

МОДЕЛИ МНОГОМЕРНЫХ ПРОЦЕССОВ С ГЕТЕРОГЕННОСТЬЮ

3.1 Модель коррелированной уязвимости

Описание рассмотренного ранее процесса выхода из строя элементов в неоднородной популяции требует обобщения на случай наблюдения нескольких характеристик. Наиболее интересным является изучение поведения не независимых, как это обычно делается, элементов, а элементов, свойства которых статистически связаны. Примером такого обобщения в биометрике является изучение продолжительности жизни и состояния здоровья родителей и детей, братьев и сестёр. В последнее время в связи с появлением результатов крупных популяционных исследований особый интерес представляет изучение пар монозиготных и дизиготных близнецов с целью выявить роль генов и условий окружающей среды в изменении состояния здоровья и рисков, опасных для жизни человека [23, 24].

Рассмотрим модель, в которой индивидуальные уязвимости двух человек коррелированы [25, 26]. Будем считать, что для каждого человека справедлива модель пропорционального риска с гамма распределённой уязвимостью, причем при фиксированных величинах уязвимости продолжительности их жизней независимы. Условное распределение продолжительностей жизни равно

$$\begin{aligned} S\{T_1 > t_1, T_2 > t_2 \mid z_1, z_2\} &= S\{T_1 > t_1 \mid z_1\}S\{T_2 > t_2 \mid z_2\} \\ &= \exp(-z_1 H_{01}(t_1) - z_2 H_{02}(t_2)), \end{aligned} \quad (3.1)$$

где T_1, T_2 - случайные величины, соответствующие длительности жизни первого и второго человека соответственно, $H_{01}(t)$ и $H_{02}(t)$ обозначают соответствующие базовые кумулятивные риски, нормированные так, чтобы среднее значение уязвимости каждого человека было равно единице. Относительно индивидуальных уязвимостей Z_1 и Z_2 сделаем предположение, что они представимы в виде $Z_i = \sum_{j=1}^n (Y_{0j} + Y_{ij})$, $i = 1, 2$, где величины Y_{0j}, Y_{1j} и Y_{2j} , $j = 1, \dots, n$ являются независимыми гамма распределёнными случайными величинами с параметрами k_{0j}, k_{1j}, k_{2j} и единым параметром l . Таким образом, случайные величины Z_1 и Z_2 имеют гамма распределение с тем же параметром l и параметром $k = l$ поскольку среднее значение уязвимости равно единице. Удобно предположить, что и $k_{1j} = k_{2j}$. Коэффициент корреляции между уязвимостями Z_1 и Z_2 вычисляется по формуле

$$\begin{aligned}
 r_z &= \frac{1}{s_z^2} E((Z_1 - E Z_1)(Z_2 - E Z_2)) \\
 &= \frac{1}{s_z^2} E\left(\sum_{i=1}^n \sum_{j=1}^n (Z_{0i} - E Z_{0i} + Z_{1i} - E Z_{1i})(Z_{0j} - E Z_{0j} + Z_{2j} - E Z_{2j})\right) \\
 &= \frac{1}{s_z^2} \sum_{j=1}^n s_{0j}^2
 \end{aligned}$$

где s_z^2 - дисперсия уязвимости, s_{0j}^2 - дисперсия j -того члена в

представлении уязвимости. Формулу для коэффициента корреляции можно переписать в виде

$$r_z = \sum_{j=1}^n r_j g_j^2,$$

где $g_j^2 = (s_{0j}^2 + s_{1j}^2) / s_z^2$ есть отношение дисперсии компоненты Y_{0j} к дисперсии уязвимости, $r_j = s_{0j}^2 / (s_{0j}^2 + s_{1j}^2)$ является коэффициентом корреляции j -тых членов в представлении уязвимостей Z_1 и Z_2 . Параметры гамма распределений связаны с соответствующими дисперсиями и коэффициентами корреляции соотношениями: $I = s_z^{-2}$, $k_{0j} = r_j (g_j / s_z)^2$, $k_{1j} = (1 - r_j) (g_j / s_x)^2$.

Усредняя двумерную функцию дожития (3.1) по гамма распределённым коррелированным уязвимостям получим двумерное распределение длительности жизни пары из гетерогенной популяции, обладающей различными кумулятивными рисками

$$\begin{aligned} S(t_1, t_2) &= \left(1 + s_z^2 (H_{01}(t_1) + H_{02}(t_2))\right)^{-s_z^{-2} \sum_{j=1}^n r_j g_j^2} \\ &\quad \times \left(1 + s_z^2 H_{01}(t_1)\right)^{-s_z^{-2} \sum_{j=1}^n (1-r_j) g_j^2} \\ &\quad \times \left(1 + s_z^2 H_{02}(t_2)\right)^{-s_z^{-2} \sum_{j=1}^n (1-r_j) g_j^2}. \end{aligned}$$

Принимая во внимание выражение для вероятности безотказной работы элемента в популяции с гамма распределённой гетерогенностью, полученное в предыдущей главе, запишем соотношение

$H_{oi}(t) = (S_i(t)^{s_z^2} - 1) / s_z^2$, $i=1,2$ из которого следует представление двумерной функции дожития через одномерные функции дожития

$$S(t_1, t_2) = \frac{S_1(t_1)^{1-r_z} S_2(t_2)^{1-r_z}}{\left(S_1(t_1)^{-s_z^2} + S_2(t_2)^{-s_z^2} - 1 \right)^{r_z/s_z^2}}.$$

Коэффициент корреляции r_z представляет собой комбинацию коэффициентов корреляции отдельных компонент уязвимости

$$r_z = \sum_{j=1}^n r_j g_j^2.$$

В практических приложениях эти коэффициенты

корреляции могут быть известны экспериментально или теоретически. Например, при анализе дожития близнецов возможно рассмотреть компоненты уязвимости, связанные с аддитивным генетическим эффектом, доминантным генетическим эффектом, эпистатическим генетическим эффектом, эффектами общего и индивидуального окружения [25]. Из постулатов количественной генетики следуют следующие значения коэффициентов корреляции для монозиготных и дизиготных близнецов [27]

эффект	аддитивный	доминантный	эпистатический	общее окружение	индивидуальное окружение
монозиготные	1.0	1.0	1.0	1.0	0.0
дизиготные	0.5	0.25	x	1.0	0.0

Используя эту информацию остальные неизвестные параметры модели определяются из экспериментальных данных.

3.2 LCM - модель латентных классов

Время безотказной работы или жизни является примером непрерывной характеристики, изучаемой в неоднородной популяции. Часто интерес представляют категориальные переменные, выражающие признаки наличия заболевания, состояния здоровья, выраженные в некоторой условной шкале и так далее. При этом удобно представлять результаты единичных наблюдений в виде M – мерных дискретных величин. Каждая координата соответствует градациям своей номинальной шкалы, совокупность всех градаций образует M – мерную таблицу сопряженности. При анализе подобных данных основной интерес представляет оценка и содержательная интерпретация вероятностей попадания наблюдений в отдельные ячейки таблицы сопряженности. В этой главе опишем модель латентных классов (LCM модель), применяемую для анализа дискретной ненаблюдаемой неоднородности [18].

Сначала допустим, что наблюдаются две дискретные переменные с I и J градациями. Результаты наблюдений обычно представляются двумерной таблицей сопряженности, состоящей из $I \times J$ ячеек. Обозначим через p_{ij} – вероятность попадания наблюдения в соответствующую ячейку. Предположим, что для некоторого K существуют векторы $\Theta_a > 0$, $p_i^{(a)} > 0$, $q_j^{(a)} > 0$ ($\alpha=1, \dots, K$; $i=1, \dots, I$; $j=1, \dots, J$)

такие, что $\sum_{a=1}^K \Theta_a = 1$, $\sum_{i=1}^I p_i^{(a)} = 1$, $\sum_{j=1}^J q_j^{(a)} = 1$ и справедливо представление

$$p_{ij} = \sum_{a=1}^K \Theta_a p_i^{(a)} q_j^{(a)} .$$

Такая запись означает, что существуют K ненаблюдаемых классов, в каждом из которых строки и столбцы таблицы сопряженности независимы (условная независимость). Каждый эксперимент, результат которого наблюдается, соответствует одному классу, вектор Θ_α задаёт априорные вероятности принадлежности эксперимента к различным классам, произведения $p_i^{(a)} q_j^{(a)}$ задают условные вероятности исхода, соответствующего ячейке ij в различных классах.

Если таблица сопряженности соответствует не двум, а M переменным, то вероятность попадания наблюдения в конкретную ячейку в рамках LCM модели может быть записана в виде

$$p_{i_1 \dots i_M} = \sum_{a=1}^K \Theta_a \prod_{j=1}^M p_{i_j}^{(a)}$$

где $p_{i_j}^{(a)}$ - вероятность попадания наблюдения по переменной j в ячейку i_j при условии, что наблюдение принадлежит классу α . Оценок параметров в модели латентных классов рассмотрены в главе 4.

3.3 GOM - модель степени принадлежности

Обобщением модели латентных классов служит модель, использующая идеи размытых множеств. При конструировании такой модели предполагается, что условия проведения каждого эксперимента соответствуют не одному, а одновременно нескольким клас-

сам. Степень принадлежности эксперимента конкретному классу определяется функцией принадлежности, которая предполагается уникальной для каждого эксперимента. При этом, удаётся избежать противоречий между реальностью и гипотезой о принадлежности эксперимента гипотетическому “чистому классу”, который может быть совершенно условным и не существовать в природе.

При использовании статистических моделей, основанных на размытых множествах, совокупность наблюдений интерпретируется как наблюдения в гетерогенной статистической совокупности, а функция принадлежности конкретного наблюдения различным классам описывает неоднородность соответствующего эксперимента или объекта. Совокупность функций принадлежности, соответствующих всем изучаемым объектам, характеризует неоднородность всей совокупности объектов или популяции. Функции принадлежности отдельных объектов характеризуют дополнительно индивидуальную гетерогенность, присущую конкретному эксперименту или индивидууму, в отличие от популяционной гетерогенности, присущей совокупности объектов – популяции.

Независимые наблюдения

Модели, основанные на размытых множествах, были введены в 70х годах для анализа категориальных данных и получили название GOM-модели (от grades of membership) [19, 20]. Следуя принятым обозначениям представим каждое наблюдение в виде целочисленного вектор с J координатами, каждая из которых принимает L_j зна-

чений. Это могут быть дискретизованные значения непрерывной величины, либо кодировки номинальных величин как при обработке анкетной информации. В последнем случае каждую координату можно понимать как “вопрос”, а значение координаты как конкретный “ответ” на это вопрос. Пусть имеется K базовых классов. Функцию принадлежности эксперимента i базовому классу k обозначим через g_{ik} . Вероятность того, что в эксперименте, принадлежащем базовому классу k , переменная j примет значение l , обозначим через I_{kjl} . За этими вероятностями закрепилось название *профильные вероятности*. На функции g_{ik} и I_{kjl} накладываются естественные ограничения

$$g_{ik} \geq 0,$$

$$\sum_{k=1}^K g_{ik} = 1, \quad i = 1, \dots, I,$$

$$I_{kjl} \geq 0,$$

$$\sum_{l=1}^{L_j} I_{kjl} = 1, \quad k = 1, \dots, K, \quad j = 1, \dots, J.$$

Применение размытых множеств в статистике основано на двух допущениях первое из которых утверждает, что вероятность исхода конкретного эксперимента является линейной комбинацией профильных вероятностей, то есть справедливо соотношение

$$P\{Y_{ijl} = 1\} = \sum_{k=1}^K g_{ik} I_{kjl},$$

где Y_{ijl} – случайная величина, принимающую значение 1 если в эксперименте i переменная j примет значение l , и значение 0 в противном случае. Второе условие имеет смысл условной независимости и означает, что для фиксированного i вероятность одновременного наблюдения значений различных переменных равна произведению вероятностей наблюдать каждое значение в отдельности

$$P\{Y_{ij_1 l_1} = 1, \dots, Y_{ij_J l_J} = 1\} = \left(\sum_{k=1}^K g_{ik} I_{k j_1 l_1} \right) \mathbf{L} \left(\sum_{k=1}^K g_{ik} I_{k j_J l_J} \right).$$

С учётом этих условий и обозначив реализацию случайной величины через y_{ijl} запишем вероятность наблюдения конкретной реализации значений всех переменных при проведении I независимых экспериментов – функцию правдоподобия

$$L = \prod_{i=1}^I \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K g_{ik} I_{kjl} \right)^{y_{ijl}}.$$

Строго говоря, эта функция является функцией правдоподобия только если индивидуальные функции принадлежности g_{ik} заданы априори. По этой причине её называют условной функцией правдоподобия в отличие от безусловной функции правдоподобия

$$L_u = \int \prod_{i=1}^I \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K I_{ijk} \mathbf{x}_k \right)^{y_{ijl}} dP(\mathbf{x}),$$

которая получается путём усреднения условной функции правдоподобия по случайной k -мерной вектор- функции принадлежности \mathbf{x} , имеющей распределение, удовлетворяющее ограничению

$x_k \geq 0, \sum_{k=1}^K x_k = 1$. Индивидуальные функции принадлежности g_{ik}

рассматриваются при этом как независимые реализации случайного вектора x .

Зависимые наблюдения

При конструировании условной функции правдоподобия предполагалось, что результаты исходов различных экспериментов независимы. При анализе данных о состоянии индивидуумов, связанных родственно и генетически, наблюдения оказываются зависимыми и большой интерес представляет изучение степени этой связи в зависимости от генетических и родственных факторов. Один из возможных подходов к применению моделей, основанных на размытых множествах, для анализа данных о родственниках использован в [27] для поиска генов, связанных с наследственным диабетом. При этом в качестве исхода “эксперимента” принималась информация о двух родственниках, сами же пары считались независимы друг от друга. Информация о неполных парах не использовалась, что является существенным ограничением применения этого подхода.

Другая возможность заключается в непосредственном учёте структуры зависимости между наблюдаемыми переменными. Пусть x^1 и x^2 обозначают функции принадлежности для двух статистически зависимых индивидуумов. Безусловная функция правдоподобия для такой пары имеет вид

$$L_u = \int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K x_k^1 I_{kjl} \right)^{y_{1jl}} \left(\sum_{k=1}^K x_k^2 I_{kjl} \right)^{y_{2jl}} dP(\mathbf{x}^1, \mathbf{x}^2).$$

Допустим, что существуют две такие случайные ненаблюдаемые величины T_1 и T_2 , называемые liability, такие, что совместное распределение вектор- функций принадлежности \mathbf{x}^1 и \mathbf{x}^2 можно представить в виде

$$P(\mathbf{x}^1, \mathbf{x}^2) = \int F(\mathbf{x}^1 | t_1) F(\mathbf{x}^2 | t_2) h(t_1, t_2) dt_1 dt_2,$$

где $F(\mathbf{x} | t)$ - условное многомерное распределение функции принадлежности при фиксированной величине ненаблюдаемой liability.

Безусловная функция правдоподобия примет вид

$$L_u = \int \int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K x_k^1 I_{kjl} \right)^{y_{1jl}} \left(\sum_{k=1}^K x_k^2 I_{kjl} \right)^{y_{2jl}} f(\mathbf{x}^1 | t_1) f(\mathbf{x}^2 | t_2) h(t_1, t_2) dt_1 dt_2 d\mathbf{x}^1 d\mathbf{x}^2.$$

Перегруппировав, запишем

$$\begin{aligned} L_u &= \int \left(\int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K x_k^1 I_{kjl} \right)^{y_{1jl}} \left(\sum_{k=1}^K x_k^2 I_{kjl} \right)^{y_{2jl}} f(\mathbf{x}^1 | t_1) f(\mathbf{x}^2 | t_2) d\mathbf{x}^1 d\mathbf{x}^2 \right) h(t_1, t_2) dt_1 dt_2 \\ &= \int \left(\int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K x_k^1 I_{kjl} \right)^{y_{1jl}} f(\mathbf{x}^1 | t_1) d\mathbf{x}^1 \right) \left(\int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K x_k^2 I_{kjl} \right)^{y_{2jl}} f(\mathbf{x}^2 | t_2) d\mathbf{x}^2 \right) h(t_1, t_2) dt_1 dt_2. \end{aligned}$$

Функция распределения величин liability $h(t_1, t_2)$ отражает зависимость между близнецами.

Следуя моделям количественной генетики [27] представим случайные величины T_1 и T_2 в виде разложения на слагаемые, ответственные за аддитивную генетическую компоненту – А, общее окружение – С, и индивидуальное окружение – Е, то есть

$T_1=A_1+C_1+E_1$, $T_2=A_2+C_2+E_2$. Возможны и более сложные представления, дополнительно учитывающие, например, доминантный и эпистатический генетические эффекты.

Компоненты А, В и С являются результатом взаимодействия большого числа генов и микро событий, в результате чего каждая из компонент может рассматриваться как нормально распределённая случайная величина. Считая компоненты С, Е и А независимыми, а компоненты A_1 и A_2 коррелированными между собой с коэффициентом корреляции r , получим, что величины T_1 и T_2 распределена по нормальному закону с корреляционной матрицей

$$\text{cov}(T_1, T_2) = \begin{pmatrix} s_a^2 + s_c^2 + s_e^2 & r * s_a^2 + s_c^2 \\ r * s_a^2 + s_c^2 & s_a^2 + s_c^2 + s_e^2 \end{pmatrix}$$

где s_a^2 , s_c^2 и s_e^2 обозначают соответственно дисперсии аддитивной генетической компоненты, общего и индивидуального окружения. Эти дисперсии считаются одинаковыми для обоих близнецов. Поскольку монозиготные близнецы имеют полностью общий хромосомный набор, а у дизиготных близнецов лишь половина хромосом общая, то коэффициент корреляции r равен единице для монозиготных близнецов и 0.5 для дизиготных [27].

Для описания условного распределения функции принадлежности можно использовать представление $x_i = \Phi(t, a_i)$ ($i=1, \dots, K-1$), где $\Phi(t, a)$ - функция распределения нормального закона со средним значением a и единичной дисперсией. При ограничениях

$x_j \geq 0$; $\sum_{j=1}^K x_j = 1$ справедливо выражение

$$F(x|t) = P\{x_1 < x_1, \dots, x_K < x_K\} = \Phi\left(\min_{j=1, \dots, K} \Phi^{-1}(x_j, a_j)\right).$$

Величины a_j являются параметрами модели, подлежащими определению по данным, причем важны не их абсолютные значения, а расстояние до среднего значения величины liability, которое, следовательно, можно принять равным нулю.

Функция правдоподобия принимает вид

$$L = L_S * L_{MZ} * L_{DZ},$$

где L_S обозначает функцию условного правдоподобия для близнецов, представленных без брата или сестры

$$L_S = \prod_{i \in I(S)} \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K g_{ik} I_{ijl} \right)^{y_{ijl}}.$$

Безусловные функции правдоподобия L_{MZ} и L_{DZ} для монозиготных и дизиготных пар близнецов определяются выражением

$$L_{zvg} = \prod_{i \in I(zvg)} \int \left(\int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K x_k I_{kjl} \right)^{y_{1jil}} f(x|t_1) dx \right) \times \left(\int \prod_{j=1}^J \prod_{l=1}^{L_j} \left(\sum_{k=1}^K y_k I_{kjl} \right)^{y_{2jil}} f(y|t_2) dy \right) N(t_1, t_2, \text{cov}_{zvg}) dt_1 dt_2$$

где $N(t_1, t_2, \text{cov})$ - двумерная плотность вероятности нормального распределения с нулевым средним и ковариационной матрицей cov.