Academy of Sciences, Moscow, 119991 Russia*

<sup>\*</sup>*e-mail: pass40@mail.ru*

Metapopulations with typical for human populations hierarchical subdivision into parts (subpopulations) are considered, corresponding to classifications based on administrative-territorial division (say, village, village council, district, region, and so on) or on a genealogical approach based on ethnogenesis, as well as on other principles of biological classification. The purpose of this work is to analyze the general properties of the distribution of the surname concentration over subpopulations in their hierarchical structure. Attention is focused on the description of the surname divergence of subpopulations, as an indicator of which the total probability of two person with different surnames encounter is considered, which is understood as the probability of the encounter in a randomly chosen subpopulation, a unit of observation, from metapopulation with a hierarchical structure. Its decomposition by hierarchy levels is obtained, which generalizes the Wahlund effect in population genetics. The total probability of non-isonymic encounters in a hierarchically subdivided metapopulation is less than the probability of random encounters in it by the sum of the average surname concentration intragroup variances corresponding to separate levels. Such properties are purely statistical characteristics of the hierarchical structure, and not a feature of a particular population system, and are not derived from the regularities of one or another model of microevolution. They are computationally formulated in the same way for any hierarchical system, although in the general case they do not coincide quantitatively. The results obtained refer to rural and urban hierarchical metapopulations as separate components of the entire population.

**Keywords:** hierarchical structure of populations, metapopulations, surname concentrations in human subpopulations, characteristics of heterogeneity of subpopulations, decomposition of the probability of encounter carriers of different surnames by hierarchy levels.