

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ ФЕДЕРАЦИИ  
МОСКОВСКИЙ ФИЗИКО-ТЕХНИЧЕСКИЙ ИНСТИТУТ (государственный университет)  
ФАКУЛЬТЕТ УПРАВЛЕНИЯ И ПРИКЛАДНОЙ МАТЕМАТИКИ  
ВЫЧИСЛИТЕЛЬНЫЙ ЦЕНТР им. А. А. ДОРОВНИЦЫНА РАН  
КАФЕДРА «ИНТЕЛЛЕКТУАЛЬНЫЕ СИСТЕМЫ»

Кузнецов Михаил Павлович

**Построение интегральных индикаторов  
в задачах с порядковыми признаками**

511656 - Математические и информационные технологии

МАГИСТЕРСКАЯ ДИССЕРТАЦИЯ

**Научный руководитель:**  
н.с. ВЦ РАН, к.ф.-м.н.  
Стрижов Вадим Викторович

Москва

2013

# Содержание

<b>1</b>	<b>Введение</b>	<b>4</b>
<b>2</b>	<b>Постановка задачи построения интегральных индикаторов</b>	<b>6</b>
2.1	Модель построения интегральных индикаторов . . . . .	6
2.2	Порядковые шкалы: частично упорядоченные множества . . . . .	6
<b>3</b>	<b>Описание порядковой шкалы конусом</b>	<b>7</b>
3.1	Соответствие полиэдрального конуса порядковому признаку . . . . .	7
3.2	Нахождение оптимальных параметров модели в случае порядковых признаков . . . . .	9
3.3	Непараметрическое описание конусов, соответствующих порядковым признакам . . . . .	11
<b>4</b>	<b>Альтернативные методы</b>	<b>14</b>
4.1	Построение интегрального индикатора в порядковых шкалах с использованием копул . . . . .	14
4.2	Алгоритм криволинейной регрессии . . . . .	21
<b>5</b>	<b>Вычислительные эксперимент</b>	<b>23</b>
5.1	Реальные данные: категоризация редких видов Красной книги РФ . . .	23
5.2	Сходимость параметров . . . . .	24
<b>6</b>	<b>Заключение</b>	<b>26</b>

### Аннотация

Рассмотрена задача построения интегральных индикаторов объектов, описанных в смешанных (линейных, номинальных, порядковых) шкалах. Множеству значений признака ставится в соответствие подмножество евклидоваго пространства специального вида. В случае порядковой шкалы, множеству значений признака ставится в соответствие конус, описывающий возможные значения этого признака. На основе параметризации этих подмножеств строится новое признаковое описание, и выполняется процедура построения интегральных индикаторов, заключающаяся в поиске ближайших векторов в конусах. Рассматриваются случаи заданного экспертами вектора параметров, значения которого принадлежат линейной и порядковой шкале. Алгоритм используется для категоризации редких видов Красной книги РФ. Производится сравнение результатов алгоритма с двумя альтернативными подходами: статистическим методом и алгоритмом криволинейной регрессии.

**Ключевые слова:** *интегральный индикатор, ранжирование альтернатив, порядковые шкалы, номинальные шкалы, обучение предпочтений.*

# 1 Введение

**Актуальность темы.** Задача построения интегральных индикаторов, или оценок качества объектов, встречается во многих прикладных областях. Особенностью этой задачи является нечисловая природа изучаемых объектов и их оценок. В зависимости от прикладной задачи, качество объектов может измеряться в различных шкалах [1]: номинальных, порядковых, интервальных, абсолютных.

Данная работа фокусируется на рассмотрении порядковых типов шкал. При этом рассматриваются отношения частичного порядка на элементах оценок объектов. Задачи подобного типа носят название задач «обучения предпочтений» (preference learning) [2]. В зависимости от постановки, в таких задачах требуется восстановить линейный или частичных порядок на множестве объектов, располагая обучающей выборкой, состоящей из матрицы «объекты-признаки» и экспертно заданными значениями интегральных индикаторов на объектах обучающей выборки.

Признаки, описывающие набор объектов, могут также иметь нечисловую природу и измеряться в различных шкалах. В данной работе рассматриваются частично упорядоченные множества значений признаков. Примером прикладной задачи подобного типа является, например, задача коллаборативной фильтрации [3].

**Цель работы.** Предложить алгоритм построения интегральных индикаторов показателей качества объектов, основываясь на экспертной информации об объектах обучающей выборки. Объекты описаны признаками, значения которых принадлежат частично упорядоченным множествам.

**Методы исследований.** При построении алгоритма использовались методы построения монотонных композиций, элементы теории бинарных отношений, методы обучения предпочтений, методы обработки нечисловой информации.

**Научная новизна.**

- Предложен способ описания порядковой шкалы конусом,
- установлено соответствие между конусом, описывающим частично упорядоченное множество, и матрицей инцидентности графа, соответствующего этому множеству,

- предложен способ корректной бинаризации различных типов шкал,
- на основе бинаризации шкал предложен алгоритм монотонной классификации, или построения интегральных индикаторов, объектов.

**Практическая ценность.** Разработан программный модуль, который

- на основе обучающей выборки строит алгоритм монотонной классификации объектов,
- оценивает качество предложенного алгоритма,
- иллюстрирует результаты.

**Положения, выносимые на защиту:**

- Метод корректной бинаризации порядковых шкал,
- алгоритм построения интегральных индикаторов объектов по их порядковому описанию.

**Апробация.** Результаты работы были использованы для определения статусов редких видов Красной книги РФ [4].

**Обзор литературы.** Задача построения интегральных индикаторов является одной из задач в области обучения предпочтений (preference learning) [2]. Всего в этой области различают три типа задач. Первый тип задач — это задачи ранжирования меток (label ranking) [5]. Помимо множества объектов, задано множество меток классов. Для каждого объекта необходимо упорядочить метки классов, задав «предпочтения» каждого объекта по классам. Частным случаем такого типа задач служит, например, задача многоклассовой классификации.

Второй тип задач в области обучения предпочтений — задачи монотонной классификации, или порядковой регрессии (instance ranking) [6]. Множество меток классов, или интегральных индикаторов, является частично упорядоченным множеством. Необходимо поставить каждому объекту в соответствие элемент этого множества, основываясь на признаковом описании объекта и обучающей информации.

Задачи третьего типа являются задачами ранжирования объектов [7]. Этот тип отличается от предыдущего тем, что не существует отдельного множества значений

классов. Ответом алгоритма на наборе объектов должна служить функция, ранжирующая эти объекты. Например, это может быть граф частичного порядка на множестве объектов.

Рассматривая в данной работе задача близка к задачам второго типа, то есть к задачам монотонной классификации. Задача монотонной классификации, или порядковой регрессии, имеет широкое применение в сфере информационного поиска [8], и потому является тщательно изученной. Принципиальное отличие данной задачи от общепринятой постановки состоит в том, что значения признаков, описывающих объекты, принадлежат частично упорядоченным множествам.

Принадлежность значений признаков частично упорядоченным множествам является одной из задач обработки нечисловой информации [9]. В задачах машинного обучения развивается направление анализа формальных понятий [10], согласно которому, объекты описываются нечисловыми атрибутами. Признаки подобного типа описываются, в том числе, т.н. решетками Галуа.

В работе рассматриваются частично упорядоченные множества значений признаков. Множества подобного типа характеризуются тем, что на них введено бинарное отношение, обладающее свойствами рефлексивности, антисимметричности и транзитивности [11]. Работа опирается на исследования в области агрегирования бинарных отношений и частичных порядков [12], активно применяющихся, например, в задачах голосования.

Частично упорядоченные множества в работе описываются выпуклыми конусами [13]. Производится процедура бинаризации частичных порядков, связанная с тем, что любая точка внутри конуса представляется выпуклой комбинацией образующих конуса. Матрица, соответствующая бинаризованному признаку, сопоставляется с матрицей инцидентности графа, соответствующему частичному порядку. Новое признаковое описание используется для построения монотонного классификатора, основанного на линейной модели.

Предлагается итеративный алгоритм оценки параметров модели. Результаты работы алгоритма сравниваются с двумя альтернативными подходами. Первый подход, статистический, использует копулы для анализа совместного распределения порядковых признаков [14]. Второй подход состоит в обобщении метода линейной регрессии на случай порядковых признаков [18].

## 2 Постановка задачи построения интегральных индикаторов

В этом разделе сформулируем постановку задачи построения интегральных индикаторов и обобщим эту постановку на случай порядковых шкал.

### 2.1 Модель построения интегральных индикаторов

Дана выборка наблюдений  $\{\mathbf{X}, \mathbf{y}_0\}$ .  $\mathbf{X}$  является матрицей плана размера  $m \times n$  и представляется набором своих столбцов:

$$\mathbf{X} = [\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n],$$

где в случае линейных шкал  $\boldsymbol{\chi}_j \in \mathbb{R}^m$ . Случай более общего типа шкал будет рассмотрен позднее.

Вектор  $\mathbf{y}_0$  является экспертной оценкой множества интегральных индикаторов объектов. Задача построения интегральных индикаторов заключается в поиске модели  $f(\mathbf{X}, \mathbf{w})$ , такой, что

$$\mathbf{y}_1 = f([\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n], \mathbf{w}).$$

Оптимальные параметры  $\hat{\mathbf{w}}$  должны минимизировать функцию ошибки

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^n} S(f([\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n], \mathbf{w}), \mathbf{y}_0).$$

Сложность задачи заключается в том, что множества значений признаков  $\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n$ , а также множество значений вектора интегральных индикаторов  $\mathbf{y}$  не являются Евклидовыми пространствами. В частности будем говорить, что множество значений признака принадлежит порядковой шкале, если оно частично упорядочено.

### 2.2 Порядковые шкалы: частично упорядоченные множества

Будем считать, что элементы  $x_{ij}$  матрицы описаний  $\mathbf{X}$  принадлежат частично упорядоченным множествам  $\mathbb{L}_j$  мощности  $L_j$ . Для упрощения обозначений, на время опустим индекс  $j$  и будем работать с одним частично упорядоченным множеством  $\mathbb{L}$  мощности  $L$ :

$$\mathbb{L} = \{l_1, \dots, l_L\}, \quad |\mathbb{L}| = L, \quad (2.1)$$

с определенной на  $\mathbb{L}$  бинарной операцией  $\succcurlyeq$ . Согласно определению частично упорядоченного множества, бинарная операция  $\succcurlyeq$  обладает тремя свойствами:

- рефлексивность,  $(l \succcurlyeq l) \forall l \in \mathbb{L}$ ;
- транзитивность,  $(l_{k_1} \succcurlyeq l_{k_2}) \wedge (l_{k_2} \succcurlyeq l_{k_3}) \Rightarrow (l_{k_1} \succcurlyeq l_{k_3}) \forall l_{k_1}, l_{k_2}, l_{k_3} \in \mathbb{L}$ ;
- антисимметричность,  $(l_{k_1} \succcurlyeq l_{k_2}) \wedge (l_{k_2} \succcurlyeq l_{k_1}) \Rightarrow (l_{k_1} = l_{k_2}) \forall l_{k_1}, l_{k_2} \in \mathbb{L}$ .

Каждому частично упорядоченному множеству  $\mathbb{L}$  соответствует ориентированный граф, вершинами которого являются элементы  $l_1, \dots, l_K$  множества  $\mathbb{L}$ . Между вершинами  $l_{k_1}, l_{k_2}$  графа существует ребро, исходящее из  $l_{k_1}$ , если выполнено  $l_{k_1} \succcurlyeq l_{k_2}$ . Такой граф описывается матрицей инцидентности  $\mathbf{L}$  размера  $L \times L$ :

$$\mathbf{L}(k_1, k_2) = \begin{cases} 1, & \text{если } l_{k_1} \succeq l_{k_2}, \\ 0, & \text{иначе.} \end{cases} \quad (2.2)$$

Матрица  $\mathbf{L}$  будет играть важнейшую роль в дальнейшем описании порядковых признаков.

### 3 Описание порядковой шкалы конусом

В этой секции поставим в соответствие частично упорядоченному множеству конус в Евклидовом пространстве. Параметризовав этот конус, поставим задачу построения интегральных индикаторов.

#### 3.1 Соответствие полиэдрального конуса порядковому признаку

Пусть матрица описаний  $\mathbf{X}$  состоит из столбцов  $\mathbf{X} = [\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n]$ , где  $\boldsymbol{\chi}_j \in \mathbb{L}_l^m$ . Здесь  $m$  — декартова степень множества  $\mathbb{L}_j$ , равная количеству объектов в выборке. Столбцу  $\boldsymbol{\chi}$ , значения которого принадлежат частично упорядоченному множеству  $\mathbb{L}$ , поставим в соответствие подмножество Евклидового пространства  $\chi \in \mathbb{R}_+^m$ :

$$\chi = \{\mathbf{x} \mid \mathbf{x} \in \mathbb{R}_+^m, \boldsymbol{\chi}_{k_1} \succcurlyeq \boldsymbol{\chi}_{k_2} \Rightarrow x_{k_1} \geq x_{k_2}\}. \quad (3.1)$$

Отметим, что множество  $\chi \subset \mathbb{R}_+^m$  и вектор  $\boldsymbol{\chi} \in \mathbb{L}_l^m$  являются разными сущностями и обозначаются одинаково с целью избежать избыточных символов.

**Теорема 3.1.** *Множество  $\chi \subset \mathbb{R}_+^m$ , определенно системой неравенств (3.1), является полиэдральным конусом в Евклидовом пространстве  $\mathbb{R}^m$  размерности, не превышающей  $L$ .*



**Доказательство.**

Во-первых, отметим, что определение (3.1) введено корректно в силу условий транзитивности, рефлексивности и антисимметричности, налагаемых на бинарное отношение  $\succsim$ .

То, что множество  $\chi$  является полиэдральным конусом, следует из определения полиэдрального конуса, задаваемого системами неравенств.

Количество различных значений компонент вектора  $\chi$  не может превышать мощности  $L$  частично упорядоченного множества  $\mathbb{L}$ . Значит, вектор  $\chi$  состоит из повторяющихся элементов (в случае  $m > L$ ), и вследствие антисимметричности в системе неравенств (3.1) появляются равенства  $\chi_{k_1} = \chi_{k_2} \Rightarrow x_{k_1} = x_{k_2}$ , разбивая множество значений признака  $\chi$  на множества одинаковых компонент. Количество таких множеств не превышает  $L$  и совпадает с размерностью конуса  $\chi$ . ■

**Представление конуса линейной комбинацией образующих.** Для того, чтобы решать задачу построения интегральных индикаторов, параметризуем конусы  $\chi_j$ , соответствующие столбцам  $\chi_j$  матрицы  $\mathbf{X}$ .

**Теорема 3.2.** *В случае, когда собственная размерность конуса  $\chi$  равна  $L$ , вектор  $\mathbf{x}$ , принадлежащий конусу  $\chi$ , соответствующий признаку  $\chi$ , представляется неотрицательной комбинацией образующих конуса*

$$\mathbf{x} = w \sum_{k=1}^L \lambda_k \zeta_k, \quad w \geq 0, \quad \sum_{k=1}^L \lambda_k = 1, \quad \lambda_k \geq 0,$$

где  $\zeta_k$  — образующая конуса  $\chi$ ,

$$\zeta_k(i) = \begin{cases} 1, & \text{если } \chi_i \succsim l_k, \\ 0, & \text{иначе,} \end{cases}$$

причем такое разложение единственно.

**Доказательство.**

Отметим, что поскольку собственная размерность конуса  $\chi$  равна  $L$ , конус определяется в точности  $L$  образующими, то есть формула разложения определена корректно. Само разложение по образующим следует из теоремы Каратеодори о выпуклой оболочке. Поскольку размерность конуса  $L$  не превышает размерности пространства  $m$ , конус  $\chi$  является симплицальным, и такое разложение единственно. ■

Отметим, что в случае размерности конуса  $\chi$  меньшей  $L$  теорема является выполненной, однако количество образующих в этом случае совпадает с размерностью конуса.

Рассмотрим пример бинаризации матрицы описаний  $\mathbf{X}$ .

**Пример 3.1.** Пусть

$$\mathbf{X} = \begin{pmatrix} l_{11} & l_{23} \\ l_{12} & l_{22} \\ l_{13} & l_{21} \end{pmatrix}, \quad \mathbb{L}_1 = \{l_{11} \succ l_{12} \succ l_{13}\}, \mathbb{L}_2 = \{l_{21} \succ l_{22} \succ l_{23}\}.$$

Тогда матрица  $\mathbf{X}$  записывается следующим образом:

$$\mathbf{X} = [\chi_1, \chi_2],$$

$$\begin{aligned} \chi_1 &= \lambda_{11}\zeta_{11} + \lambda_{12}\zeta_{12} + \lambda_{13}\zeta_{13} = \lambda_{11} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \lambda_{12} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} + \lambda_{13} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \\ \chi_2 &= \lambda_{21}\zeta_{21} + \lambda_{22}\zeta_{22} + \lambda_{23}\zeta_{23} = \lambda_{21} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \lambda_{22} \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} + \lambda_{23} \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \end{aligned}$$

при выполнении условий

$$\sum_{k=1}^{L_j} \lambda_{jk} = 1, \quad \lambda_{jk} \geq 0. \quad (3.2)$$

Согласно определению образующих  $\zeta_k$ , элементы этих образующих являются, по сути, элементами матрицы инцидентов графа, соответствующего частичному порядку. Проиллюстрируем это на примере.

**Пример 3.2.**

$$\begin{array}{ccc} l_1 & \longrightarrow & l_2 \\ & \searrow & \\ & & l_3 \end{array} \quad (3.3)$$

Рассмотрим пример частично упорядоченного множества  $\mathbb{L}$ , задаваемого графом (3.3). Это множество состоит из трех элементов  $l_1, l_2, l_3$ , для которого выполня-

ются условия  $l_1 \succ l_2$ ,  $l_1 \succeq l_3$ . Согласно определению (2.2), матрица  $\mathbf{L}$ , соответствующая графу (3.3),

$$\mathbf{L} = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Строки этой матрицы используются при формировании матрицы  $\mathbf{Z}$  для заданного набора объектов.

Из теоремы 3.2 и примера 3.2 видно, что для бинаризации частичного порядка и составления образующих конуса достаточно воспользоваться матрицей смежности графа, соответствующего частично упорядоченному множеству.

### 3.2 Нахождение оптимальных параметров модели в случае порядковых признаков

Задача построения интегральных индикаторов заключается в поиске модели  $f(\mathbf{X}, \mathbf{w})$ , такой, что

$$\mathbf{y}_1 = f([\chi_1, \dots, \chi_n], \mathbf{w}).$$

Оптимальные параметры  $\hat{w}$  должны минимизировать функцию ошибки

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^n} S(f([\chi_1, \dots, \chi_n], \mathbf{w}), \mathbf{y}_0).$$

Будем рассматривать линейные модели, то есть вектор  $\mathbf{y}_1$  представляется линейной комбинацией

$$\mathbf{y}_1 = w_1 \chi_1 + \dots + w_n \chi_n.$$

Множества  $\chi_j$  являются подмножествами  $\mathbb{R}^m$ , однако благодаря параметризации, задаваемой теоремой 3.2, линейная модель запишется в следующем виде:

$$\mathbf{y}_1 = w_1 \boldsymbol{\lambda}_1 \mathbf{Z}_1 + \dots + w_n \boldsymbol{\lambda}_n \mathbf{Z}_n, \quad (3.4)$$

где векторы  $\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n$  являются векторами элементов выпуклых комбинаций конусов  $\chi_j$ ,

$$w_j \in \mathbb{R}, \quad \boldsymbol{\lambda}_j \in \Lambda = \{\boldsymbol{\lambda} | \boldsymbol{\lambda} \geq \mathbf{0}, \|\boldsymbol{\lambda}\|_1 = 1\},$$

а столбцы матриц  $\mathbf{Z}_1, \dots, \mathbf{Z}_n$  являются образующими конусов  $\chi_1, \dots, \chi_n$ , определяемыми теоремой 3.2.

Оптимальные параметры  $\mathbf{w}$  должны минимизировать функцию ошибки:

$$S(\mathbf{f}([\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n], \mathbf{w}), \mathbf{y}_0) \rightarrow \min_{\substack{\mathbf{w} \in \mathbb{R}^n; \\ \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n, \boldsymbol{\lambda}_0 \in \Lambda}} .$$

В линейном случае (3.4), а также в случае квадратичной функции ошибки,

$$S(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|,$$

задача поиска оптимальных параметров запишется следующим образом:

$$(\hat{\mathbf{w}}, \hat{\boldsymbol{\lambda}}_1, \dots, \hat{\boldsymbol{\lambda}}_n, \hat{\boldsymbol{\lambda}}_0) = \arg \min_{\substack{\mathbf{w} \in \mathbb{R}^n; \\ \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n, \boldsymbol{\lambda}_0 \in \Lambda}} \|\boldsymbol{\lambda}_0 \mathbf{Z}_0 - \sum_{j=1}^n w_j \boldsymbol{\lambda}_j \mathbf{Z}_j\|. \quad (3.5)$$

Отметим, что в данном случае экспертная оценка  $\mathbf{y}_0$  выставлена в ранговой шкале, и ее параметризация также дается теоремой 3.2 с матрицей образующих  $\mathbf{Z}_0$  и соответствующим вектором выпуклой комбинации  $\boldsymbol{\lambda}_0$ .

**Алгоритм вычисления оптимальных параметров.** Отметим, что оптимальные параметры, соответствующие минимизации (3.5), не вычисляются явным образом. Предлагается итеративный алгоритм. На каждом четном шаге при фиксированных параметрах выпуклой комбинации  $\hat{\boldsymbol{\lambda}}_1, \dots, \hat{\boldsymbol{\lambda}}_n, \hat{\boldsymbol{\lambda}}_0$  вычисляется вектор параметров модели  $\hat{\mathbf{w}}$ , затем при фиксированном  $\hat{\mathbf{w}}$  вычисляются  $\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n, \boldsymbol{\lambda}_0$ .

Принимается начальное приближение алгоритма:  $\hat{\lambda}_{jk} = 1/L_j$ , то есть образующие конусов имеют одинаковый вес, и рассматриваются центральные точки всех конусов.

На каждом следующем шаге производится процедура последовательной минимизации функции ошибки:

1.  $\hat{\mathbf{w}} := \arg \min_{\mathbf{w} \in \mathbb{R}} \|\hat{\boldsymbol{\lambda}}_0 \mathbf{Z}_0 - \sum_{j=1}^n w_j \hat{\boldsymbol{\lambda}}_j \mathbf{Z}_j\|,$
2.  $\hat{\boldsymbol{\lambda}}_j = \arg \min_{\boldsymbol{\lambda}_j \in \Lambda} \|\hat{\boldsymbol{\lambda}}_0 \mathbf{Z}_0 - \sum_{j=1}^n \hat{w}_j \hat{\boldsymbol{\lambda}}_j \mathbf{Z}_j\|,$  для всех  $j = 1, \dots, n$ .

Выполнена теорема о то скорости нахождения алгоритмом оптимального решения.

**Теорема 3.3.** *В случае линейной модели 3.4, итеративный алгоритм находит оптимальное решение не более чем за  $\Theta$  итераций, где*

$$\Theta = L_0 - 1 + L_1 - 1 + \dots + L_n - 1 + n = \sum_{j=0}^n L_j - 1 \quad -$$

*собственная размерность пространства параметров.*

**Доказательство.**

Идея доказательства заключается в следующем: на каждой итерации алгоритм находит в соответствующем конусе вектор, ближайший к линейной комбинации остальных. В самом трудном случае, когда конусы не пересекаются, оптимальные векторы лежат на образующих или гранях конусов. На каждом шаге алгоритм выбирает в каждом конусе соответствующую грань, не ухудшая решение. Всего количество граней всех конусов совпадает с собственной размерностью пространства признаков  $\Theta$ , то есть и количество шагов алгоритма не превосходит  $\Theta$ . ■

**Множество параметров модели.** Рассмотрим более детально множество параметров  $\mathbf{w}$  модели  $f$ . Согласно теореме 3.2, вектор  $\mathbf{w} \in \mathbb{R}_+$  является неотрицательным, однако минимизация ошибки 3.5 проводится по всему множеству  $\mathbb{R}$ . Такое обобщение позволяет учитывать также отрицательно монотонную зависимость интегральных индикаторов  $\mathbf{y}_1$  от признаков. Если предполагается исключительно неотрицательная зависимость, то шаг 1 итеративного алгоритма переписывается следующим образом:

$$\hat{\mathbf{w}} := \arg \min_{\mathbf{w} \in \mathbb{R}_+} \left\| \hat{\boldsymbol{\lambda}}_0 \mathbf{Z}_0 - \sum_{j=1}^n w_j \hat{\boldsymbol{\lambda}}_j \mathbf{Z}_j \right\|.$$

На множестве параметров модели могут быть заданы экспертные предпочтения. Например, эксперт может задать линейный порядок на элементах множества параметров, указав важность признаков, например, так:

$$w_1 \geq \dots \geq w_n \geq 0.$$

В этом случае, минимизация параметров модели  $w$  может проводиться по экспертно заданному множеству  $\mathcal{R} = \{w_1 \geq \dots \geq w_n \geq 0\}$ , и шаг 1 алгоритма переписывается следующим образом:

$$\hat{\mathbf{w}} := \arg \min_{\mathbf{w} \in \mathcal{R}} \left\| \hat{\boldsymbol{\lambda}}_0 \mathbf{Z}_0 - \sum_{j=1}^n w_j \hat{\boldsymbol{\lambda}}_j \mathbf{Z}_j \right\|.$$

В более общем случае, множество  $\mathcal{R}$  может быть произвольным экспертно заданным множеством  $\mathcal{R} \in \mathbb{R}^n$ .

### 3.3 Непараметрическое описание конусов, соответствующих порядковым признакам

Рассмотрим метод построения интегрального индикатора  $\mathbf{y}_1$  по матрице описаний  $\mathbf{X} = [\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n]$  и экспертной оценке  $\mathbf{y}_0$ , не вводя множество дополнительных

параметров [14].

В нашем случае, столбцы матрицы  $\mathbf{X}$  являются конусами  $\chi_1, \dots, \chi_n$ . Любой полиэдральный конус определяется системой линейных неравенств

$$\chi_j = \{\mathbf{x}_j | A_j \mathbf{x}_j \leq 0\}, j = 1, \dots, m.$$

В общем случае, множество  $\chi_j$  является многогранником, определяемым матрицей  $A_j$ . Например, в случае линейного порядка на элементах вектора  $\mathbf{x}_j$ ,  $x_{1j} \geq x_{2j} \geq \dots \geq x_{mj} \geq 0$ , матрица

$$A_j = \begin{pmatrix} -1 & 1 & 0 & \dots & 0 \\ 0 & -1 & 1 & \dots & 0 \\ 0 & 0 & 0 & \dots & -1 \end{pmatrix},$$

эта матрица определяет выпуклый многогранный конус  $\chi_j$ .

Принимается линейная модель построения интегрального индикатора. Это означает, что вектор  $\mathbf{y}_1$  представляется в виде суммы:

$$\mathbf{y}_1 = \mathbf{x}_1 + \dots + \mathbf{x}_m, \mathbf{x}_1 \in \chi_1, \dots, \mathbf{x}_m \in \chi_m.$$

Для дальнейшего рассмотрения напомним понятие суммы Минковского двух множеств. Суммой Минковского двух подмножеств  $L_1$  и  $L_2$  линейного пространства  $L$  называется множество  $L'$ , состоящее из всевозможных сумм векторов из  $L_1$  и  $L_2$ . Из этих двух определений логично вытекает определение области значений линейной модели. Областью значений для линейной модели назовем сумму Минковского

$$\chi = \chi_1 + \dots + \chi_m.$$

Заметим, что введенное таким образом определение линейной модели с ранговым описанием объектов обобщает линейную модель в ее стандартном понимании:

$$\mathbf{y} = \sum_{j=1}^m \mathbf{X}_j w_j,$$

поскольку каждое слагаемое  $\mathbf{X}_j w_j$  для всех  $w_j$  является подмножеством конуса  $\chi_j$ . Следовательно, в случае рангового описания исчезают веса признаков.

Таким образом, для построения интегрального индикатора требуется построить сумму Минковского выпуклых многогранников. Для ее построения воспользуемся методом точного построения суммы Минковского двух многогранников, заданных системами линейных неравенств [15].

Итак, пусть дано два выпуклых многогранника, задающихся системами неравенств:

$$\chi_1 = \{\mathbf{x}_1 | A_1 \mathbf{x}_1 \leq \mathbf{b}_1\}, \chi_2 = \{\mathbf{x}_2 | A_2 \mathbf{x}_2 \leq \mathbf{b}_2\}.$$

Их суммой Минковского будет являться вектор  $\mathbf{x}$ , являющийся решением системы:

$$\begin{cases} \mathbf{x} - \mathbf{x}_1 - \mathbf{x}_2 = 0, \\ A_1 \mathbf{x}_1 \leq \mathbf{b}_1, \\ A_2 \mathbf{x}_2 \leq \mathbf{b}_2, \end{cases}$$

которая заменой переменной  $\mathbf{x}_1 = \mathbf{x} - \mathbf{x}_2$  преобразуется в систему:

$$\begin{cases} A_1 \mathbf{x} - A_2 \mathbf{x}_2 \leq \mathbf{b}_1, \\ A_2 \mathbf{x}_2 \leq \mathbf{b}_2. \end{cases}$$

Справедливо следующее утверждение: вектор  $\mathbf{x} \in \chi$  тогда и только тогда, когда найдется вектор  $\mathbf{x}_2$ , удовлетворяющий системе 3.3.

Таким образом, задача поиска вектора  $\mathbf{x}$  сводится к решению системы линейных неравенств:

$$C \mathbf{x}_2 \leq \mathbf{d}, C = \begin{pmatrix} -A_1 \\ A_2 \end{pmatrix}, \mathbf{d} = \begin{pmatrix} \mathbf{b}_1 - A_1 \mathbf{x} \\ \mathbf{b}_2 \end{pmatrix}. (2)$$

Для решения этой системы используем вариант леммы Минковского-Фаркаша [16], который формулируется следующим образом.

Пусть  $A$  и  $b$  — матрица и вектор. Разрешимость системы линейных неравенств  $Ax \leq b$  эквивалентна тому, что  $yb \geq 0$  для любого вектора-строки  $y \geq 0$  со свойством  $yA = 0$ . В нашем случае, эта лемма записывается так:

$$\exists \mathbf{x}_2 : C \mathbf{x}_2 \leq \mathbf{d} \Leftrightarrow \forall \mathbf{z} : C^T \mathbf{z} = 0, \mathbf{z} \geq 0 \rightarrow (\mathbf{d}, \mathbf{z}) \geq 0.$$

Пусть  $V$  является фундаментальной системой решений для этого случая, причем  $V = \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$ , где  $V_i$  — ФСР, соответствующая матрице  $A_i$ . Тогда условие  $(\mathbf{d}, \mathbf{z}) \geq 0$  переписывается в виде:

$$V_1^T (\mathbf{b}_1 - A_1 \mathbf{x}) + V_2^T \mathbf{b}_2 \geq 0.$$

Таким образом, принимая  $A = V_1^T A_1$ ,  $\mathbf{b} = V_1^T \mathbf{b}_1 + V_2^T \mathbf{b}_2$ , мы получим параметры  $A, \mathbf{b}$  системы неравенств, описывающей сумму Минковского  $\chi_1 + \chi_2$ .

Заметим, что отдельной трудностью является нахождение фундаментальной системы решений  $V$ , каждый столбец которой должен быть  $\geq 0$ . Алгоритмы отыскания такой системы описан в [17].

Проекцией точки  $x$  на множество  $D$  называется вектор  $P(x) \in D$ , удовлетворяющий условию:

$$P_D(x) = \arg \min_{y \in D} \|x - y\|.$$

Построив область значения модели  $\chi$ , определим уточненный интегральный индикатор  $\mathbf{y}_1$  как проекцию экспертной оценки  $\mathbf{y}_0$  на  $\chi$ :

$$\mathbf{y}_1 = P_\chi(\mathbf{y}_0).$$

Отметим, что эта проекция является единственной в силу выпуклости множества  $\chi$ .

Рассмотрим небольшой пример. Пусть матрица

$$X = \begin{pmatrix} 3 & 1 \\ 2 & 2 \\ 1 & 3 \end{pmatrix}.$$

Тогда для нее матрицы

$$A_1 = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{pmatrix}, A_2 = \begin{pmatrix} 0 & 1 & -1 \\ 1 & -1 & 0 \\ -1 & 0 & 0 \end{pmatrix}.$$

Используя матрицу  $C = \begin{pmatrix} 0 & -1 & 1 \\ -1 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & -1 \\ 1 & -1 & 0 \\ -1 & 0 & 0 \end{pmatrix}$ , получаем

$$A = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ -1 & 1 & -1 \\ 0 & 0 & -1 \end{pmatrix}$$

и неравенство для суммы:  $\chi = \{\mathbf{x} | A\mathbf{x} \leq 0\}$ .

Таким образом, область значений линейной модели записывается в виде:

$$\chi = \begin{cases} \mathbf{x}_1 \geq 0, \\ \mathbf{x}_2 \geq 0, \\ \mathbf{x}_3 \geq 0, \\ \mathbf{x}_1 + \mathbf{x}_3 \geq \mathbf{x}_2. \end{cases}$$



Если, например, экспертная оценка  $\mathbf{y}_0 = (1, 5, 2)$ , то ее проекция на область значения  $\chi$ :  $\mathbf{y}_1 = (\frac{5}{3}, \frac{16}{3}, \frac{11}{3})$ .

## 4 Альтернативные методы

В этом разделе рассмотрим альтернативные подходы к решению задачи построения интегральных индикаторов объектов, описанных в порядковых шкалах. Рассмотрим статистический подход, основанный на совместном распределении критериев [18], а также алгоритм криволинейной регрессии [19]. В этой секции будем рассматривать только линейно упорядоченные множества значений признаков.

### 4.1 Построение интегрального индикатора в порядковых шкалах с использованием копул

Предлагается найти отображение  $a : X \rightarrow Y$ , минимизирующее функционал среднего риска. Минимум среднего риска достигается алгоритмом

$$a(x) = \arg \max_{y \in Y} P(y|\mathbf{x}_i),$$

где  $P(y|\mathbf{x}_i)$  — апостериорная вероятность класса  $y$  для объекта  $\mathbf{x}$ . Эта вероятность является условной по  $\mathbf{x}$ . Для оценки апостериорной вероятности  $P(y|\mathbf{x}_i)$  будем использовать копулы.

**Свойства копул, используемые для оценки условной вероятности.**

**Определение 4.1.** Функция  $C : [0, 1]^d \rightarrow [0, 1]$  называется копулой размерности  $d$ , если выполняются следующие условия:

$$C(u_1, \dots, u_{i-1}, 0, u_{i+1}, \dots, u_d) = 0,$$

$$C(1, \dots, 1, u, 1, \dots, 1) = u,$$

$$B = \prod_{i=1}^d [a_i, b_i] \subseteq [0, 1]^d : \int_B dC(u) \geq 0.$$

Выполнение этих свойств означает, что функция  $C$  является функцией распределения многомерной случайной величины  $[u_1, \dots, u_d]$ , такой, что одномерное распределение каждого из  $u_i$  равномерно на интервале  $[0, 1]$ .

Важным фактом, позволяющим применять копулы для построения регрессионных моделей, является следующая теорема.

**Теорема 4.1.** *Многомерная функция распределения случайной величины:*

$$H(x_1, \dots, x_d) = P[X_1 \leq x_1, \dots, X_d \leq x_d]$$

случайного вектора  $(X_1, \dots, X_d)$  с одномерными функциями распределения

$$F_i(x) = P[X_i \leq x_i]$$

может быть записана в виде:

$$H(x_1, \dots, x_d) = C(F_1(x_1), \dots, F_d(x_d)).$$

Таким образом, для оценивания совместного распределения  $H$  случайных величин  $X_1, \dots, X_d$  достаточно оценить их одномерные распределения  $F_i(x_i)$  и функцию копулы, связывающую эти случайные величины.

Следующая теорема утверждает, что функция копулы не изменяется при действии на случайные величины любых монотонных преобразований.

**Теорема 4.2.** *Пусть  $X, Y$  — две случайные величины с совместной функцией распределения  $H(x, y)$ . Пусть также  $\varphi, \psi$  — две монотонных функции, преобразующие случайные величины  $X$  и  $Y$  в*

$$Z = \varphi(X), \quad T = \psi(Y)$$

с совместной функцией распределения  $H'(Z, T)$ . Тогда копула, связывающая случайные величины  $Z$  и  $T$ :

$$C'(F'(z), G'(t)) = H'(z, t) = C(F'(z), G'(t)),$$

то есть,

$$C' = C.$$

Таким образом, чтобы оценить функцию копулы, описывающую связь между случайными величинами  $X_1, \dots, X_d$ , достаточно знать только ранговые соотношения этих случайных величин. Абсолютные значения величин  $X_1, \dots, X_d$  используются только при оценивании их одномерных распределений.

Для решения задачи классификации таксонов необходимо знать апостериорную вероятность (4.1). Эта вероятность выражается через частную производную функции копулы  $C$ , о чем утверждает следующая теорема.

**Теорема 4.3.** Пусть  $X, Y$  — две случайные величины с одномерными функциями распределения  $F(X), G(Y)$ . Тогда условная вероятность  $P(Y \leq y | X = x)$  равна частной производной копулы:

$$P(Y \leq y | X = x) = \frac{\partial}{\partial v} C(u, v) |_{(G(y), F(x))},$$

взятой в точке

$$u = G(y), \quad v = F(X).$$

В нашей задаче имеется  $n$  случайных величин, соответствующих признакам, и случайная величина  $Y$ .

Для оценки условной вероятности необходимо ввести некоторые дополнительные обозначения.

Имеется набор объектов  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ . Каждый объект описывается  $n$  признаками. Обозначим одномерные функции распределения  $y$  и всех компонент многомерной случайной величины  $\mathbf{x}$ :

$$G_Y^0(y), G_{X^1}^1(x^1), \dots, G_{X^n}^n(x^n).$$

Обозначим совместные функции распределения упорядоченных поднаборов - векторов  $\mathbf{x}^k = (x^1, \dots, x^k)$  размерности от 1 до  $n$ :

$$F_{\mathbf{X}^k}^k(\mathbf{x}^k), \quad \mathbf{x}^k = (x^1, \dots, x^k), \quad k = 1, \dots, n.$$

Для нахождения условной вероятности  $P(Y \leq y | \mathbf{x}_i)$ , воспользуемся частной производной копулы  $C(u, v)$  по переменной  $u$ :

$$P(Y \leq y | \mathbf{x}_i) = \frac{\partial}{\partial u} C(u, v) |_{F_{\mathbf{X}^n}^n(\mathbf{x}_i^n), G_Y(y)},$$

взятой в точке

$$u = F_{\mathbf{X}^n}^n(\mathbf{x}_i^n), \quad v = G_Y(y).$$

Неизвестной в этой формуле является функция совместного распределения  $F_{\mathbf{X}^n}^n$ . Чтобы найти эту функцию, воспользуемся теоремой 3:

$$F_{\mathbf{X}^n}^n(\mathbf{x}^n) = C^{n-1}(u, v) |_{F_{\mathbf{X}^{n-1}}^{n-1}(\mathbf{x}^{n-1}), G_{X_n}^n(x_n)},$$

...

$$F_{\mathbf{X}^i}^i(\mathbf{x}^i) = C^{i-1}(u, v) |_{F_{\mathbf{X}^{i-1}}^{i-1}(\mathbf{x}^{i-1}), G_{X_i}^i(x_i)},$$

...

$$F_{X^1, X^2}^2(x^1, x^2) = C^1(u, v)|_{G_{X^1}^1(x^1), G_{X^2}^2(x^2)}.$$

Таким образом, чтобы оценить апостериорную вероятность  $P(Y \leq y | \mathbf{x}_i)$ , необходимо оценить все  $n + 1$  одномерные распределения  $y$  и компонент случайного вектора  $\mathbf{x}$ , а также  $n$  копул  $C, C^1, \dots, C^{n-1}$ .

**Копулы, используемые при построении интегрального индикатора.** Для решения задачи (4.1) предлагается использовать Архимедовскую копулу:

**Определение 4.2.** Копула  $C(u_1, \dots, u_d)$  называется архимедовской, если для нее выполнены следующие условия:

$$C(u_1, \dots, u_d) = \psi(\psi^{-1}(u_1) + \dots + \psi^{-1}(u_d)),$$

где функция  $\psi$  называется генератором, и для нее должны быть выполнены:

$$(-1)^k \psi^{(k)}(x) \geq 0$$

для всех  $x \geq 0$  и  $k = 0, 1, \dots, d - 2$ . А также, функция

$$(-1)^{d-2} \psi^{d-2}(x)$$

должны быть невозрастающей и выпуклой.

Будем использовать частные случаи Архимедовской копулы, задаваемые следующими функциями-генераторами:

копула Клейтона,

$$\psi(t) = (1 + \theta t)^{-\frac{1}{\theta}}, \quad \theta \in \Theta = (0, \infty)$$

и копула Гумбеля,

$$\psi(t) = \exp(-t^{\frac{1}{\theta}}), \quad \theta \in \Theta = [1, \infty).$$

Отметим, что эти семейства копул зависят только от одного параметра  $\theta$ , что значительно упрощает задачу в вычислительном смысле.

В случае копулы Гумбеля, частная производная имеет следующий вид:

$$\frac{\partial}{\partial u} C(u, v) = \left( \frac{\ln u}{\ln C} \right)^{\theta-1} \frac{C}{u}.$$

**Оценка параметров копулы.** Как было сказано выше, для оценки параметра  $\theta \in \Theta$  копулы используются не сами случайные величины  $X, Y$ , а последовательности рангов этих величин. Выборкам  $X$  и  $Y$  соответствуют последовательности рангов:

$R_x = (R_{x_1}, \dots, R_{x_m})$ , где  $R_{x_i}$  – ранг  $i$ –го объекта в вариационном ряду выборки  $X$ ,

$R_y = (R_{y_1}, \dots, R_{y_n})$ , где  $R_{y_i}$  – ранг  $i$ –го объекта в вариационном ряду выборки  $Y$ .

Отметим, что наиболее часто используемым методом оценки параметров распределения является метод максимизации правдоподобия, который в случае копул записывается следующим образом:

$$L(\theta) = \sum_{i=1}^m \log \left( c_\theta(F(X_i), G(Y_i)) \right),$$

$$c_\theta(u, v) = \frac{\partial^2}{\partial u \partial v} C_\theta(u, v).$$

Вместо значений функций одномерных распределений  $F(X_i), G(Y_i)$  можно подставить их эмпирические значения, получив таким образом функцию псевдоправдоподобия [?]:

$$L'(\theta) = \sum_{i=1}^m \left( \log c_\theta \left( \frac{R_i}{m+1}, \frac{S_i}{m+1} \right) \right).$$

Заметим, что функция  $L'$  зависит только от самой копулы  $C_\theta$ , то есть, в нашем случае, только от параметра  $\theta$ , и максимизация этой функции не представляет собой большой вычислительной сложности.

Благодаря этому способу, задача оценки распределений  $F_{\mathbf{X}^i}^i(\mathbf{x}^i)$  распадается на два независимых этапа: оценка параметра  $\theta_i$  копул  $C^i$  путем максимизации псевдоправдоподобия и оценка параметров одномерных распределений  $G_Y^0(y), G_{X^1}^1(x^1), \dots, G_{X^n}^n(x^n)$  с помощью метода максимума правдоподобия.

**Алгоритм оценки апостериорного распределения.** Приведем подробный алгоритм оценки распределений  $F_{\mathbf{X}^i}^i(\mathbf{x}^i)$ . Как было сказано выше, необходимо оценить  $n+1$  одномерное распределение  $G_Y^0(y), G_{X^1}^1(x^1), \dots, G_{X^n}^n(x^n)$  и  $n$  функций копулы  $C, C^1, \dots, C^n$ .

1. Оцениваются одномерные распределения  $G_{X^1}^1(x^1), G_{X^2}^2(x^2)$ . Все функции  $G_{X^i}^i(x^i)$  будем искать в классе бета-распределений. То есть, распределение случайной

величины  $X$  задается плотностью вероятности  $g_X$ , имеющей вид:

$$\begin{cases} g_X(x) &= \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \\ B(\alpha, \beta) &= \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx. \end{cases}$$

Параметры  $\alpha$  и  $\beta$  для этого распределения оцениваются методом моментов. Для этого численно решается система уравнений:

$$\begin{cases} E(X) &= \frac{\alpha}{\alpha+\beta}, \\ D(X) &= \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}. \end{cases}$$

2. Оценим копулу  $C^1(u, v)$ , связывающую переменные  $x^1$  и  $x^2$ , максимизируя функцию псевдоправдоподобия:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L'_{12}(\theta) = \sum_{i=1}^m \left( \log c_{\theta} \left( \frac{R_{x_1}}{m+1}, \frac{R_{x_2}}{m+1} \right) \right).$$

3. Оценив одномерные распределения  $G_{X^1}^1(x^1)$ ,  $G_{X^2}^2(x^2)$  и копулу  $C^1(u, v)$ , получаем оценку функции совместного распределения  $F_{X^1, X^2}^2(x^1, x^2)$ . Повторяем шаги 1-2, каждый раз прибавляя по одному новому признаку  $x^i$  и оценивая на шаге 3 функцию  $F_{\mathbf{X}^i}^i(\mathbf{x}^i)$ .

4. повторив  $n$  раз шаги 1-2, получим функцию совместного распределения всех признаков  $F_{\mathbf{X}^n}^n$ . На последнем шаге оценим функцию распределения  $G_Y^0(y)$ , копулу  $C(u, v)$ , связывающую  $Y$  и  $X$ , и найдем  $\hat{y}$ , доставляющий максимум апостериорной вероятности:

$$\hat{y} = \arg \max_y P(Y \leq y | \mathbf{x}_i) = \arg \max_y \frac{\partial}{\partial u} C(u, v) |_{F_{\mathbf{X}^n}^n(\mathbf{x}_i^n), G_Y(y)},$$

взятой в точке

$$u = F_{\mathbf{X}^n}^n(\mathbf{x}_i^n), \quad v = G_Y(y),$$

где

$$F_{\mathbf{X}^n}^i(\mathbf{x}_i^i) = C^{i-1}(u, v) |_{F_{\mathbf{X}^{i-1}}^{i-1}(\mathbf{x}^{i-1}), G_{X_i}^i(x_i)},$$

взятой в точке

$$u = F_{\mathbf{X}^{i-1}}^{i-1}(\mathbf{x}^{i-1}), \quad v = G_{X_i}^i(x_i)$$

для всех

$$i = 2, \dots, n.$$

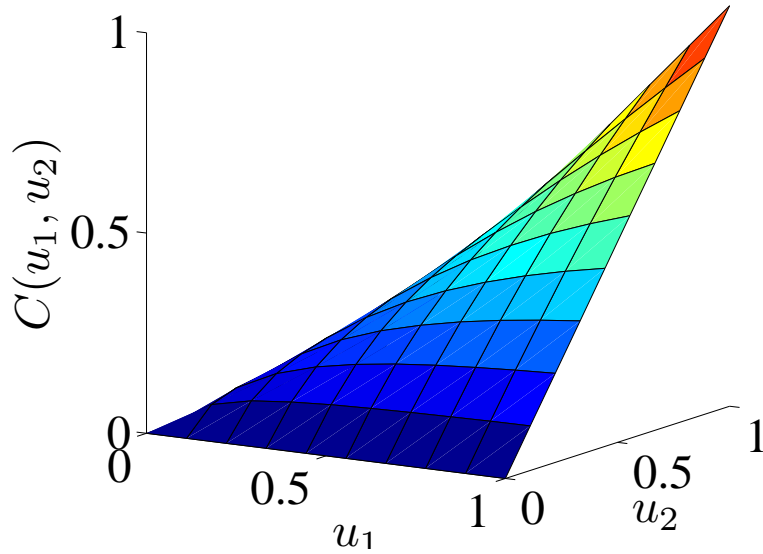


Рис. 1: Копула Клейтона

**Выбор признаков.** Для повышения качества классификации на контрольной подвыборке, предлагается выбрать наиболее информативные признаки. Множество индексов признаков, включенных в функцию вероятности 4.1, назовем активным набором и обозначим  $\mathcal{A} \subseteq \mathcal{J}$ .

Для того, чтобы выбрать наиболее информативные признаки, предлагается использовать следующий эвристический алгоритм. Информационными будем считать те признаки, которые имеют наибольшую ранговую связь со случайной величиной  $Y$ . Чтобы понять, какие признаки имеют наибольшую связь, рассмотрим некоторые свойства копул о ранговой связи.

**Утверждение 4.1.** *Случайные величины  $X$  и  $Y$  являются независимыми тогда и только тогда, когда*

$$C(u, v) = uv, \quad u, v \in [0, 1],$$

где

$$C(F(x), G(y)) = H(x, y),$$

где  $H(x, y)$  — совместная функция распределения случайных величин  $X$  и  $Y$ .

**Утверждение 4.2.** *Границы Фреше для копулы:*

$$W(u, v) \leq C(u, v) \leq M(u, v), \quad u, v \in [0, 1],$$

где

$$W(u, v) = \max(0, u + v - 1)$$

— минимальная копула,

$$M(u, v) = \min(u, v)$$

— максимальная копула.

Причем, если  $C(u, v) = W(u, v)$ , то  $Y$  — монотонно убывающая функция  $X$ , если  $C(u, v) = M(u, v)$ , то  $Y$  — монотонно возрастающая функция  $X$ .

Для примера, рассмотрим копулу Гумбеля (4.1):

$$C_\theta(u, v) = \exp \left[ \left( (-\log(u))^\theta + (-\log(v))^\theta \right)^{\frac{1}{\theta}} \right] \quad \theta \geq 1.$$

При стремлении параметра копулы  $\theta \rightarrow 1$ ,  $C_\theta(u, v) \rightarrow uv$ , то есть, случайные величины являются независимыми. При стремлении параметра  $\theta \rightarrow \infty$ , ранговая связь между случайными величинами возрастает. Таким образом, ранговая связь изменяется монотонно при варьировании параметра копулы. Для решения задачи отбора признаков будем отбирать те из них, для которых параметр копулы со случайной величиной  $Y$  является наибольшим.

Исходя из этого рассуждения, предлагается следующий алгоритм.

1. Примем пустое множество активных признаков

$$\mathcal{A} = \emptyset.$$

2. Для всех  $j = 1, \dots, n$  вычислим параметры  $\theta_j$  для копул  $C_{\theta_j}(F_i(x^j), G(y))$  и включим в набор

$$\mathcal{A} = \mathcal{A} \cup \{k\}$$

тот признак  $k$ , для которого

$$k = \arg \max_{j \in \mathcal{J}} \theta_j.$$

Обозначим множество оставшихся признаков

$$\mathcal{J}' = \mathcal{J} \setminus \mathcal{A}.$$

3. Для всех признаков  $j \in \mathcal{A}$  и всех  $k_1, \dots, k_{\mathcal{A}}$  вычислим параметры  $\theta_j$  для копул

$$C_{\theta_j}(F_i(x_i), H_{k_1, \dots, k_{\mathcal{A}}}(x^{k_1}, \dots, x^{k_{\mathcal{A}}}))$$

и включим в набор

$$\mathcal{A} = \mathcal{A} \cup \{k\}$$



тот признак  $k$ , для которого

$$k = \arg \max_{j \in \mathcal{J}'} \theta_j.$$

4. Будем повторять шаг 3, пока значение ошибки на контрольной выборке не стабилизируется.

## 4.2 Алгоритм криволинейной регрессии

**Постановка задачи.** Криволинейная модель  $f(\mathbf{w}, \mathbf{x}_i)$  имеет вид

$$f(\mathbf{w}, \mathbf{x}_i) = \xi(\mathbf{b}_0, h(\mathbf{w}, \mathbf{x}_i)), \quad (4.1)$$

$$h(\mathbf{w}, \mathbf{x}_i) = \sum_{j \in \mathcal{J}} u_j g(\mathbf{b}_j, x_{ij}). \quad (4.2)$$

где вектор параметров  $\mathbf{w} = [\mathbf{b}_0; \mathbf{b}_1; \dots; \mathbf{b}_n; \mathbf{u}] = [\mathbf{b}_0^T, \mathbf{b}_1^T, \dots, \mathbf{b}_n^T, \mathbf{u}^T]^T$  состоит из векторов  $\mathbf{b}_j$  — параметров монотонной коррекции  $j$ -го признака  $\chi_j$  и весовых коэффициентов признаков  $\mathbf{u} = [u_1, \dots, u_j, \dots, u_n]^T$ . Функция  $g$  монотонной коррекции задана следующим образом:

$$g(j, \chi) : \chi \mapsto \mathbf{b}_j = \begin{cases} 1 \mapsto b_{j1}, \\ 2 \mapsto b_{j2}, \\ \dots \\ k_j \mapsto b_{jk_j}. \end{cases}$$

При этом соблюдается условие монотонности параметров,

$$\text{Ord}(\mathbf{b}_j) : 0 < b_{j1} < b_{j2} < \dots < b_{jk_j} < 1 \quad \text{для } j = 1, \dots, n \quad \text{и} \quad (4.3)$$

$$\text{Ord}(\mathbf{b}_0) : b_{01} < b_{02} < \dots < b_{0k_0}.$$

Функция  $\xi(\mathbf{b}_0, h(\mathbf{w}, \mathbf{x}_i))$  определяет для числа  $h(\mathbf{w}, \mathbf{x}_i)$  ближайшую по модулю компоненту вектора  $\mathbf{b}_0$ :

$$\xi(\mathbf{b}_0, h(\mathbf{w}, \mathbf{x}_i)) = \arg \min_{j \in \mathcal{J}} |b_{0j} - h(\mathbf{w}, \mathbf{x}_i)|.$$

Введя обозначение для матрицы скорректированных экспертных оценок

$$G = [g_{ij}] = [g(\mathbf{b}_j, x_{ij})], \quad i \in \mathcal{I}, j \in \mathcal{J},$$

перепишем (4.1) и (4.2) в виде модели интегрального индикатора

$$f(\mathbf{w}, \mathbf{x}_i) = \xi(\mathbf{b}_0, [G\mathbf{u}]_i). \quad (4.4)$$

Назначим функцией ошибки модели сумму квадратов регрессионных остатков,

$$S(\mathbf{w}) = \|\mathbf{f}(\hat{\mathbf{w}}, X) - \mathbf{y}\|_2^2 + \lambda \|\hat{\mathbf{u}}\|_2^2,$$

включающую регуляризующее слагаемое с фиксированным коэффициентом  $\lambda$ , где  $\hat{\mathbf{w}}$  и  $\hat{\mathbf{u}}$  — параметры, которые необходимо оценить.

**Оценивание параметров модели.** Оценивание параметров  $\mathbf{w}$  модели  $\mathbf{f}$  выполняется итеративно. Перед началом итераций значения векторов  $\mathbf{b}_0, \mathbf{b}_1, \dots, \mathbf{b}_n$  назначены таким образом, что функция  $g$  является тождественной,  $g = \text{id}$ . Оценивание параметров выполняется в три этапа. Сначала при фиксированных значениях векторов  $\hat{\mathbf{b}}_0, \dots, \hat{\mathbf{b}}_n$  оцениваются весовые коэффициенты

$$\hat{\mathbf{u}} = \arg \min_{\mathbf{u} \in \mathbb{R}^n} S([\hat{\mathbf{b}}_0; \dots; \hat{\mathbf{b}}_n; \mathbf{u}]).$$

Затем при фиксированных значениях коэффициентов  $\hat{\mathbf{u}}$  оцениваются параметры монотонной коррекции

$$[\mathbf{b}_1; \dots; \mathbf{b}_n] = \arg \min_{\text{Ord}(\mathbf{b}_1), \dots, \text{Ord}(\mathbf{b}_n)} S([\hat{\mathbf{b}}_0; \dots; \hat{\mathbf{b}}_n; \mathbf{u}])$$

с учетом требования монотонности значений этих параметров. На последнем этапе оценивается вектор  $\mathbf{b}_0$

$$\mathbf{b}_0 = \arg \min_{\text{Ord}(\mathbf{b}_0)} S([\hat{\mathbf{b}}_0; \dots; \hat{\mathbf{b}}_n; \mathbf{u}]).$$

Итерации выполняются до стабилизации функции ошибки  $S$ .

Рассмотрим эти три этапа более подробно. За начальное приближение примем столбцы матрицы  $G$

$$\hat{G} = [\mathbf{g}(\hat{\mathbf{b}}_1, \boldsymbol{\chi}_1), \dots, \mathbf{g}(\hat{\mathbf{b}}_n, \boldsymbol{\chi}_n)] = [\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n],$$

поскольку, как было сказано выше,  $g = \text{id}$ , и вектор  $\hat{\mathbf{y}} = \mathbf{y}$ . Таким образом, векторы  $\hat{\mathbf{b}}_0, \hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n$  в начальном приближении в качестве элементов содержат элементы множеств  $\mathbb{L}_0, \mathbb{L}_1, \dots, \mathbb{L}_n$ .

**Шаг 1.** Найдем  $\hat{\mathbf{u}}$  при фиксированных  $\hat{\mathbf{b}}_0, \dots, \hat{\mathbf{b}}_n$ :

$$\hat{\mathbf{u}} = \arg \min_{\mathbf{u}} \|\hat{\mathbf{y}} - \hat{G}\mathbf{u}\| + \lambda \|\mathbf{u}\|.$$

Решение на шаге 1 имеет вид:

$$\hat{\mathbf{u}} = (\hat{G}^T \hat{G} + \lambda I)^{-1} \hat{G}^T \hat{\mathbf{y}}.$$

**Шаг 2.** При фиксированных  $\hat{\mathbf{b}}_0, \hat{\mathbf{u}}$  оценим скорректированную матрицу описаний

$$G = [\mathbf{g}(\mathbf{b}_1, \boldsymbol{\chi}_1), \dots, \mathbf{g}(\mathbf{b}_n, \boldsymbol{\chi}_n)] = [\mathbf{g}_1, \dots, \mathbf{g}_n].$$

Для каждого  $\mathbf{g}_j \in \mathbb{R}^m$  будем вычислять вектор  $\hat{\mathbf{g}}_j$ , являющийся монотонной коррекцией исходного вектора  $\mathbf{g}_j$ :

$$\begin{cases} [\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_n] = \arg \min \|\boldsymbol{\xi}(\mathbf{b}_0, G\hat{\mathbf{u}}) - \hat{\mathbf{y}}\|_2^2, \\ \text{из } g_{ij_1} \leq g_{ij_2} \text{ следует } \hat{g}_{ij_1} \leq \hat{g}_{ij_2} \quad i \in \mathcal{I}, j_1, j_2 \in \mathcal{J}, \\ g_{ij} \in [0, 1] \quad i \in \mathcal{I}, j \in \mathcal{J}, \text{ согласно (4.3)}. \end{cases}$$

По векторам  $\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_n$  затем однозначно восстанавливаются векторы  $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n$  как упорядоченные векторы, содержащие различные элементы  $\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_n$ . Для решения этой задачи используется алгоритм градиентного спуска, описанный в [20].

**Шаг 3.** Наконец, при фиксированных  $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n, \hat{\mathbf{u}}$  оценим вектор  $\mathbf{b}_0$  и  $\hat{\mathbf{y}} = \mathbf{g}(\mathbf{b}_0, \mathbf{y})$ :

$$\hat{\mathbf{b}}_0 = \arg \min_{\text{Ord}(\mathbf{b}_0)} \|\boldsymbol{\xi}(\mathbf{b}_0, \hat{G}\hat{\mathbf{u}}) - \mathbf{g}(\mathbf{b}_0, \mathbf{y})\|_2^2.$$

**Выбор признаков при классификации.** Множество индексов признаков, включенных в модель, назовем активным набором и обозначим  $\mathcal{A} \subseteq \mathcal{J}$ .

Поставим задачу выбора наиболее информативных признаков следующим образом. Разобьем выборку  $\mathcal{D}$  на две подвыборки, обучающую и тестовую. Обозначим индексы элементов этих подвыборок соответственно  $\mathcal{L} \sqcup \mathcal{T} = \mathcal{I}$ . Для некоторого активного набора признаков  $\mathcal{A}$  найдем на обучающей подвыборке  $\mathcal{D}_{\mathcal{L}}$  оптимальные, согласно заданной функции ошибки  $S$ , параметры  $\hat{\mathbf{w}}_{\mathcal{A}}$ ,

$$\hat{\mathbf{w}}_{\mathcal{A}} = \arg \min_{\mathbf{w}} S(\mathbf{w}_{\mathcal{A}} | \mathcal{D}_{\mathcal{L}}).$$

Затем выберем наиболее информативные признаки — активный набор  $\hat{\mathcal{A}}$  по всем поднаборам индексов признаков  $\mathcal{A} \subseteq \mathcal{J}$ , доставляющий на тестовой выборке  $\mathcal{D}_{\mathcal{T}}$  минимум функции ошибки:

$$\hat{\mathcal{A}} = \arg \min_{\mathcal{A} \subseteq \mathcal{J}} S(\hat{\mathbf{w}}_{\mathcal{A}} | \mathcal{D}_{\mathcal{T}}).$$

Для выбора наиболее информативного подмножества признаков используется итеративный алгоритм добавления признаков.

На первом шаге этого алгоритма принимается активное множество информативных признаков  $\hat{\mathcal{A}} = \emptyset$ . На каждом следующем шаге к множеству  $\hat{\mathcal{A}}$  добавляется признак с индексом  $\hat{j}$ , такой что

$$\hat{j} = \arg \min_{j \in \mathcal{T} \setminus \hat{\mathcal{A}}} S(\hat{\mathbf{w}}_{\hat{\mathcal{A}} \cup \{j\}} | \mathcal{D}_{\mathcal{T}}).$$

Эта процедура продолжается итеративно до тех пор, пока значение функции ошибки  $S$  на контрольной выборке  $\mathcal{D}_{\mathcal{T}}$  не достигнет минимума.

## 5 Вычислительные эксперимент

### 5.1 Реальные данные: категоризация редких видов Красной книги РФ

Рассматривается задача построения интегрального индикатора в ранговых шкалах. В качестве практического приложения рассматривается проблема определения статуса угрожаемых видов животных, входящих в список Красной книги РФ. В Красной книге РФ принята следующая категоризация редкости видов (таксонов) по степени угрозы их исчезновения. Имеется шесть различных категорий статуса (меток классов) таксонов: 0 — вероятно исчезнувшие, 1 — находящиеся под угрозой исчезновения, 2 — сокращающиеся в численности, 3 — редкие, 4 — неопределенные по статусу, 5 — восстанавливаемые и восстанавливающиеся. Эта категоризация является монотонной: метки классов ранжированы по возрастанию биологического разнообразия.

Каждый таксон описан набором признаков, отражающих его состояние. Эксперт, владеющий информацией о таксоне, выставляет оценку для каждого признака в ранговой шкале. Таким образом, задана матрица «объект-признак», состоящая из описаний таксонов и вектор меток классов таксонов. Требуется построить модель, восстанавливающую класс таксона из Красной книги РФ по его описанию.

Задача ревизии Красной книги РФ и построения модели вычисления интегрального индикатора является актуальной из-за постоянного пополнения книги новыми записями о таксонах.

## 5.2 Сходимость параметров

**Алгоритм параметризованных конусов.** Сходимость алгоритма в случае малого количества параметров показана на рис. 2. Видно, что значение ошибки на обучающей и контрольной выборке монотонно уменьшается. На рис. 3 показана сходимость алгоритма для случай количества признаков, сравнимого с количеством объектов. Из-за большой подгонки под обучающие данные, значение ошибки на контроле сначала уменьшается, затем начинает расти. На рис. 4 показан случай большого количества параметров. Алгоритм переобучается с первой же итерации, и ошибка на контроле начинает расти сразу.

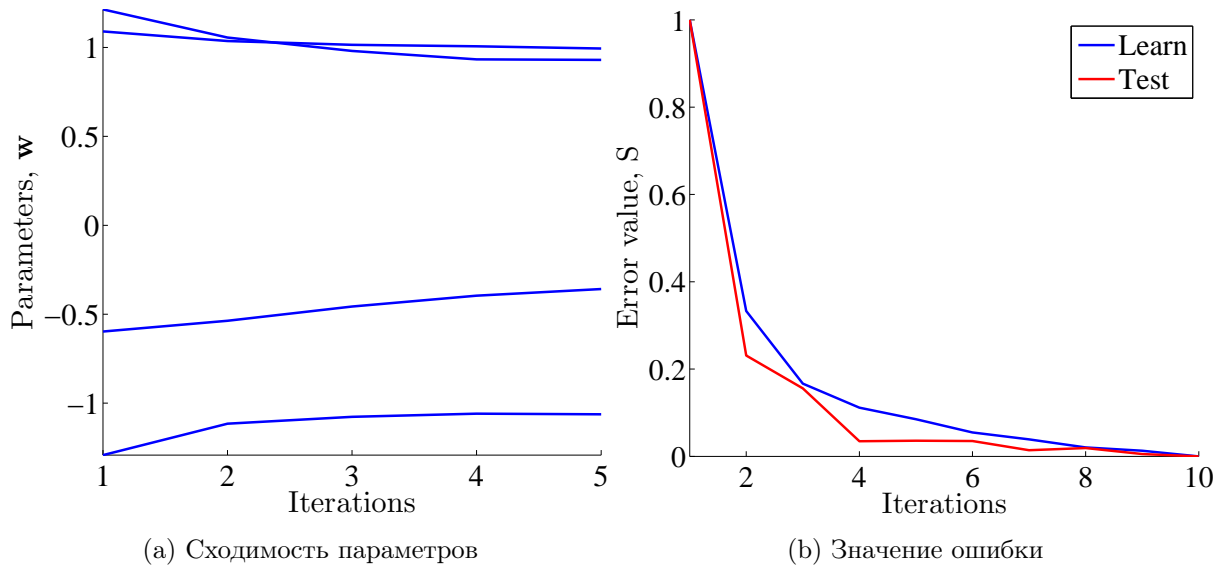


Рис. 2: Случай малого количества параметров

**Алгоритм криволинейной регрессии.** Сходимость параметров алгоритма криволинейной регрессии показана на рис. 5, 6, 7.

На рис. 5 показана сходимость весов регрессии  $u$ , на рис. 6 — сходимость элементов вектора  $b_0$ . По оси абсцисс отложено количество итераций, по оси ординат — количественное значение каждого признака. Видно, что сходимость наблюдается на десятой итерации.

На рис. 7 показана зависимость функции ошибки от количества выбираемых признаков. Видно, что ее минимум достигается при семи признаках, и значение средней ошибки равно  $Q = 0.75$ .

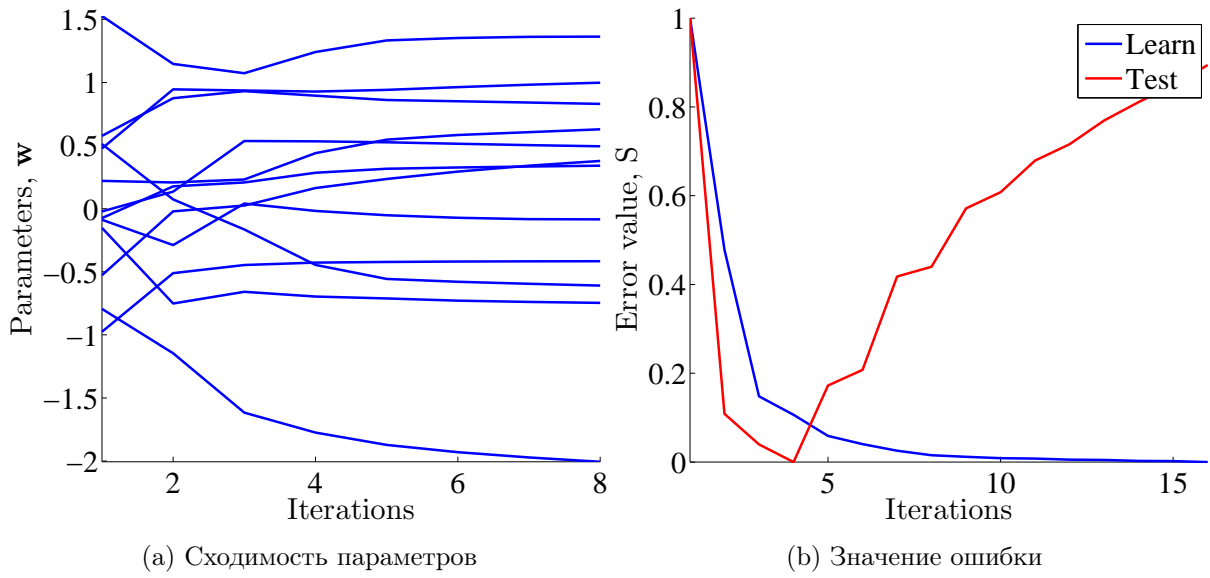


Рис. 3: Случай среднего количества параметров

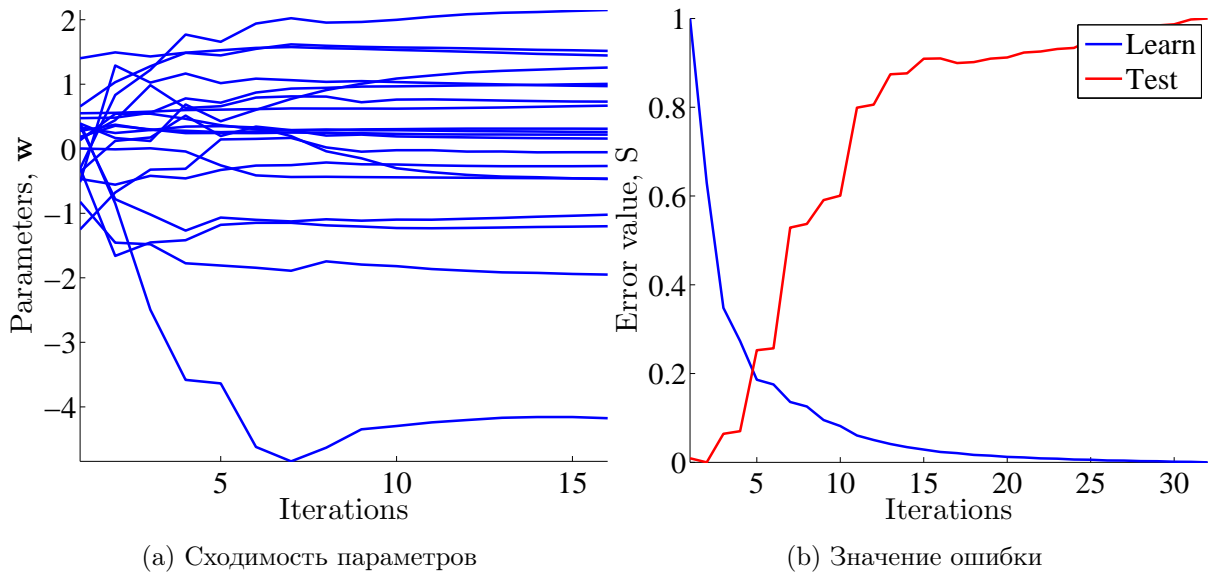


Рис. 4: Случай большого количества параметров, переобучение

**Статистический алгоритм.** На рис. 8а показана зависимость ошибки классификации от количества выбранных признаков. Оптимальное значение достигается при  $|\mathcal{A}| = 4$ . В исходной таблице данных эти признаки индексированы номерами 22, 24, 23 и 20.

На рис. 8б показана зависимость параметра копулы от количества выбранных признаков. Видно, что значение параметра монотонно убывает с ростом количества признаков.

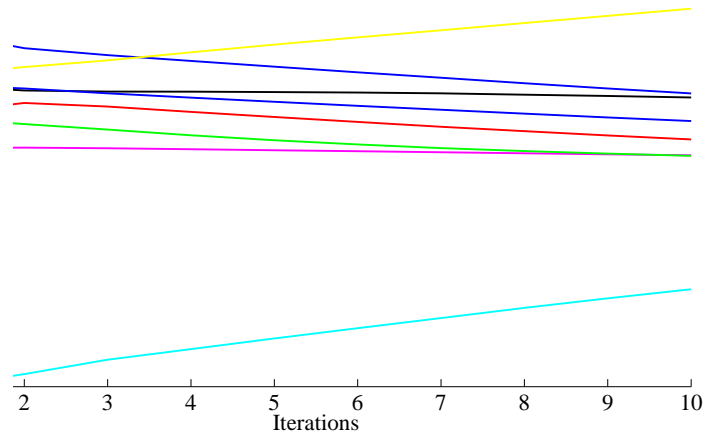


Рис. 5: Сходимость весов признаков

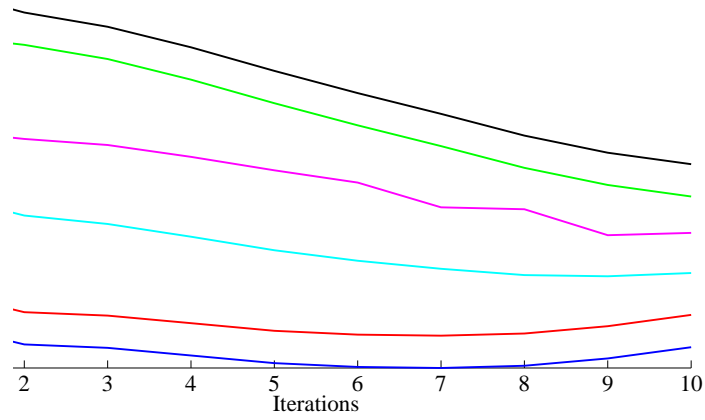


Рис. 6: Сходимость элементов вектора  $\mathbf{b}_0$ .

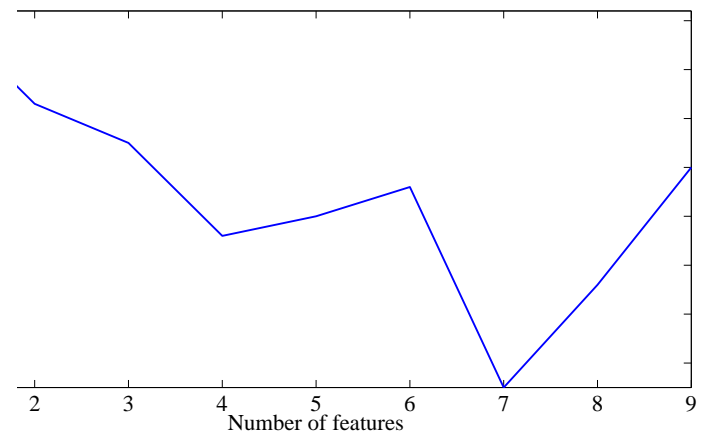


Рис. 7: Зависимость функции ошибки от количества выбираемых признаков.

**Сравнение результатов.** В ходе вычислительного эксперимента сравнивались результаты алгоритма, основанного на параметризации конусов, и двух альтернатив-

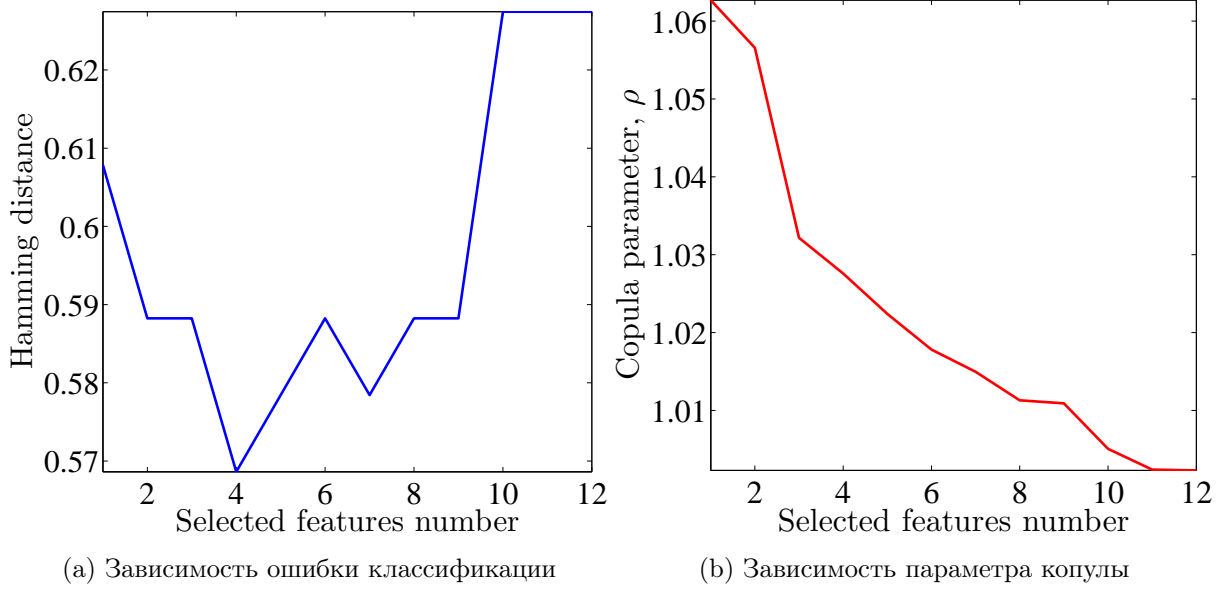


Рис. 8: Зависимость классификации от количества выбранных признаков

ных алгоритмов.

Оценивалась ошибка на обучении и контроле, контрольная выборка выбиралась методом Leave-One-Out. Результаты эксперимента представлены в табл. 1.

Таблица 1: Сравнение алгоритмов классификации

Алгоритм	Средняя ошибка на обучении	LOO	Время построения модели, сек
Конусы	0.29	0.58	1.2
Криволинейная регрессия	0.57	0.71	3.6
Копулы	0.57	0.61	0.25

## 6 Заключение

В работе предложен алгоритм построения интегральных индикаторов объектов, описанных в смешанных шкалах. Предложено описание порядковой шкалы частично упорядоченным множеством. В соответствие частично упорядоченному множеству поставлен конус. Установлено соответствие между образующими этого конуса



и столбцами матрицы инцидентности графа, соответствующего частично упорядоченному множеству. Предложено параметризовать конус, задающий множество значений признака. Предложен итеративный алгоритм поиска оптимальных параметров путем минимизации функции ошибки, доказана его сходимость к оптимальному решению. Также предложен непараметрический подход, основанный на построении допустимого множества значений интегральных индикаторов как суммы конусов входящих в модель признаков.

Для оценки качества предлагаемых алгоритмов рассмотрено два альтернативных подхода к решению задачи. Первым является статистический подход, в котором взаимосвязь между порядковыми признаками устанавливается функцией копулы. Вторым подход является обобщением линейной регрессии на случай порядковых признаков. Приведены иллюстрации сходимости параметров алгоритмов и сравнение результатов на примере категоризации редких видов Красной книги РФ.

## Список литературы

- [1] А. И. Орлов. *Эконометрика*. Экзамен, 2002.
- [2] Johannes Fuernkranz and Eyke Huellermeier. *Preference learning*. Springer, 2011.
- [3] J. Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen. Collaborative filtering recommender systems. *Lecture Notes in Computer Science*, 4321:291–324, 2007.
- [4] *Красная книга Российской Федерации (животные)*. М: АСТ Астрель, 2001.
- [5] Eyke Huellermeier, Johannes Fuernkranz, Weiwei Cheng, and Klaus Brinker. Label ranking by learning pairwise preferences. *Artificial Intelligence*, 172:1897–1916, November 2008.
- [6] Weiwei Cheng, Michael Rademaker, Bernard De Baets, and Eyke Huellermeier. Predicting partial orders: Ranking with abstention. *Machine Learning and Knowledge Discovery in Databases Lecture Notes in Computer Science*, 6321:215–230, 2010.
- [7] Frank Wilcoxon. Individual comparisons by ranking methods. *Breakthroughs in Statistics*, pages 196–202, 1992.
- [8] Andrew Trotman. Learning to rank. *Information Retrieval*, 8:381, 2005.
- [9] А. И. Орлов. *Нечисловая статистика*. МЗ-Пресс, 2004.
- [10] Sergei Obiedkov Sergei Kuznetsov. Comparing performance of algorithms for generating concept lattices. *Journal of Experimental & Theoretical Artificial Intelligence*, 14:189–216, 2002.
- [11] П. С. Александров. *Введение в теорию множеств и общую топологию*. Наука, 1977.
- [12] Gunther Schmidt. *Relational Mathematics (Encyclopedia of Mathematics and its Applications)*. Cambridge University Press, 2010.
- [13] Saroj B. Malik Sujit Kumar Mitra, P. Bhimasankaram. *Matrix Partial Orders, Shorted Operators and Applications*. 2010.

- [14] М.П. Кузнецов and В.В. Стрижов. Построение интегрального индикатора с использованием ранговой матрицы описаний. *Интеллектуализация обработки информации. Доклады 9-й международной конференции*, pages 130–132, 2012.
- [15] М.П. Кузнецов. Построение интегрального индикатора в ранговых шкалах с использованием копул для анализа совместного распределения критериев. *Машинное обучение и анализ данных*, 1:411–419, 2012.
- [16] М. В. Уханов. Алгоритм построения суммы многогранников. *Вестник ЮУрГУ*, pages 39–44, 2001.
- [17] А. Схрейвер. Теория линейного и целочисленного программирования. *Ж. вычисл. матем. и матем. физ.*, 1:360, 1991.
- [18] Н. Б. Черникова. Алгоритм для нахождения общей формулы неотрицательных решений системы линейных уравнений. *Ж. вычисл. матем. и матем. физ.*, 4:733–738, 1964.
- [19] М. П. Кузнецов, В. В. Стрижов, and М.М Медведникова. Алгоритм многоклассовой классификации объектов, описанных в ранговых шкалах. *Научно-технический вестник СПб ГПУ. Информатика. Телекоммуникации. Управление*, 5, 2012.
- [20] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.