

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«САМАРСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ АКАДЕМИКА С.П. КОРОЛЕВА»  
(САМАРСКИЙ УНИВЕРСИТЕТ)

*Р.О. МИШАНОВ*

# ИСПОЛЬЗОВАНИЕ МЕТОДОВ СТАТИСТИЧЕСКОЙ КЛАССИФИКАЦИИ ДЛЯ ОПРЕДЕЛЕНИЯ КАЧЕСТВА КОМПОНЕНТОВ БОРТОВЫХ РЭС

Рекомендовано редакционно-издательским советом федерального государственного автономного образовательного учреждения высшего образования «Самарский национальный исследовательский университет имени академика С.П. Королева» в качестве учебного пособия для обучающихся по основной образовательной программе высшего образования по направлению подготовки 11.04.03 Конструирование и технология электронных средств

САМАРА  
Издательство Самарского университета  
2021

УДК 621.37(075)

ББК 32.844я7

М710

Рецензенты: д-р техн. наук, проф. Л. Д. Л о ж к и н ,  
д-р техн. наук, доц. А. И. Д а н и л и н

*Мишанов, Роман Олегович*

**М710** **Использование методов статистической классификации для определения качества компонентов бортовых РЭС: учебное пособие / Р.О. Мишанов.** – Самара: Издательство Самарского университета, 2021. – 80 с.: с ил.

**ISBN 978-5-7883-1696-3**

В пособии приведена необходимая информация для понимания методов классификации ЭРИ на этапе изготовления РЭС, также приведены методы классификации, разработанные на основе положений кластерного анализа, искусственных нейронных сетей.

Учебное пособие предназначено для магистрантов факультета электроники и приборостроения, изучающих дисциплину «Кластерный анализ качества электронных средств». Оно может быть полезно широкому кругу специалистов, связанных с изготовлением и эксплуатацией радиоэлектронной аппаратуры.

Подготовлено на кафедре «Конструирование и технология электронных систем и устройств» Самарского университета.

УДК 621.37(075)

ББК 32.844я7

ISBN 978-5-7883-1696-3

© Самарский университет, 2021

# ОГЛАВЛЕНИЕ

<b>Введение</b> .....	5
<b>1. Обеспечение надёжности бортовых РЭС</b> .....	6
1.1. Особенности применения ЭРИ в бортовых РЭС.....	6
1.2. Технические методы обеспечения надёжности бортовых РЭС .....	8
1.3. Повышение качества и надёжности бортовых РЭС за счет прогнозирования состояния используемой ЭКБ .....	9
<b>2. Основные понятия в теории статистической обработки данных</b> .....	11
2.1. Понятие случайной величины, генеральной совокупности и выборки .....	11
2.2. Меры центральной тенденции .....	12
2.3. Меры разброса данных .....	13
2.4. Нормальное распределение случайной величины .....	13
2.5. Шкалирование данных (нормализация и стандартизация данных).....	15
<b>3. Основные положения теории распознавания образов</b> .....	17
3.1. Теория распознавания образов.....	17
3.2. Классификация и кластеризация.....	18
3.3. Пример использования теории распознавания образов для классификации ЭРИ .....	19
<b>4. Использование методов кластерного анализа для классификации ЭРИ</b> .....	21
4.1. Кластерный анализ.....	21
4.2. Характеристики кластеров .....	22
4.3. Близость объектов (метод сходства, функции расстояний) .....	23
4.4. Методы кластерного анализа .....	25

4.5. Иерархические методы кластеризации .....	27
4.6. Метод объединения (правила объединения) .....	28
4.7. Итеративные методы кластеризации.....	29
4.8. Алгоритм k-средних (k-means clustering).....	30
4.9. Блок-схема алгоритма иерархической кластеризации для классификации ЭРИ.....	32
4.10. Блок-схема алгоритма k-средних для классификации ЭРИ.....	34
<b>5. Использование нейросетевых структур для классификации ЭРИ.....</b>	<b>36</b>
5.1. Искусственные нейронные сети (ИНС) и их применение.....	36
5.2. Использование персептрона для решения задачи классификации ЭРИ .....	44
5.3. Самоорганизующиеся системы (карты самоорганизации) и их применение для решения задачи классификации ЭРИ...	47
<b>Библиографический список .....</b>	<b>52</b>
<b>Приложение А. Пример решения задачи классификации методом иерархической кластеризации .....</b>	<b>57</b>
<b>Приложение Б. Пример решения задачи классификации методом k-средних кластерного анализа .....</b>	<b>68</b>

## **ВВЕДЕНИЕ**

Любое изделие электронной техники характеризуется качеством, т.е. набором свойств, существенно отличающих это изделие от других. Более того, качество характеризует степень пригодности изделия к использованию по назначению. При эксплуатации вследствие влияния дестабилизирующих факторов окружающей среды, процессов старения, износа материалов характеристики изделия изменяются, что влечет за собой изменение качества этого изделия. При этом поддержание определенного уровня качества и надёжности аппаратуры является довольно обширной и сложной задачей.

Огромное влияние на состояние изделия оказывают комплектующие изделия, и в особенности, заложенная элементная база, которая определяет функциональность этого изделия. Более того, элементная база постоянно усложняется, что требует внедрения новых методов контроля комплектующих изделий, методов прогнозирования состояния изделий для повышения качества и надёжности выпускаемой продукции, особенно в такой специфичной области как космическое приборостроение.

# 1. ОБЕСПЕЧЕНИЕ НАДЁЖНОСТИ БОРТОВЫХ РЭС

## 1.1. Особенности применения ЭРИ в бортовых РЭС

Надёжность радиоэлектронных средств (РЭС) космического назначения в значительной степени зависит от качества используемой элементной базы. А с учетом тенденций к разработке космических аппаратов (КА) с длительными сроками активного существования (САС) 10–15 лет и более, обеспечение бесперебойного функционирования является сложной и дорогостоящей задачей.

В электронных изделиях ракетно-космической техники разрешена к применению квалифицированная элементная база категорий качества «ОС» и «ОСМ», в некоторых случаях допускается применение элементной базы категории «ВП», но с проведением дополнительных испытаний, подтверждающих её качество. У таких ЭРИ минимальная наработка до отказа составляет не менее 135000–140000 ч. Стоит отметить, что вся элементная база должна пройти обязательный 100% диагностический неразрушающий контроль (ДНК) и входной контроль (ВК). Стандарт [1] устанавливает, что в проектах КА с длительными САС разрешено использовать ЭРИ иностранного производства (ЭРИ ИП) уровня качества «Spacе» с соответствующим сертификатом.

Такая же ситуация наблюдается и для КА с САС менее 10–15 лет, с оговоркой, что разрешается использование ЭРИ ИП уровня качества «Military», а для КА с САС менее 5 лет разрешено использовать ЭРИ ИП уровня качества «Industrial» [2].

Также стоит отметить, что разработчик обязан ограничивать номенклатуру применяемой элементной базы, которая должна со-

ответствовать отраслевому ограничительному перечню («Перечень ЭКБ 1–22 Минпромторг России»). ЭРИ ИП должны проходить сертификационные испытания на соответствие требованиям отечественных стандартов с учетом модели внешних воздействующих факторов (ВВФ). Необходимость проведения сертификационных испытаний ЭРИ ИП обусловлена:

- отсутствием возможности контроля технологического процесса производства ЭРИ ИП;
- отсутствием нормативно-технической документации на поставляемые ЭРИ ИП (технических условий на ЭРИ ИП);
- отсутствием обоснованного ограничительного перечня ЭРИ ИП, разрешенных к применению в РЭА;
- рассогласованием в требованиях отечественных и зарубежных стандартов (проблема гармонизации) [3].

Также существует практика использования неквалифицированных ЭРИ (не имеющих категорию качества «ОС» и «ОСМ») и несертифицированных ЭРИ (категории качества «Industrial», «Commercial»), но с обязательным проведением контрольных и дополнительных отбраковочных испытаний [4–10].

Стоит отметить, что работы по дополнительным испытаниям ЭРИ (входной контроль, отбраковочные испытания, ДНК, выборочный разрушающий физический анализ (РФА)) являются сложными и дорогостоящими, и как правило, проводятся в испытательных технических центрах (ИТЦ), имеющих специализированное оборудование и разрабатывающих методики ДНК в ходе собственных исследовательских работ [11].

На данный момент в ракетно-космической отрасли определен курс на импортозамещение ЭКБ ИП на ЭКБ ОП, являющийся одновременно актуальной, но сложной задачей. Более того, на некоторые типы ЭРИ ИП на данный момент не существует аналогов среди ЭРИ ОП с необходимыми тактико-техническими характери-

стиками, что приводит к ускоренному проведению НИОКР по разработке таких аналогов [12].

## **1.2. Технические методы обеспечения надёжности бортовых РЭС**

Обеспечение надёжности РЭС космического назначения является важнейшей и трудоёмкой задачей разработчика и изготовителя сложной аппаратуры. Технические методы обеспечения надёжности РЭС космического назначения классифицируют на:

- схемно-конструктивные;
- производственно-технологические;
- эксплуатационные [13].

Схемно-конструктивные методы используют на этапе разработки изделия. К ним относят использование систем автоматизированного проектирования (САПР) с использованием углубленного математического моделирования различных конструктивно-технологических вариантов (КТВ) исполнения с определением оптимальных условий и режимов эксплуатации, определением производственных и эксплуатационных запасов по основным параметрам изделий.

Производственно-технологические методы основаны на использовании статистического контроля и статистического регулирования технологических процессов изготовления изделий, применении отбраковочных испытаний изделий с целью выявления экземпляров со скрытыми дефектами, а также проведения физико-технического анализа дефектных изделий с целью выработки решений, влияющих на технологический процесс изготовления.

К эксплуатационным методам обеспечения надёжности РЭС относят методы, позволяющие снизить интенсивность отказов, например, использование изделия в ограниченном температурном

диапазоне согласно указанному в технических условиях (ТУ). Также к ним относят проведение предремонтного анализа отказывающихся изделий, переход с планового ремонта на ремонт по фактическому состоянию.

К техническим методам обеспечения надёжности применительно к космической технике также можно отнести описанные ранее:

- квалификационные испытания ЭРИ, при которых проводят подтверждение соответствия заданным требованиям к надёжности, условиям функционирования в аппаратуре;
- дополнительные отбраковочные испытания, проводимые по специальным программам и предназначенные для исключения возможных производственных дефектов [13].

### **1.3. Повышение качества и надёжности бортовых РЭС за счет прогнозирования состояния используемой ЭКБ**

Прогнозирование – исследовательский процесс, в результате которого получают вероятностные данные о будущем состоянии прогнозируемого объекта (процесса) [14]. Методы прогнозирования широко применяются при разработке бортовых технических устройств:

- для обоснования необходимого уровня надежности на этапе разработки, анализе требований ТЗ, при проработке технических предложений;
- в случае невозможности применения других методов расчета надежности, например, при отсутствии полной информации на ранних стадиях проектирования;
- для расчета интенсивностей отказов элементов изделий разных типов с учетом их качества изготовления, нагруженности, области применения аппаратуры [15].

Прогнозирование состояния ЭРИ на определенный срок функционирования РЭС космического назначения – один из способов повышения качества и надёжности сложной аппаратуры. Наиболее успешно эту задачу можно решить с помощью методов индивидуального прогнозирования, позволяющих по измеренным на данный момент времени параметрам конкретного экземпляра ЭРИ провести прогноз состояния этого экземпляра с использованием математической модели с заранее выбранным упреждением.

Прогнозирование может быть качественным и количественным [14]. Под качественным прогнозированием понимается прогноз в случае необходимости получения информации о качестве прогнозируемого показателя (параметра) объекта (процесса), например, к какому классу к времени прогноза  $t_{np}$  при оценке по прогнозируемому параметру будет относиться технический объект: к классу годных или к классу потенциально дефектных. Под количественным прогнозированием понимается прогноз в случае, если необходимо знание величины прогнозируемого показателя (параметра) объекта (процесса). Количественный прогноз тесно связан с вероятностью, с которой произойдет то или иное событие в будущем, а также с некоторыми количественными характеристиками этого события. Применительно к техническим устройствам, такой прогноз уместен при необходимости получения информации о величине прогнозируемого показателя надёжности к времени прогноза  $t_{np}$ , по которому можно судить о состоянии изделия. Следует отметить, что качественный анализ можно получить через цепь логических суждений (качественно), либо через количественный анализ.

## **2. ОСНОВНЫЕ ПОНЯТИЯ В ТЕОРИИ СТАТИСТИЧЕСКОЙ ОБРАБОТКИ ДАННЫХ**

### **2.1. Понятие случайной величины, генеральной совокупности и выборки**

Случайная величина – переменная, которая в результате испытания в зависимости от случая принимает одно из возможного множества своих значений [16]. Действительно, интересующие нас параметры ЭРИ, измеряемые в разные моменты времени, в разных условиях и с разной точностью измерений с высокой степенью вероятности не покажут один и тот же результат. Поэтому измеряемые параметры можно считать дискретными случайными величинами.

Статистическая совокупность – множество единиц изучаемого явления, объединенных одной качественной основой, но отличающихся друг от друга отдельными свойствами [17, 18].

Изучение объектов может охватывать всю совокупность или ограничиваться исследованием некоторой части. В первом случае исследование называют полным (сплошным), во втором – частичным (выборочным). Конечно, полное исследование позволяет получить исчерпывающую информацию об объектах, но к такому исследованию прибегают редко, т.к. это требует огромных затрат времени и ресурсов, а иногда такое исследование вообще нереализуемо с практической точки зрения.

Совокупность, из которой отбирают определенную часть её членов для совместного изучения, называют генеральной совокупностью. Отобранная часть этой генеральной совокупности называется выборкой.

Выборочный метод – основной метод изучения статистических закономерностей.

Основное требование, предъявляемое к выборке, заключается в том, чтобы эта выборка наиболее полно отражала состояние генеральной совокупности, т.е. быть представительной или репрезентативной. Другими словами, выборка должна отражать те же свойства, что представлены в генеральной совокупности, и в том же процентном соотношении [19].

## 2.2. Меры центральной тенденции

Среднее арифметическое значение (выборочное среднее,  $M$ ) – средняя оценка изучаемого параметра:

$$M = \frac{\sum_{i=1}^n x_i}{n},$$

где  $M$  – среднее арифметическое значение,  $n$  – количество экземпляров объектов;  $x_i$  – значение наблюдаемого параметра  $i$ -го экземпляра. Стоит отметить, что при стремлении объёма выборки к бесконечности, среднее арифметическое значение приближается к математическому ожиданию наблюдаемой генеральной совокупности.

Мода ( $M_o$ ) – наиболее часто встречающееся значение наблюдаемого параметра. В интервальном частотном распределении мода является серединой интервала с максимальной частотой.

Медиана ( $M_e$ ) – некоторое значение, находящееся в середине последовательностей показателей в ранжированном по возрастанию или убыванию ряду значений. Особенностью медианы является то, что от значения медианы слева и справа находятся по 50% от всего количества наблюдаемых результатов.

Знание значений моды и медианы необходимо для того, чтобы установить, является ли распределение значений изучаемого пара-

метра симметричным и приближено ли оно к нормальному распределению. В таком случае среднее арифметическое значение, мода и медиана совпадают (в идеальном случае), либо незначительно разнятся [19].

### 2.3. Меры разброса данных

Размах распределения – разность между максимальным и минимальным значением наблюдаемого параметра. Такая мера разброса данных является неустойчивой, т.к. зависит всего от двух значений. Поэтому для более точной оценки необходимо учитывать разность между каждым значением в выборке. Такой мерой разброса является дисперсия  $\sigma^2$ :

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n},$$

где  $\sigma^2$  – дисперсия;  $x_i$  – значение наблюдаемого параметра  $i$ -го экземпляра;  $\bar{x}$  – среднее значение параметра в выборке.

Наиболее часто на практике используют среднее квадратическое отклонение  $\sigma$ , т.к. оно выражается в тех же единицах, что и наблюдаемый параметр.

### 2.4. Нормальное распределение случайной величины

В практических задачах наиболее часто встречается нормальный закон распределения случайной величины. Нормальный закон возникает в случае, когда на наблюдаемый параметр оказывает влияние большое число независимо действующих случайных факторов, каждый из которых незначительно влияет на наблюдаемый параметр [20].

Особенностью нормального закона распределения является то, что он является предельным законом, к которому приближаются другие законы распределения при часто встречающихся типичных условиях.

На рис. 1 изображен график нормального закона распределения случайной величины. По оси абсцисс отложены значения случайной величины  $x$ , по оси ординат –  $f(x)$  – вероятность появления того или иного значения случайной величины (в процентах или долях единицы). Среднее, наиболее вероятное значение случайной величины – математическое ожидание  $M(x)$  соответствует максимуму кривой распределения, т.е. её «горбу». Ширина кривой распределения отражает изменчивость, варьирование случайной величины, что характеризуется дисперсией (средним квадратическим отклонением). Площади, расположенные под участками кривой распределения, показывают, какое количество случайных величин попадает в эти зоны.

По рисунку 1 видно, что почти все значения (а именно 99,72%) при нормальном законе распределения лежат в области от  $-3\sigma$  до  $3\sigma$ . Эта закономерность получила название «трёх сигм». Эти соотношения имеют особое значение для стандартизации или нормирования наблюдаемых показателей [19].

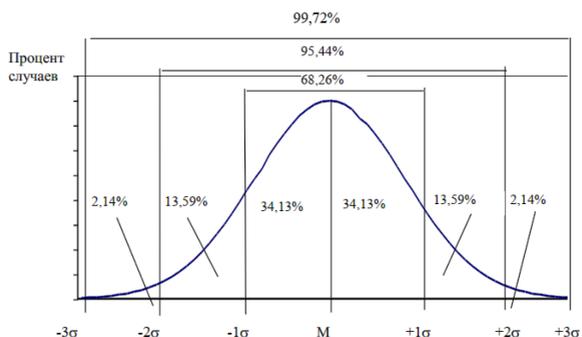


Рис. 1. График нормального закона распределения случайной величины

## **2.5. Шкалирование данных (нормализация и стандартизация данных)**

Нормализация – метод предобработки данных с целью приведения их к некоторой общей шкале без потери информации о различии диапазонов измерений [21]. Нормализация позволяет изменять диапазон данных без изменения формы распределения. На практике наиболее часто встречаются следующие методы нормализации данных: десятичное масштабирование, минимаксная нормализация, нормализация средним. В общем случае нормализация позволяет определить все значения выборки в диапазоне от 0 до 1 [22].

Стандартизация является более широким понятием, чем нормализация. Целью стандартизации является обеспечение возможности корректного сравнения значений наблюдений, собранных одними и теми же методами, но в различных условиях [21]. Стандартизация данных представляет собой предобработку данных с целью их приведения к определенному формату и представлению для корректного применения в последующем многомерном анализе. Стандартизация изменяет форму распределения данных, приводя её к нормальному распределению, причем большинство стандартизированных параметров получают путём линейного или нелинейного преобразования. Линейные преобразования используют в случае, если первичные значения параметров распределены по нормальному закону или близкому к нормальному. При этом соотношения между первичными значениями сохраняются, т.к. из каждого первичного значения вычитается одна и та же величина с последующим делением на другую постоянную величину. Стандартизация применяется в тех случаях, когда для используемого в дальнейшем многомерного анализа необходимо измерять расстояния между образцами.

Применительно к классификации изделий по заранее определенным характеристикам, стандартизация необходима для приве-

дения к единому масштабу измерений. Например, если классификация цифровых микросхем производится по критическому пониженному питающему напряжению и времени задержки распространения сигнала, то измеренный первый параметр в вольтах и второй, измеренный в наносекундах, имеют разный масштаб измерений. Кроме того, наблюдения, собранные в различных условиях, могут происходить из различных вероятностных распределений с различными параметрами.

При стандартизации происходит формирование стандартизированной шкалы, причем такой, что среднее значение выборки равно 0, а среднее квадратическое отклонение равно 1. При использовании такой шкалы место каждого значения в наборе данных определено по отклонению от среднего значения в единицах стандартного отклонения. Для перехода к стандартизированной шкале ( $Z$ -оценка) необходимо определить новое значение параметра:

$$z_i = \frac{x_i - \bar{x}}{\sigma_x},$$

где  $z_i$  – значение параметра  $i$ -го экземпляра в стандартизированной шкале;  $x_i$  – значение наблюдаемого параметра  $i$ -го экземпляра;  $\bar{x}$  – среднее значение параметра в выборке;  $\sigma_x$  – среднее квадратическое отклонение.

## **3. ОСНОВНЫЕ ПОЛОЖЕНИЯ ТЕОРИИ РАСПОЗНАВАНИЯ ОБРАЗОВ**

### **3.1. Теория распознавания образов**

Образ – совокупность данных о реальном или абстрактном объекте (процессе, явлении), позволяющая выделять его из всего множества анализируемых данных и группировать с другими объектами в соответствии с требованиями решаемой задачи [23].

Распознавание образов – научная дисциплина, целью которой является классификация объектов по нескольким категориям или классам. Рассматриваемые объекты принято называть образами [24].

Теория распознавания образов – раздел информатики и смежных дисциплин, развивающий основы метода классификации и идентификации объектов (предметов, явлений, процессов, сигналов, ситуаций), которые характеризуются конечным набором некоторых свойств и признаков.

Распознавание образов является относительно новой интенсивно развивающейся областью математической кибернетики, методы которой успешно используются в экономике, медицине, технике, физике, биологии, социологии, психологии и других разделах науки, где требуется на основе полученных данных прийти к определенному решению. Распознавание образов является основой искусственного интеллекта [25].

Для успешного применения описываемой теории необходимо решить задачу подготовки априорной информации, т.к. качество классификации образов напрямую зависит от точности имеющихся сведений [26]. Задача предварительной обработки является до-

вольно сложной, не имеет простого математического описания и заключается в определении вектора признаков объекта, представляющего собой образ в признаковом пространстве.

Наибольший практический интерес для рассмотрения задач классификации представляют многомерные статистические исследования. Сложность объекта исследования и глубина анализа прямо пропорциональны размерности информационного поля [27].

Вектор признаков – набор характеристик объектов, необходимых для решения конкретной задачи распознавания.

Классификация основывается на прецедентах. Прецедентом является образ, правильная классификация которого известна. Такой образ принимается за образец при решении задач классификации.

### 3.2. Классификация и кластеризация

В зависимости от того, имеется ли прецедентная информация или она отсутствует, различают классификацию с обучением и без обучения.

Классификация с обучением (с учителем) – классификация на основе имеющейся прецедентной информации. Например, имеется  $n$  векторов, классификация которых известна, и существует  $(n+1)$ -ый вектор, класс которого необходимо определить. В таком случае  $n$  векторов называются обучающими, а  $(n+1)$ -ый вектор – испытуемым или классифицируемым.

Классификация без обучения (без учителя, самообучение, кластеризация) – классификация в случае отсутствия информации о правильной классификации образов. В этом случае ставится задача разделения образов на классы по сходству соответствующих векторов признаков. Отличие от приведенного ранее примера заключается в том, что классификация  $n$  векторов неизвестна.

Одной из ведущих теорий в области распознавания образов является кластерный анализ, благодаря которому решение задач классификации было осуществлено несложными компьютерными методами, а также были получены легко интерпретируемые результаты [27].

### 3.3. Пример использования теории распознавания образов для классификации ЭРИ

Пусть существует система классификации резисторов, структурная схема которой показана на рис. 2. Задачей распознавания является отнесение каждого экземпляра резистора из всей исследуемой совокупности к определенному классу: годных  $K_1$  и потенциально дефектных изделий  $K_2$ .

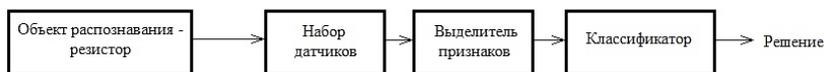


Рис. 2. Структурная схема классификации резисторов по качеству и надёжности

Первоначально принято решение определять качество резисторов по набору параметров: температурному коэффициенту сопротивления (ТКС)  $x_1$ , рассеиваемой мощности  $x_2$ , процентному изменению сопротивления при испытаниях  $x_3$ , эквивалентному напряжению шума  $x_4$ , генерируемого резистором. Обеспечение условий измерения и непосредственно само измерение параметров возложено на измерительные схемы с входящими в них датчиками. Выделитель признаков при этом уменьшает объем данных, предназначенный для корректной классификации резисторов. Затем измеренные значения признаков поступают на классификатор,

который производит оценку представленных данных и формулирует решение относительно качества объекта. Таким образом, на классификатор поступает вектор признаков:

$$x = \{x_1, x_2, x_3, x_4\}.$$

Задачей является разделение пространства признаков на две области, все точки одной из которых соответствуют классу K1, а другой – классу K2. В приведенном примере такое пространство признаков будет четырехмерным. Допустим, была получена разделяющая гиперплоскость, выше которой расположились образы экземпляров, относящиеся к классу K1, а ниже – к классу K2. Хотя разделение получено достаточно хорошим, нет никакой гарантии, что такое правило будет адекватным для новых экземпляров. Очевидно, что можно было бы взять новые экземпляры и проверить, насколько успешно правило работает для новых экземпляров. Таким образом, поставленная задача включает в себя элементы вероятности (т.е. статистики), а значит, может потребоваться поиск такой процедуры классификации, которая сделает вероятность ошибки минимальной [28].

## 4. ИСПОЛЬЗОВАНИЕ МЕТОДОВ КЛАСТЕРНОГО АНАЛИЗА ДЛЯ КЛАССИФИКАЦИИ ЭРИ

### 4.1. Кластерный анализ

Кластерный анализ – многомерная статистическая процедура, выполняющая сбор данных, содержащих информацию о выборке объектов, и затем упорядочивающая объекты в сравнительно однородные группы. Другими словами, кластерный анализ применяется для решения задачи классификации данных и выявления соответствующей структуры в ней.

Кластерный анализ получил распространение для решения 4 основных задач:

- разработка типологии или классификации (основное направление);
- исследование полезных концептуальных схем группирования объектов;
- порождение гипотез на основе исследования данных;
- проверка гипотез или исследования для определения, действительно ли типы (группы), выделенные тем или иным способом, присутствуют в имеющихся данных [29].

Достоинства кластерного анализа:

- позволяет производить разбиение объектов не по одному параметру, а по целому набору признаков. Кроме того, кластерный анализ в отличие от большинства математико-статистических методов не накладывает никаких ограничений на вид рассматриваемых объектов и позволяет рассматривать множество исходных данных практически произвольной природы.

- позволяет рассматривать достаточно большой объем информации и резко сокращать, сжимать большие массивы информации и делать их наглядными.

К недостаткам относят:

- зависимость состава и количества кластеров от выбираемых критериев разбиения;
- при сведении исходного массива данных к компактному виду теряются индивидуальные черты объектов и заменяются характеристиками обобщенных значений.

### ***Задача кластерного анализа***

На основании данных, содержащихся во множестве  $X$ , необходимо разбить множество объектов  $G$  на  $m$  кластеров (подмножеств)  $Q_1, Q_2, \dots, Q_m$  так, чтобы каждый объект  $G_j$  принадлежал одному и только одному подмножеству разбиения. Таким образом, объекты одного кластера будут сходными, а принадлежащие разным кластерам – разнородными [30].

Решением задачи кластерного анализа является разбиение, удовлетворяющее критерию оптимальности.

## **4.2. Характеристики кластеров**

Центр кластера (центроид) – среднее геометрическое место точек объектов кластера в пространстве признаков.

Радиус кластера – максимальное расстояние точек от центра кластера.

Объект является принадлежащим кластеру, если расстояние от объекта до центра кластера меньше радиуса кластера.

Размер кластера определяется по:

- радиусу кластера;

- среднему квадратическому отклонению расстояний от объектов до центра этого кластера.

Возможна ситуация, что у одного и того же объекта расстояние до двух кластеров менее радиусов этих кластеров. В таком случае объект является спорным, т.е. по мере сходства его можно отнести к разным кластерам.

Для проведения кластерного анализа над массивом данных необходимо:

- чтобы рассматриваемые признаки объекта в принципе допускали разбиение на кластеры;
- учитывать масштаб измерения признаков (см. п. 2.5).

### **4.3. Близость объектов (метод сходства, функции расстояний)**

Для корректного проведения кластерного анализа необходимо решить 2 задачи:

- определить единый масштаб измерений всех рассматриваемых признаков;
- определить метод оценки близости.

Первая задача решается подготовкой данных с помощью нормализации и/или стандартизации данных (см. п. 2.5).

Для решения второй задачи вводится понятие метрики. В этом случае наблюдения представляются точками в пространстве признаков, причем сходства и различия между точками находятся в соответствии с метрическими расстояниями между ними.

Для того, чтобы стать метрикой для меры сходства необходимо выполнение 4-х критериев [31]:

- симметрия. Даны 2 объекта  $a$  и  $b$ . Расстояние между ними удовлетворяет условию:

$$d(a, b) = d(b, a) \geq 0;$$

- неравенство треугольника. Даны три объекта  $a$ ,  $b$ ,  $c$ . Расстояние между ними удовлетворяет условию:

$$d(a, b) \leq d(a, c) + d(b, c);$$

- различимость нетождественных объектов. Если

$$d(a, b) \neq 0, \text{ то } a \neq b;$$

- неразличимость идентичных объектов. Если  $a$  и  $a'$  являются идентичными объектами, то:

$$d(a, a') = 0.$$

Мера близости  $\mu$  – величина, имеющая предел и возрастающая с возрастанием близости объектов, а также удовлетворяющая условиям:

- $\mu$  непрерывна,
- $\mu_{ab} = \mu_{ba}$ ,
- $1 \leq \mu_{ab} \leq 0$ .

$$\mu(a, b) = \frac{1}{1 + d(a, b)},$$

где  $a, b$  – объекты;  $m$  – количество признаков;  $x_i$  –  $i$ -ый признак.

В табл. 1 приведены основные методы сходства, наиболее часто применимые в методах кластерного анализа.

Евклидово расстояние – геометрическое расстояние в многомерном пространстве. Использование Евклидова расстояния способствует объединению объектов в шарообразные кластеры.

Квадрат Евклидова расстояния – используется для придания больших весов наиболее удаленным друг от друга объектам.

Обобщенное степенное расстояние (расстояние Минковского) – является интерпретацией универсальной метрики.

Расстояние Чебышева – используется для классификации в случае, если необходимо распознать два объекта как различные при значительном их отличии по одному признаку.

Таблица 1. Методы сходства для количественных шкал

Наименование	Формула
Евклидово расстояние	$d_{E ab} = \left( \sum_{l=1}^m (x_a^l - x_b^l)^2 \right)^{1/2}$
Квадрат Евклидова расстояния	$d_{E ab}^2 = \sum_{l=1}^m (x_a^l - x_b^l)^2$
Обобщенное степенное расстояние Минковского	$d_P ab = \left( \sum_{l=1}^m (x_a^l - x_b^l)^P \right)^{1/P}$
Расстояние Чебышева	$d_{ab} = \max_{1 \leq a, 1 \leq b}  x_a - x_b $
Манхэттенское расстояние (Расстояние городских кварталов)	$d_H(x_a, x_b) = \sum_{l=1}^k  x_a^l - x_b^l $

Манхэттенское расстояние (расстояние городских кварталов, «Хэммингово», «сити-блок») – расстояние вычисляется как среднее разностей по координатам. В большинстве случаев эта мера расстояния приводит к результатам, подобным при использовании Евклидова расстояния, но влияние отдельных выбросов меньше, т.к. отсутствует возведение в квадрат [27].

#### 4.4. Методы кластерного анализа

Существует несколько десятков теоретических алгоритмов кластерного анализа. Но далеко не все получили практическую реализацию в виде компьютерных алгоритмов. Кроме того, не все методы проработаны детально, также некоторые методы не позволяют получить удовлетворяющее исследователя разбиение. Поэтому такие методы и не получили широкого применения. Выбор алгоритма субъективен и во многом зависит от:

- математической культуры;
- структуры информационного поля, свойственного рассматриваемой задаче;
- имеющихся технических и программных средств реализации;
- возможности верификации результатов другими методами и т.д.

Из наиболее изученных и применяемых методов в различных областях науки выделяют следующие:

- иерархические агломеративные методы;
- иерархические дивизивные (дивазивные, дивизионные) методы;
- итеративные методы группировки;
- методы поиска модальных значений плотности;
- факторные методы;
- методы сгущений;
- методы, использующие теорию графов [31].

На рис. 3 приведена классификация методов кластерного анализа, получивших широкое практическое применение.

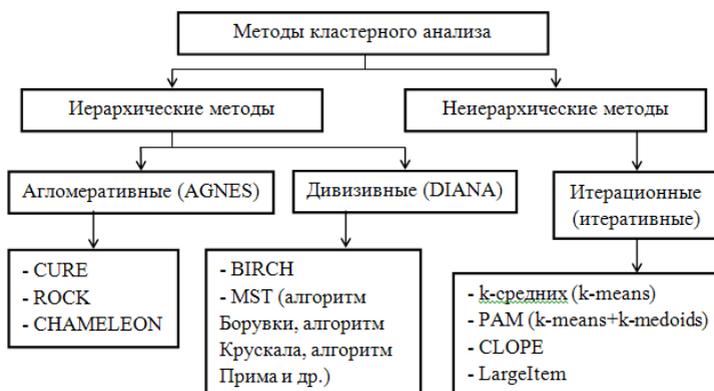


Рис. 3. Наиболее используемые методы кластерного анализа

## 4.5. Иерархические методы кластеризации

AGNES – Agglomerative Nesting – иерархические агломеративные. Такие методы характеризуются последовательным объединением исходных элементов и, соответственно, уменьшением числа кластеров.

В начале алгоритма все объекты рассматриваются в качестве самостоятельного кластера, состоящего из одного объекта. Далее алгоритм ищет наиболее похожие объекты и объединяет их в один кластер с пересчетом центроида и т.д. Алгоритм продолжается до тех пор, пока все объекты не объединятся в один кластер.

DIANA – Divisive Analysis – иерархические дивизивные методы. Такие методы характеризуются последовательным разделением единственного исходного кластера. Алгоритм подразумевает увеличение числа кластеров с каждым шагом. Таким образом, образуется последовательность расщепляющих групп.

*Преимущества:*

Иерархические методы кластеризации являются самыми распространенными в практике.

Возможность получения неперекрывающихся кластеров, являющихся вложенными на более высоком уровне сходства.

Иерархические методы кластеризации просты в освоении, т.к. зависят лишь от метода сходства и метода объединения.

*Недостатки:*

Процесс кластеризации медленный, т.к. объекты распределяются по кластерам лишь за один проход, поэтому если начальное разбиение произошло не качественно, то на последующих шагах его изменить нельзя.

Возможность порождения различных решений в результате простого переупорядочивания объектов, что влияет на устойчивость классификации. Особенно при малых выборках объектов.

Признаком иерархических методов кластеризации является графическая интерпретация полученных результатов – дендрограмма (древовидная диаграмма).

Иерархические методы используют правила объединения. Это такие правила, по которым происходит последовательное объединение объектов в кластеры и кластеров в более крупные кластеры.

#### **4.6. Метод объединения (правила объединения)**

*Метод одиночной связи* (метод ближайшего соседа). В этом случае объект-кандидат будет присоединен к уже существующему кластеру, если, по крайней мере, один из элементов кластера находится на том же уровне сходства, что и объект-кандидат. Другими словами, расстояние между кластерами определяется расстоянием между двумя наиболее близкими объектами (ближайшими соседями) в соседних кластерах. Применение такого метода имеет тенденцию к формированию «цепочечных» кластеров («цепной эффект»).

*Метод полной связи* (метод наиболее удаленных соседей). В этом случае сходство между объектом-кандидатом и объектом кластера не должно быть меньше некоторого порогового уровня. Другими словами, расстояния между кластерами определяются минимальным из наибольших расстояний между любыми двумя объектами в различных кластерах (т.е. «наиболее удаленными соседями»). Этот метод обычно работает очень хорошо, когда объекты происходят из реально различных «рощ». Если же кластеры имеют в некотором роде удлиненную форму или их тип является «цепочечным», то этот метод непригоден.

*Метод средней связи* (метод межгрупповой связи). В этом случае вычисляется среднее сходство объекта-кандидата на включение со всеми объектами существующего кластера, затем, если найденное среднее значение сходства достигает или превосходит

некоторый заданный пороговый уровень сходства, то объект включают в кластер. Другими словами, расстояние между двумя различными кластерами вычисляется как среднее расстояние между всеми парами объектов в них. Метод эффективен, когда объекты в действительности формируют различные «рощи», однако он работает одинаково хорошо и в случаях кластеров «цепочного» вида. То, что объединение кластеров в методе средней связи происходит при расстоянии большем, чем в методе одиночной связи, но меньшем, чем в методе полной связи, и объясняет промежуточное положение этого метода.

**Метод Уорда (Варда).** Метод направлен на минимизацию внутригрупповой суммы квадратов для любых двух гипотетических кластеров, которые могут быть сформированы на каждом шаге. Таким образом, метод выбирает такие объекты или кластеры для присоединения, которые дают минимальное приращение внутригрупповой суммы квадратов. Метод позволяет получать кластеры приблизительно равных размеров и имеющих гиперсферическую форму.

#### **4.7. Итеративные методы кластеризации**

При большом количестве наблюдений иерархические методы кластерного анализа обеспечивают громоздкие вычисления, и дендрограмма становится сложной для восприятия, что нивелирует все преимущества таких методов. В таких случаях используют неиерархические методы, основанные на разделении, которые представляют собой итеративные методы дробления исходной совокупности. В процессе деления новые кластеры формируются до тех пор, пока не будет выполнено правило остановки.

Такая неиерархическая кластеризация состоит в разделении набора данных на определенное количество отдельных кластеров.

Существует два подхода.

Первый заключается в определении границ кластеров как наиболее плотных участков в многомерном пространстве исходных данных, т.е. определение кластера там, где имеется большое «сгущение точек».

Второй подход заключается в минимизации меры различия объектов

#### **4.8. Алгоритм k-средних (k-means clustering)**

Наиболее распространен среди неиерархических методов алгоритм k-средних, также называемый быстрым кластерным анализом. Полное описание алгоритма можно найти в работе Хартигана и Вонга (Hartigan and Wong, 1978). В отличие от иерархических методов, которые не требуют предварительных предположений относительно числа кластеров, для возможности использования этого метода необходимо иметь гипотезу о наиболее вероятном количестве кластеров.

Алгоритм k-средних строит k кластеров, расположенных на возможно больших расстояниях друг от друга. Основной тип задач, которые решает алгоритм k-средних – наличие предположений (гипотез) относительно числа кластеров, при этом они должны быть различны настолько, насколько это возможно. Выбор числа k может базироваться на результатах предшествующих исследований, теоретических соображениях или интуиции.

Основная идея такого метода заключается в следующем. Задано фиксированное число k кластеров наблюдения сопоставляются кластерам так, что средние в кластере (для всех переменных) максимально возможно отличаются друг от друга.

Алгоритм k-средних заключается в следующем.

1. Первоначальное распределение объектов по кластерам. Выбирается число  $k$ , и на первом шаге эти точки считаются «центрами тяжести» (центроидами, т.е. по координатным средним кластеров) кластеров. Каждому кластеру соответствует один центр.

Выбор начальных центроидов может осуществляться следующим образом:

- выбор  $k$ -наблюдений для максимизации начального расстояния;
- случайный выбор  $k$ -наблюдений;
- выбор первых  $k$ -наблюдений.

В результате каждый образ объекта пространства признаков будет относиться к кластеру, центр тяжести которого ближе центров тяжести остальных кластеров.

2. Итеративный процесс. Так как каждый кластер теперь состоит из некоторого количества образов, то положение центра тяжести меняется. Вычисляются новые центры тяжести. Затем образы опять перераспределяются по кластерам.

Процесс вычисления центров и перераспределения образов продолжается до тех пор, пока не будет выполнено одно из условий:

- кластерные центры стабилизировались, т.е. все наблюдения принадлежат кластеру, которому принадлежали до текущей итерации (это называется сходимостью);
- число итераций равно максимальному числу итераций.

Выбор числа кластеров является довольно сложной задачей. Если нет предположений относительно этого числа, рекомендуют создать 2 кластера, затем 3, 4, 5 и т.д., сравнивая полученные результаты.

После получения результатов кластерного анализа методом  $k$ -средних следует проверить правильность кластеризации (т.е. оценить, насколько кластеры отличаются друг от друга). Для этого рассчитываются координаты центра для каждого кластера. При хорошей кластеризации должны быть получены сильно от-

личающиеся значения для всех измерений или хотя бы большей их части.

Достоинства алгоритма  $k$ -средних:

- простота использования (высокая скорость выполнения, эффективность при работе с большим количеством данных);
- быстрота использования;
- понятность и прозрачность алгоритма [32].

Недостатки алгоритма  $k$ -средних:

- алгоритм слишком чувствителен к выбросам, которые могут исказить среднее. Возможным решением этой проблемы является использование модификации алгоритма – алгоритм  $k$ -медианы;
- алгоритм может медленно работать на больших базах данных. Возможным решением данной проблемы является использование выборки данных.
- необходимость изменения алгоритма для работы с дискретными значениями данных.

Метод  $k$ -средних получил распространение как метод предварительного разбиения большого количества данных на группы, после чего для каждой группы проводится кластерный анализ более мощными методами.

#### **4.9. Блок-схема алгоритма иерархической кластеризации для классификации ЭРИ**

Блок-схема алгоритма представлена на рис. 4 и включает в себя несколько этапов:

- определение исходных данных по результатам обучающего эксперимента;
- выбор параметров алгоритмов кластерного анализа;
- преобразование исходных данных;

- построение и анализ дендрограмм;
- определение состава групп кластеров;
- определение точности классификации.

В качестве методов сходства и методов объединения необходимо выбирать такие варианты, которые подчеркивают разделение кластеров. Также при выборе метода сходства нужно учитывать, что рассматриваемые параметры ЭРИ одинаково значимы для классификации.

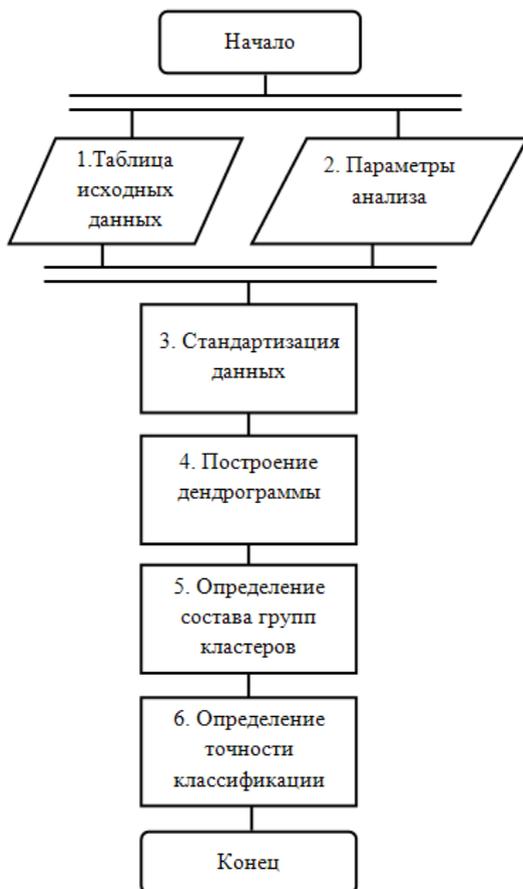


Рис. 4. Блок-схема алгоритма классификации с помощью иерархических методов кластеризации

#### 4.10. Блок-схема алгоритма k-средних для классификации ЭРИ

Блок-схема алгоритма k-средних для классификации ЭРИ представлена на рис. 5.

При задании исходных данных, требуют решения следующие вопросы:

- в качестве признаков использовать только информативные параметры или информативные и прогнозируемый параметры;
- определить количество изначально заданных кластерных групп;
- определить количество итераций процесса построения.

Если алгоритм используется в качестве проверочного, то имеет место в качестве исходных данных использовать информативные и прогнозируемый параметры, чем увеличивается информационное поле. Если целью методики является определение состава кластерных групп, к одной из которых в дальнейшем будет относиться экземпляр, то в качестве исходных данных необходимо использовать только информативные параметры.

Особенностью метода k-средних по сравнению с иерархическими методами кластеризации является указание определенного количества кластеров. Такие кластеры представляют собой центры, вокруг которых группируются объекты с наиболее близкими параметрами. Так как рассматривается разбиение выборки ЭРИ на класс годных и класс потенциально дефектных, то число кластеров определено в количестве двух.

Количество итераций построения зависит от количества наблюдений и признаков, т.е. от «сложности» исходных данных и определяется практическим путем.

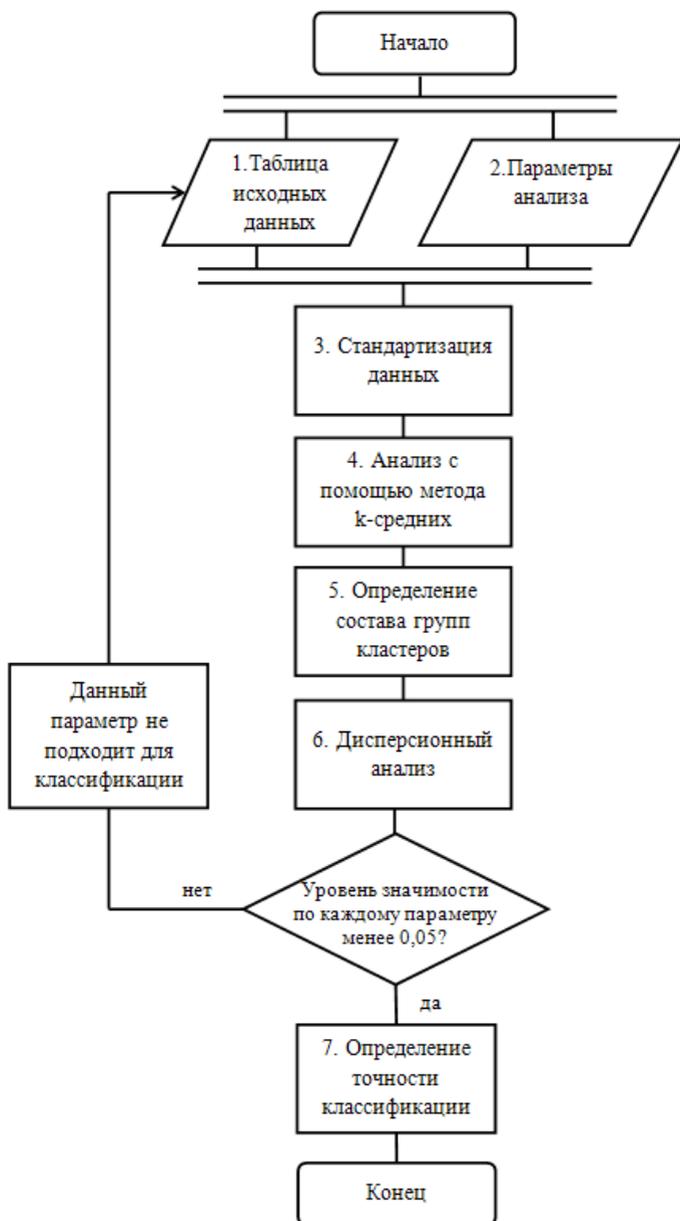


Рис. 5. Блок-схема алгоритма k-средних для классификации ЭРИ

## 5. ИСПОЛЬЗОВАНИЕ НЕЙРОСЕТЕВЫХ СТРУКТУР ДЛЯ КЛАССИФИКАЦИИ ЭРИ

### 5.1. Искусственные нейронные сети (ИНС) и их применение

Искусственные нейронные сети (ИНС) используются для решения сложных задач, которые требуют аналитических вычислений подобных тем, что делает человеческий мозг.

ИНС находят широкое применение в таких областях как моделирование, распознавание образов, обработка сигналов, управление благодаря одному важному свойству – способности обучаться на основе данных при участии учителя или без его вмешательства. Таким образом, ИНС являются гибким инструментом решения большого количества различных задач.

Самыми распространенными применениями ИНС является:

*Классификация* – распределение данных по параметрам. Например, существует набор людей и нужно решить, кому из них можно выдать кредит, а кому нет. Эту работу может сделать нейронная сеть, анализируя такую информацию как возраст, платежеспособность, кредитная история и т.д.

*Предсказание* – возможность предсказывать следующий шаг. Например, можно предсказать рост или падение акций, основываясь на ситуации на фондовом рынке.

*Распознавание* – в настоящее время, самое широкое применение нейронных сетей. Для решения этой задачи существует целое направление, называемое глубоким обучением (Deep Learning). Примером решения задачи распознавания может служить ассоциа-

тивный поиск информации, при котором пользователь формирует набор слов, а нейросеть выдает результат поиска. Также примером может служить распознавание текста, изображений, голоса, прокладка оптимального маршрута и др.

Более подробно про сферы применения ИНС написано в [Хаб].

Преимущества ИНС заключаются в:

- распараллеливании обработки информации;
- способности обучаться, т.е. создавать обобщения (generalization), т.е. получать обоснованный результат на основании данных, которые не встречались в процессе обучения [33].

Стоит упомянуть, что ИНС не могут обеспечить готовые решения, поэтому целесообразно их интегрировать в сложные системы. Другими словами комплексную задачу, требующую решения, можно разбить на несколько простых, часть из которых можно подвергнуть решению с помощью ИНС.

ИНС также не лишены недостатков. Во-первых, ИНС не способны давать точные и однозначные ответы. Всегда существует некоторая вероятность неправильного ответа. С другой стороны, не все задачи требуют точного ответа, довольно часто достаточно наиболее вероятного ответа. Во-вторых, т.к. ИНС состоит из нейронов, преобразующих сигнал по определенному закону, не зависят друг от друга, следует, что рассматриваемая задача решается в один заход, и результат не пересматривается этой же сетью. Таким образом, ИНС не способны решать задачу по шагам. В-третьих, ИНС не способны решать вычислительные задачи, т.к. это требует задания очередности использования нейронов сети и пошагового вычисления.

Исходя из преимуществ и недостатков ИНС, можно сделать вывод, что они являются хорошим, но дополнительным инструментом решения большого количества прикладных задач. Также можно сделать вывод, что ИНС являются хорошим инструментом

для решения задачи классификации ЭРИ при условии правильного определения набора исходных данных, построения структуры сети, оценки ошибки полученных результатов [34].

### Структура простой ИНС.

**Нейронная сеть** – это система соединенных и взаимодействующих между собой простых процессоров (искусственных нейронов).

Нейрон является простейшим элементом сети, который получает информацию, производит над ней вычисления и передает ее дальше.

Нейроны бывают 3 видов: входными (input); скрытыми (hide); выходными (output);

Совокупность нейронов на одном уровне образуют слой.

Нейрон характеризуется 2 параметрами:

- входными данными;
- выходными данными.

Нейроны оперируют с числами из диапазона  $[0,1]$  или  $[-1,1]$ .

На рис. 6 приведена структура простой ИНС с одним скрытым слоем, состоящим из 4 нейронов.

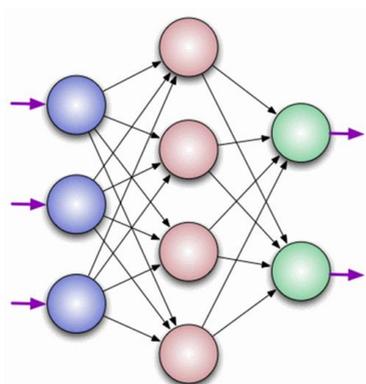


Рис. 6. Модель ИНС

Существует два подхода к машинному обучению:

- обучение без учителя;
- обучение с учителем.

При первом подходе на вход сети подается вектор входных данных, но сеть не получает пояснений по этим данным. В таком случае сеть лишь выделяет статистические закономерности.

При использовании второго подхода для некоторых (или всех) векторов входных данных известен выходной результат, который мы хотим получить. В этом случае задачей является настройка сети таким образом, чтобы для обучающей выборки выявить закономерности, связывающие входные и выходные данные.

## Модели нейронов

**Нейрон** – единица обработки информации в нейронной сети.

Нелинейная модель нейрона изображена на рис. 7.

В каждом нейроне можно выделить 3 составляющие:

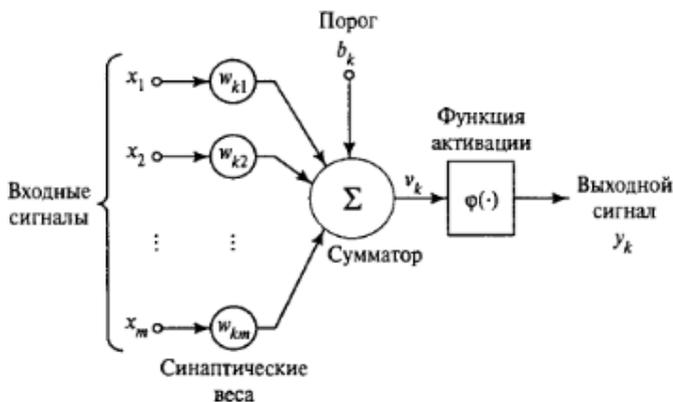


Рис. 7. Модель нелинейного нейрона

- набор синапсов (или связей), каждый из которых характеризуется своим весом. Синаптический вес может принимать как положительные, так и отрицательные значения;

- сумматор, который складывает входные сигналы, взвешенные относительно соответствующих синапсов нейрона (линейная комбинация);
- функция активации (сжатия), которая ограничивает амплитуду выходного сигнала. Значения выходного сигнала обычно лежат в диапазоне  $[0,1]$  или  $[-1,1]$ .

Также на схеме присутствует пороговый элемент  $b_k$  (bias), предназначенный для отражения информации о входном сигнале (увеличении или уменьшении), подаваемого на функцию активации. Порог  $b_k$  является внешним параметром искусственного нейрона.

Функционирование нейрона можно описать математически:

$$u_k = \sum_{j=1}^m w_{kj} x_j,$$

$$v_k = u_k + b_k,$$

$$y_k = \varphi(v_k),$$

где  $x_1, x_2, \dots, x_m$  – входные сигналы;  $w_{k1}, w_{k2}, \dots, w_{km}$  – синаптические веса нейрона  $k$ ;  $u_k$  – линейная комбинация входных воздействий;  $v_k$  – индуцированное локальное поле (потенциал активации);  $b_k$  – порог;  $\varphi$  – функция активации;  $y_k$  – выходной сигнал нейрона.

Рассмотрим пример функционирования нейронной сети со структурой, показанной на рис. 8.

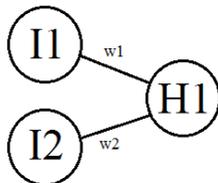


Рис. 8. Простая нейронная сеть

Математически опишем процессы, происходящие в такой сети:

$$u = x_1 w_1 + x_2 w_2;$$

$$v = u + b; \quad b = 0;$$

$$y = \varphi(v).$$

Функции активации [33, 35]:

1. Функция единичного скачка (пороговая функция, функция Хэвисайда, вид функции представлен на рис. 9).

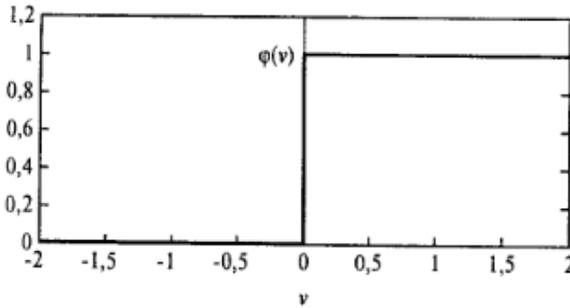


Рис. 9. Функция единичного скачка

$$\varphi(v) = \begin{cases} 1, & \text{если } v \geq 0; \\ 0, & \text{если } v < 0. \end{cases}$$

Тогда выходной сигнал с такого нейрона:

$$y_k = \begin{cases} 1, & \text{если } v_k \geq 0; \\ 0, & \text{если } v_k < 0, \end{cases}$$

где  $v_k$  – это индуцированное локальное поле нейрона, т.е.

$$v_k = \sum_{j=1}^m w_{kj} x_j + b_k.$$

2. Кусочно-линейная функция (вид функции представлен на рис. 10).

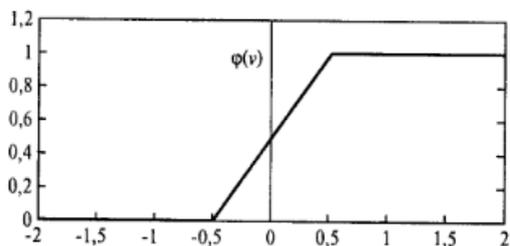


Рис. 10. Кусочно-линейная функция

$$\varphi(v) = \begin{cases} 1, & \text{если } v \geq +1/2; \\ |v|, & \text{если } +\frac{1}{2} > v > -1/2; \\ 0, & \text{если } v \leq -1/2. \end{cases}$$

Такую функцию называют функцией аппроксимации нелинейного усилителя.

3. Сигмоидальная функция (вид функции представлен на рис. 11).

Самая распространенная функция. Ее особенность заключается в стремительном росте при поддержке баланса между линейным и нелинейным поведением.

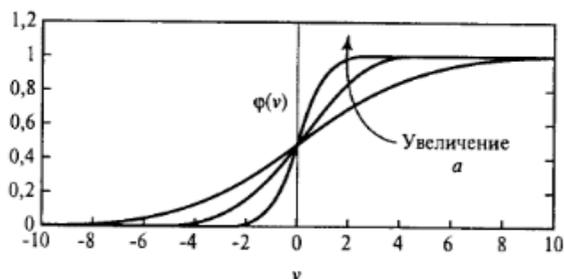


Рис. 11. Сигмоидальная функция

$$\varphi(v) = \frac{1}{1 + \exp(-av)},$$

где  $a$  – параметр наклона, от которой зависит крутизна функции.

Такая функция является дифференцируемой, что важно в теории ИНС.

У рассмотренных функций активаций диапазон значений  $[0,1]$ . Иногда требуется рассмотрение диапазона значений функции активации  $[-1,1]$ . В таком случае используют сигнум-функцию.

4. Сигнум. Функция симметрична относительно начала координат и имеет форму гиперболического тангенса (вид функции представлен на рис. 12).

$$\varphi(v) = \begin{cases} 1, & \text{если } v > 0; \\ 0, & \text{если } v = 0; \\ -1, & \text{если } v < 0. \end{cases}$$

$$\varphi(v) = \tanh(v).$$

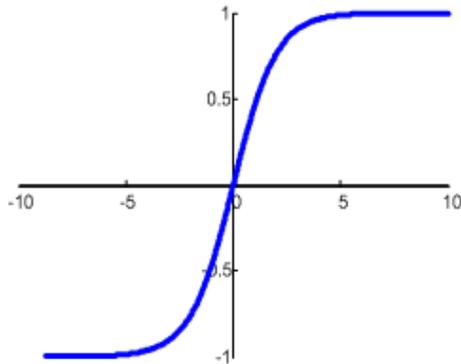


Рис. 12. Сигнум-функция

Ошибка сети вычисляется 3 способами:

– среднеквадратическая ошибка (MSE – Mean Squared Error):

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2;$$

– корень среднеквадратической ошибки (RMSE – Root Mean Squared Error):

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2};$$

– Arctan:

$$ARCTAN = \frac{1}{n} \sum_{i=1}^n \arctan^2(Y_i - \hat{Y}_i).$$

## **5.2. Использование перцептрона для решения задачи классификации ЭРИ**

Перцептрон (англ. Perceptron – восприятие) – математическая или компьютерная модель восприятия информации мозгом (кибернетическая модель мозга).

Перцептрон представляет собой простейшую форму нейронной сети, предназначенную для классификации линейно-разделимых образов. Было доказано, что если образы (векторы), используемые для обучения перцептрона, выбраны из двух линейно-разделимых классов, то алгоритм перцептрона сходится и формирует поверхность решений в форме гиперплоскости, разделяющей эти два класса [33].

Перцептрон состоит из одного или нескольких слоёв скрытых нейронов с настраиваемыми синаптическими весами и порогами. В первом случае перцептрон называют однослойным, во втором – многослойным.

Перцептрон, построенный на одном нейроне, решает задачу разделения совокупности на два класса. Увеличивая размерность внутреннего слоя перцептрона и включая в него несколько нейронов, увеличивается количество классов, на которые можно разделить совокупность.

Рассмотрим однослойный перцептрон с 2 входами, задачей которого является классификация совокупности на два класса. В качестве функции активации используется жесткая пороговая функция. Эквивалентный граф такого перцептрона представлен на рис. 13.

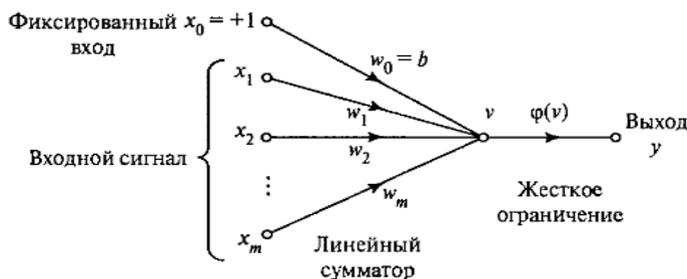


Рис. 13. Эквивалентный граф однослойного перцептрона

Тогда областью решений является двумерная плоскость с разделяющей прямой (см. рис. 14).

Для корректного функционирования перцептрона два класса должны быть линейно-разделимыми, т.е. образы, принадлежащие различным классам, должны быть значительно удалены друг от друга. В случае нелинейно-разделимых классов перцептрон работает с ошибками (рис. 15).

Многослойный перцептрон состоит из одного входного слоя (множество сенсорных элементов), одного или нескольких вычислительных слоёв (скрытые слои) и одного выходного слоя.

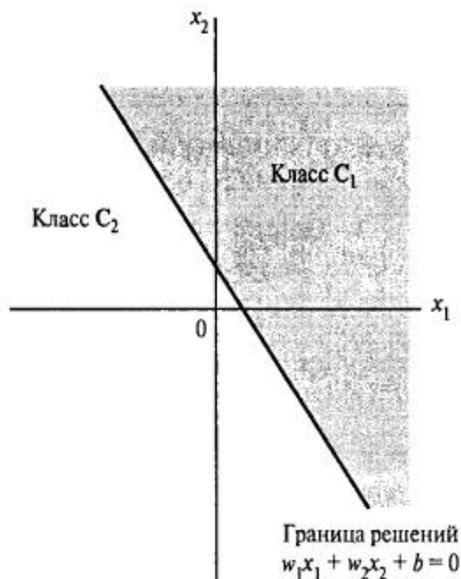


Рис. 14. Области решений при классификации

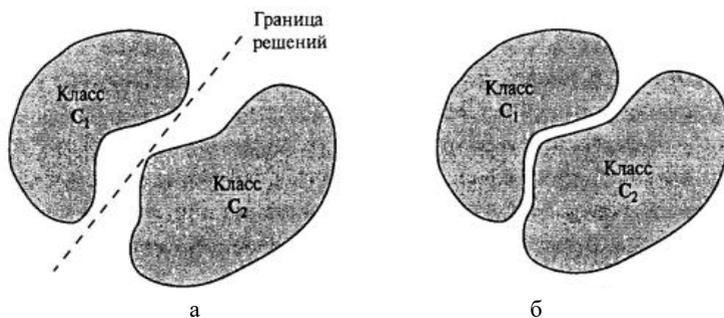


Рис. 15. Линейно-разделимые и нелинейно-разделимые образы

В отличие от однослойного персептрона в многослойном персептроне:

- каждый нейрон имеет нелинейную функцию активации (чаще всего, сигмоидальную);

- сеть имеет один или несколько слоев скрытых нейронов. Эти нейроны позволяют сети обучаться последовательно, извлекая важные признаки из вектора входных данных;
- высоким уровнем связности, зависящим от синаптических весов. Изменение уровня связности происходит при изменении большого количества весовых коэффициентов.

Однослойный перцептрон не может классифицировать линейно-неразделимые входные образы, что успешно решает многослойный перцептрон.

Многослойный перцептрон способен выполнять более общие классификации, отделяя те образы, которые содержатся в выпуклых ограниченных или неограниченных областях. Область называется выпуклой, если для любых двух ее точек соединяющий их отрезок целиком лежит в области. Область называется ограниченной, если ее можно заключить в некоторый круг. Неограниченную область невозможно заключить внутрь круга (например, область между двумя параллельными линиями) [36].

Для того, чтобы перцептрон начал функционировать, его нужно обучить на прецедентных примерах. В процессе обучения происходит настройка значений весов  $w$  и порогов  $b$ . Обучение перцептрона является обучением с учителем [33, 36]. Наиболее часто перцептрон обучают методом коррекции ошибки с использованием дельта-правила [37].

Пример использования перцептрона с различной конфигурацией скрытого слоя для классификации ЭРИ приведен в работе [38].

### **5.3. Самоорганизующиеся системы (карты самоорганизации) и их применение для решения задачи классификации ЭРИ**

Самоорганизующиеся карты могут использоваться для решения таких задач, как моделирование, прогнозирование, поиск за-

кономерностей в больших массивах данных, выявление наборов независимых признаков и сжатие информации.

Самоорганизующиеся системы основаны на конкурентном обучении. Отдельные нейроны выходного слоя соревнуются друг с другом за право активации, в результате чего активным становится один нейрон из всего выходного слоя [33].

В картах самоорганизации нейроны располагаются в узлах одномерной или двухмерной решетки. Карты большей размерности используются реже. В ходе обучения сети и конкуренции нейронов, нейроны-победители избирательно настраиваются на различные входные образы. Положение таких нейронов упорядочивается по отношению друг к другу так, что на решетке создается значимая система координат. Таким образом формируются топографические карты входных образов. В полученных картах пространственное положение выходных нейронов соответствует конкретной области признаков данных.

Структура самоорганизующихся сетей схожа со строением человеческого мозга. В нем отдельные сенсорные входы топологически упорядочены в определенных областях церебральной коры мозга. Другими словами, нейроны, работающие с близко расположенными областями информации (в нашем случае, более схожими входными векторами признаков), также расположены близко друг к другу и взаимодействуют друг с другом посредством коротких синаптических связей.

Основной задачей самоорганизующихся карт является преобразование поступающих векторов сигналов, имеющих произвольную размерность, в одно- или двухмерную дискретную карту.

На рис. 16 изображена схематическая диаграмма двумерной решетки выходного слоя.

Все нейроны на выходном слое связаны со всеми нейронами входного слоя. Эта сеть имеет структуру прямого распространения с одним вычислительным слоем, состоящим из нейронов, упорядоченных в столбцы и в строки.

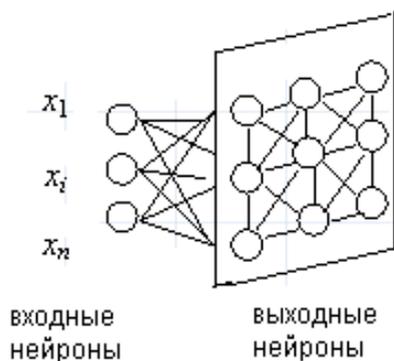


Рис. 16. Схема самоорганизующейся сети

Алгоритм, ответственный за формирование самоорганизующихся карт, начинается с инициализации синаптических весов сети. Изначально задаются случайные малые веса, а карта признаков не имеет какого-либо порядка. Затем при обучении запускаются 3 процесса:

- конкуренция. Для каждого входного образа нейроны сети вычисляют относительные значения дискриминантной функции. Эта функция является основой конкуренции среди нейронов;
- кооперация. Победивший нейрон определяет пространственное положение топологической окрестности нейронов, обеспечивая тем самым базис для кооперации между этими нейронами;
- синаптическая адаптация. Этот механизм позволяет возбужденным нейронам увеличивать собственные значения дискриминантных функций по отношению к входным образам посредством соответствующих корректировок синаптических весов. Коррекция производится таким образом, чтобы отклик нейрона-победителя на последующее применение аналогичных примеров усиливался.

Представителем самоорганизующихся карт является модель Кохонена, схема которой представлена на рис. 17.

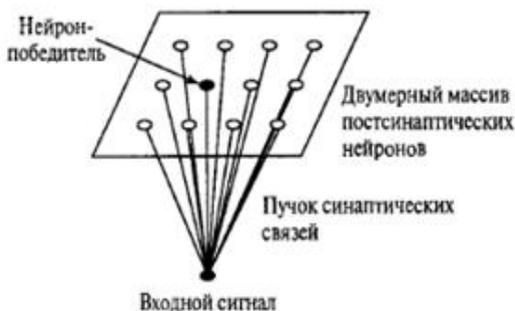


Рис. 17. Модель самоорганизующейся сети Кохонена

Наиболее распространенное применение сетей Кохонена – решение задачи классификации без учителя, т.е. кластеризации.

Напомним, что при такой постановке задачи нам дан набор объектов, каждому из которых сопоставлена строка таблицы (вектор значений признаков). Требуется разбить исходное множество на классы, т.е. для каждого объекта найти класс, к которому он принадлежит.

В результате получения новой информации о классах возможна коррекция существующих правил классификации объектов.

Также сети Кохонена используют для разведочного анализа данных и обнаружения новых явлений.

Разведочный анализ данных. Сеть Кохонена способна распознавать кластеры в данных, а также устанавливать близость классов. Таким образом, пользователь может улучшить свое понимание структуры данных, чтобы затем уточнить нейросетевую модель. Если в данных распознаны классы, то их можно обозначить, после чего сеть сможет решать задачи классификации. Сети Кохонена можно использовать и в тех задачах классификации, где

классы уже заданы. Тогда преимущество будет в том, что сеть сможет выявить сходство между различными классами.

Обнаружение новых явлений. Сеть Кохонена распознает кластеры в обучающих данных и относит все данные к тем или иным кластерам. Если после этого сеть встретится с набором данных, непохожим ни на один из известных образцов, то она не сможет классифицировать такой набор и тем самым выявит его новизну.

Пример классификации ЭРИ на класс годных и потенциально дефектных с помощью сетей Кохонена приведен в работе [39].

## БИБЛИОГРАФИЧЕСКИЙ СПИСОК

1. *ГОСТ Р 56526-2015*. Требования надежности и безопасности космических систем, комплексов и автоматических космических аппаратов единичного (мелкосерийного) изготовления с длительными сроками активного существования. – Москва: Стандартинформ, 2016. – 46 с.
2. *ГОСТ Р 56649-2015*. Техника ракетно-космическая. Электронная компонентная база иностранного производства. Порядок применения. – Москва: Стандартинформ, 2016. – 56 с.
3. Левин, Р.Г. Некоторые проблемы обеспечения электронной компонентной базой в процессе жизненного цикла образцов специальной техники / Р.Г. Левин, А.В. Уханов // Анализ и прогнозирование систем управления: тр. IX Междунар. науч.-практич. конф. молодых ученых, студентов и аспирантов. Ч. II. – Санкт-Петербург: СЗТУ, 2008. – 287 с.
4. Данилин, Н. Проектирование и разработка космических бортовых приборов, ориентированных на современную зарубежную электронную компонентную базу / Н. Данилин, С. Белослудцев // Современная электроника. – 2008. – № 4. – С. 54–59.
5. Урличич, Ю. Дополнительные отбраковочные испытания современной космической электронной компонентной базы / Ю. Урличич, Н. Данилин, Д. Чернов [и др.] // Современная электроника. – 2007. – № 2. – С. 8–11.
6. Герасимов, О.Н. Способ организации производственного контроля и диагностики РЭС с заданным уровнем остаточного ресурса / О.Н. Герасимов, А.В. Затылкин, Н.К. Юрков // Надежность и качество сложных систем. – 2016. – № 1(13). – С. 94–98.

7. Байда, Н.К. Эволюция отказоустойчивых БЦВК и направления их развития на однокристалльных микро-ЭВМ / Н.К. Байда, А.И. Кривонос, И.В. Лысенко [и др.] // Системы обработки информации. – 2001. – Вып. 4(14). – С. 217–225.
8. LaBel, K.A. Commercial Microelectronics Technologies for Applications in the Satellite Radiation Environment / K.A. LaBel, M.M. Gates, A.K. Moran. – URL: <http://radhome.gsfc.nasa.gov/radhome/papers/aspen.htm> (дата обращения: 26.03.2021).
9. Howard, J. Synopsys V1.3 Proton Dose and Single Event Effects Testing of the Intel Pentium III (P3) and AMD K7 Microprocessors / J. Howard, E. Webb, K. LaBel, M. Carts, R. Stattel, C. Rogers. URL: <http://radhome.gsfc.nasa.gov/radhome/papers/i062100.pdf> (дата обращения: 26.03.2021).
10. Powell, D. GUARDS: a generic upgradable architecture for real-time dependable systems / D. Powell, J. Arlat, L. Beus-Dukic, A. Bondavalli, P. Coppola, A. Fantechi, E. Jenn, C. Rebejac, A. Wellings // IEEE Transactions on Parallel and Distributed Systems, 1999. Vol. 10, Issue 6. P. 580-599.
11. Булаев, И.Ю. Методы и средства обнаружения скрытых дефектов КМОП-микросхем / И.Ю. Булаев // Ракетно-космическое приборостроение и информационные системы. – 2015. – Т. 2, Вып. 3. – С. 88–91.
12. Ушакова, К.О. Риски и угрозы при решении вопросов импортозамещения в авиационно-космическом сегменте военно-промышленного комплекса / К.О. Ушакова, В.Г. Исаев // Информационно-технологический вестник. – 2019. – № 3(21). – С. 3–14.
13. Севастьянов, Н.Н. Основы управления надёжностью космических аппаратов с длительными сроками эксплуатации / Н.Н. Севастьянов, А.И. Андреев; под общ. ред. Н.Н. Севастьянова. – Томск: Издательский Дом ТГУ, 2015. – 266 с.

14. Чуев, Ю.В. Прогнозирование количественных характеристик процессов / Ю.В. Чуев, Ю.Б. Михайлов, В.И. Кузьмин. – Москва: Советское Радио, 1975. – 400 с.
15. ГОСТ 27.301-95. Надежность в технике. Расчет надежности. Основные положения. – Москва: Издательство стандартов, 1996. – 16 с.
16. Кремер, Н.Ш. Теория вероятностей и математическая статистика: учебник для студентов вузов, обучающихся по экономическим специальностям. – 3-е изд., перераб. и доп. – Москва: ЮНИТИ-ДАНА, 2010. – 551 с.
17. Полякова, В.В. Основы теории статистики: учебное пособие / В.В. Полякова, Н.В. Шаброва. – 2-е изд., испр. и доп. – Екатеринбург: Изд-во Урал. ун-та, 2015. – 148 с.
18. Кириллов, А.В. Статистика. Ч. 1. Общая теория статистики: учебное пособие / А.В. Кириллов. – Самара: Изд-во СГАУ, 2012. – 112 с.
19. Рукавишникова, Н.Г. Статистический анализ данных и способы представления результатов исследования: учебно-методическое пособие к курсам «Экспериментальная психология» и «Психодиагностика» / Н.Г. Рукавишникова, Е.Г. Заверткина. – Ярославль: Изд-во ЯГПУ, 2000. – 47 с.
20. Ивченко, Г.И. Введение в математическую статистику: учебник / Г.И. Ивченко, Ю.И. Медведев. – Москва: Издательство ЛКИ, 2010. – 600 с.
21. Энциклопедия по бизнес-анализу. Стандартизация данных (Data standardization). – URL: <https://wiki.loginom.ru/articles/data-standartization.html> (дата обращения: 28.03.2021).
22. Корпоративные курсы Python в Big Data и Machine Learning. 4 шага к моделированию Machine Learning: практические примеры на Python. – URL: <https://python-school.ru/preprocessing-in-ml-4-steps> (дата обращения: 28.03.2021).

23. Чабан, Л.Н. Методы и алгоритмы распознавания образов в автоматизированном дешифрировании данных дистанционного зондирования: учебное пособие / Л.Н. Чабан. – Москва: МИИ-ГАиК, 2016. – 94 с.
24. Местецкий М.Л. Математические методы распознавания образов. Курс лекций / М.Л. Местецкий. – Москва: Изд-во МГУ, 2004. – 85 с.
25. Мазуров, В.Д. Математические методы распознавания образов: учебное пособие / В.Д. Мазуров. – 2-е изд., доп. и перераб. – Екатеринбург: Изд-во Урал. ун-та, 2010. – 101 с.
26. Патрик, Э. Основы теории распознавания образов / под ред. Б.Р. Левина; пер с англ. – Москва: Сов. Радио, 1980. – 408 с.
27. Гитис, Л.Х. Статистическая классификация и кластерный анализ / Л.Х. Гитис. – Москва: Издательство Московского государственного горного университета, 2003. – 157 с.
28. Дуда, Р. Распознавание образов и анализ сцен / Р. Дуда, П. Харт; под ред. В.Л. Стефанюка; пер. с англ. – Москва: Изд-во «МИР», 1976. – 511 с.
29. Национальный открытый университет «ИНТУИТ». Курс Data Mining. Лекция 13. Методы кластерного анализа. Иерархические методы. – URL: <https://intuit.ru/studies/courses/6/6/lecture/182> (дата обращения: 28.03.2021).
30. Дюран, Б. Кластерный анализ / Б. Дюран, П. Оделл; под ред. А.Я. Боярского; пер. с англ. Е.З. Демиденко. – Москва: Статистика, 1977. – 128 с.
31. Факторный, дискриминантный и кластерный анализ / Дж.-О. Ким, Ч.У. Мьюллер, У.Р. Клекка [и др.]; под ред. И.С. Енюкова; пер. с англ. – Москва: Финансы и статистика, 1989. – 2015 с.
32. Логинов, П.С. Применение метода  $k$ -средних и диаграмм Вороного для кластерного анализа базовых станций в телекоммуникациях / П.С. Логинов // Перспективы науки. – 2016. – № 2(77). – С. 59–63.

33. Хайкин, С. Нейронные сети: полный курс / Хайкин С.; пер. с англ. – 2-е изд. – Москва: Издание «Вильямс», 2006. – 1104 с.
34. Сообщество IT-специалистов. Нейронные сети: практическое применение. – URL: <https://habr.com/ru/post/322392/> (дата обращения: 2.04.2021).
35. Гафаров, Ф.М. Искусственные нейронные сети и их приложения: учебное пособие / Ф.М. Гафаров, А.Ф. Галимянов. – Казань: Изд-во Казан. ун-та, 2018. – 121 с.
36. Уоссермен, Ф. Нейрокомпьютерная техника: теория и практика / Ф. Уоссермен; пер. с англ. Ю.А. Зуева, В.А. Точенова. – Москва: Мир, 1992. – 237 с.
37. Портал искусственного интеллекта. Обучение персептрона. Дельта-правило. – URL: <http://www.aiportal.ru/articles/neural-networks/perceptron-learning.html> (дата обращения: 3.04.2021).
38. Мишанов, Р.О. Использование однослойного персептрона для решения задачи классификации электрорадиоизделий с целью повышения качества и надёжности бортовой аппаратуры / Р.О. Мишанов // Надёжность и качество сложных систем. – 2020. – № 2(30). – С. 106–114.
39. Мишанов, Р.О. Применение самоорганизующихся карт Кохонена для классификации электрорадиоизделий и повышения надёжности бортовой аппаратуры / Р.О. Мишанов // Сборник трудов ИТНТ-2018, 24–27 апреля 2018 г. – Самара: Новая техника, 2018. – С. 2311–2318.

**ПРИМЕР РЕШЕНИЯ ЗАДАЧИ КЛАССИФИКАЦИИ  
МЕТОДОМ ИЕРАРХИЧЕСКОЙ КЛАСТЕРИЗАЦИИ**

**Задача:** дана партия стабилитронов, для описания свойств которой осуществили выборку в объеме 8 экземпляров. Информативные параметры экземпляров выборки представлены в табл. А1. Необходимо определить состав двух кластеров так, чтобы эти кластеры были в наибольшей степени различны.

Для кластеризации использовать метод иерархической кластеризации, за меру расстояния принять Евклидово расстояние, за правило объединения – метод одиночной связи (ближайшего соседа).

Таблица А1. Значения информативных параметров  
экземпляров выборки

№ экз.	1	2	3	4	5	6	7	8
X1	15	18	8	5	9	10	15	8
X2	3	9	2	6	4	10	7	8

Решение задачи включает в себя выполнение следующих этапов.

1. Изобразим расположение образов экземпляров на двухмерной координатной плоскости (рис. А1). По вертикальной оси отложим значения параметра X1, по горизонтальной – значения параметра X2. Число возле изображения образа объекта – номер экземпляра.

2. Для подготовки данных и приведения значения параметров к единому масштабу измерений проведем стандартизацию выборки по следующим формулам:

$$X_{cp} = \frac{\sum_{i=1}^n X_i}{n},$$

$$Z_i = \frac{X - X_{\text{cp}}}{\sigma},$$

$$\sigma = \sqrt{D} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - X_{\text{cp}})^2},$$

где  $X_{\text{cp}}$  – среднее значение параметра  $X$ ;

$X_i$  – значение параметра  $i$ -го экземпляра;

$n$  – объем выборки;

$\sigma$  – среднее квадратичное отклонение;

$Z_i$  – преобразованное значение параметра  $X$   $i$ -го экземпляра;

$D$  – дисперсия.

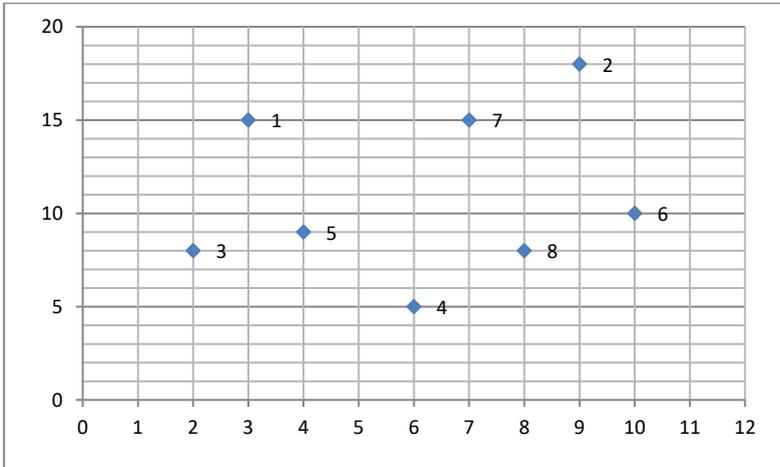


Рис. А1. Изображение образов объектов на двухмерной координатной плоскости

Таким образом, получим следующие таблицы для параметров  $X_1$  и  $X_2$  (табл. А2, табл. А3). Округление значений произведем до 4 знака после запятой.

По результатам вычислений сформируем сводную таблицу преобразованных данных (табл. А4) и изобразим образы объектов на двумерной координатной плоскости, по вертикальной оси которой отложим значения преобразованного параметра Z1, по горизонтальной – значения параметра Z2 (рис. А2). Число возле изображения образа объекта – номер экземпляра.

Таблица А2. Таблица преобразования данных параметра X1

№ экз.	X1	X1 – X <sub>cp</sub>	(X1 – X <sub>cp</sub> ) <sup>2</sup>	Z1
1	15	4	16	0,8944
2	18	7	49	1,5652
3	8	-3	9	-0,6708
4	5	-6	36	-1,3416
5	9	-2	4	-0,4472
6	10	-1	1	-0,2236
7	15	4	16	0,8944
8	8	-3	9	-0,6708
X <sub>cp</sub>	11	$\sigma$	4,4721	

Таблица А3. Таблица преобразования данных параметра X2

№ экз.	X2	X2 – X <sub>cp</sub>	(X2 – X <sub>cp</sub> ) <sup>2</sup>	Z2
1	3	-3,125	9,7656	-1,0775
2	9	2,875	8,2656	0,9913
3	2	-4,125	17,0156	-1,4224
4	6	-0,125	0,0156	-0,0431
5	4	-2,125	4,5156	-0,7327
6	10	3,875	15,0156	1,3362
7	7	0,875	0,7656	0,3017
8	8	1,875	3,5156	0,6465
X <sub>cp</sub>	6,125	$\sigma$	2,9001	

Таблица А4. Сводная таблица значений преобразованных данных

№ экз.	1	2	3	4	5	6	7	8
<b>Z1</b>	0,8944	1,5652	-0,6708	-1,3416	-0,4472	-0,2236	0,8944	-0,6708
<b>Z2</b>	-1,0775	0,9913	-1,4224	-0,0431	-0,7327	1,3362	0,3017	0,6465

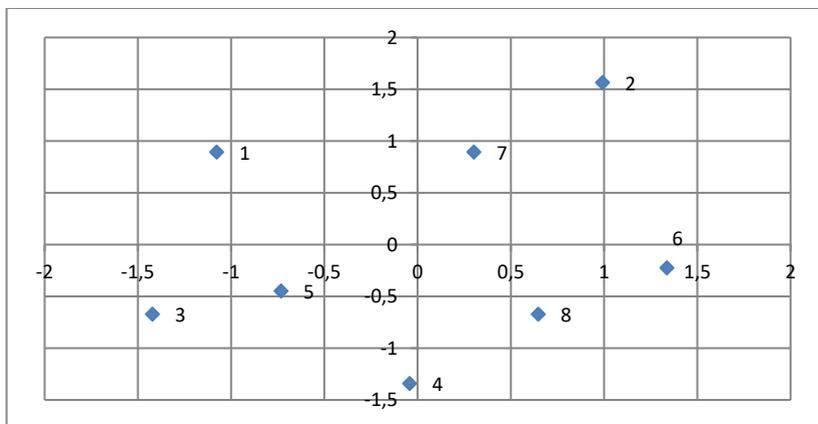


Рис. А2. Изображение образов объектов на двухмерной координатной плоскости преобразованных данных

На двухмерной координатной плоскости преобразованных данных взаимное расположение образов объектов осталось прежним. Это означает, что связи между объектами сохранились.

3. Построим матрицу расстояний с использованием заданной меры расстояний (Евклидово расстояние). Матрица представляет собой таблицу, в которой наименование строк и столбцов – номер экземпляра, а пересечение строк и столбцов – расстояние между соответствующими экземплярами в соответствии с выбранной мерой расстояния. Очевидно, что на пересечении строк и столбцов, соответствующих одинаковым экземплярам, расстояние будет равно 0. Матрица расстояний представлена табл. А5.

Таблица А5. Табличное представление матрицы расстояний

№ Эжз.	1	2	3	4	5	6	7	8
1	0	2,1749	1,6028	2,4637	1,3852	2,6601	1,3793	2,3286
2	2,1749	0	3,2903	3,0855	2,6500	1,8218	0,9621	2,2625
3	1,6028	3,2903	0	1,5337	0,7250	2,7945	2,3286	2,0689
4	2,4637	3,0855	1,5337	0	1,1294	1,7755	2,2625	0,9621
5	1,3852	2,6500	0,7250	1,1294	0	2,0809	1,6941	1,3973
6	2,6601	1,8218	2,7945	1,7755	2,0809	0	1,5232	0,8219
7	1,3793	0,9621	2,3286	2,2625	1,6941	1,5232	0	1,6028
8	2,3286	2,2625	2,0689	0,9621	1,3973	0,8219	1,6028	0

Кроме того, в таблице относительно диагональных ячеек матрица отзеркалена. Поэтому, в дальнейшем для удобства матрица будет представлена в виде полуматрицы.

4. В соответствии с заданным правилом объединения на каждом шаге объединение образов экземпляров или сформированных на прошлых шагах кластеров будет происходить по минимальному расстоянию. Таким образом, согласно табл. А5 минимальное расстояние 0,725 наблюдается между образами 3 и 5 экземпляров. Эти экземпляры объединяются в кластер. Из этого следует, что матрица расстояний преобразовывается в вид, представленный в табл. А6.

5. Т.к. в соответствии с правилом объединения метод одиночной связи, заключающийся в том, что объединению подвергаются кластеры, между которыми наблюдается минимальное расстояние, а расстоянием между кластерами является минимальное расстояние между любыми двумя образами, принадлежащим разным кластерам, можно сделать вывод, что в преобразованную матрицу расстояний необходимо записать то значение, которое является минимальным между рассматриваемым объектом (образом, кластером) и образом 3 или 5 объекта.

6. Согласно табл. А6 минимальное расстояние 0,8219 наблюдается между образами 6 и 8 экземпляров. Эти экземпляры объединяются в кластер. Из этого следует, что матрица расстояний преобразовывается в вид, представленный в табл. А7.

7. Согласно табл. А7 минимальное расстояние 0,9621 наблюдается сразу у 2 вариантов: между образами 2 и 7 экземпляров и между образом 4 экземпляра и кластером, состоящим из образов 6, 8. Для удобства сначала рассмотрим первый случай, преобразуем матрицу расстояний (табл. А8), а затем рассмотрим второй случай и преобразуем матрицу расстояний, представленную в табл. А8 в матрицу, представленную в табл. А9.

Таблица А6. Табличное представление матрицы расстояний  
после объединения 3 и 5 экземпляра

№ Экз.	1	2	3,5	4	5	6	7	8
1	0	2,1749	1,3852	2,4637	-	2,6601	1,3793	2,3286
2	-	0	2,6500	3,0855	-	1,8218	0,9621	2,2625
3,5	-	-	0	1,1294	-	2,0809	1,6941	1,3973
4	-	-	-	0	-	1,7755	2,2625	0,9621
-	-	-	-	-	0	-	-	-
6	-	-	-	-	-	0	1,5232	0,8219
7	-	-	-	-	-	-	0	1,6028
8	-	-	-	-	-	-	-	0

Таблица А7. Табличное представление матрицы расстояний после объединения 6 и 8 экземпляра

№ Экз.	1	2	3,5	4	5	6,8	7	8
1	0	2,1749	1,3852	2,4637	-	2,3286	1,3793	-
2	-	0	2,6500	3,0855	-	1,8218	0,9621	-
3,5	-	-	0	1,1294	-	1,3973	1,6941	-
4	-	-	-	0	-	0,9621	2,2625	-
-	-	-	-	-	0	-	-	-
6,8	-	-	-	-	-	0	1,5232	-
7	-	-	-	-	-	-	0	-
8	-	-	-	-	-	-	-	0

8. Согласно табл. А9 минимальное расстояние 1,1294 наблюдается между двумя кластерами: кластером, образованным образцами 3 и 5 экземпляров и кластером, образованным образцами 4, 6 и 8 экземпляров. Эти кластеры объединяются в один кластер. Из этого следует, что матрица расстояний преобразовывается в вид, представленный в табл. А10.

9. Согласно табл. А10 минимальное расстояние 1,3793 наблюдается между образцом 1 экземпляра и кластером, состоящим из образцов 2, 7. Эти образцы объединяются в один кластер. Из этого следует, что матрица расстояний преобразовывается в вид, представленный в табл. А11.

10. Из табл. А11 видно, что все образцы экземпляров объединяются в единый кластер на расстоянии 1,3852.

11. Так как все образцы экземпляров объединились в один кластер, а также известны расстояния, на которых образцы объектов объединялись (табл. А12), то можно построить дендрограмму, представленную на рис.А3.

**Таблица А8. Табличное представление матрицы расстояний после объединения 2 и 7 экземпляра**

<b>№ Экз.</b>	<b>1</b>	<b>2,7</b>	<b>3,5</b>	<b>4</b>	<b>6,8</b>
<b>1</b>	0	1,3793	1,3852	2,4637	2,3286
<b>2,7</b>	-	0	2,6500	3,0855	1,5232
<b>3,5</b>	-	-	0	1,1294	1,3973
<b>4</b>	-	-	-	0	0,9621
<b>6,8</b>	-	-	-	-	0

**Таблица А9. Табличное представление матрицы расстояний после объединения образа 4 экземпляра и кластером образов 6, 8 экземпляров**

<b>№ Экз.</b>	<b>1</b>	<b>2,7</b>	<b>3,5</b>	<b>4,6,8</b>
<b>1</b>	0	1,3793	1,3852	2,3286
<b>2,7</b>	-	0	2,6500	1,5232
<b>3,5</b>	-	-	0	1,1294
<b>4,6,8</b>	-	-	-	0

**Таблица А10. Табличное представление матрицы расстояний после объединения двух кластеров (образов 3 и 5 экземпляров и образов 4, 6, 8 экземпляров)**

<b>№ Экз.</b>	<b>1</b>	<b>2,7</b>	<b>3,4,5,6,8</b>
<b>1</b>	0	1,3793	1,3852
<b>2,7</b>	-	0	1,5232
<b>3,4,5,6,8</b>	-	-	0

**Таблица А11. Табличное представление матрицы расстояний после объединения образа 1 экземпляра и кластером образов 2, 7 экземпляров**

<b>№ Экз.</b>	<b>1,2,7</b>	<b>3,4,5,6,8</b>
<b>1,2,7</b>	0	1,3852
<b>3,4,5,6,8</b>	-	0

**Таблица А12. Соответствие объединенных объектов и расстояния объединения**

<b>Объединение объектов</b>	<b>Расстояние</b>
3 и 5	0,725
6 и 8	0,8219
2 и 7	0,9621
4 и (6, 8)	0,9621
(3, 5) и (4, 6, 8)	1,1294
1 и (2, 7)	1,3793
(1...8)	1,3852

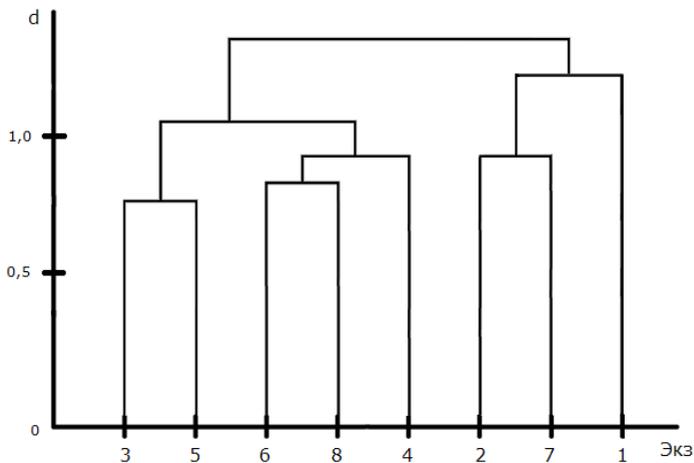


Рис. А3. Вертикальная дендрограмма

**Вывод:** дендрограмма, изображенная на рис. А3, показывает состав кластеров при различных расстояниях объединения, причем эти кластеры в наибольшей степени различны. Два кластера образуются в диапазоне 1,3793...1,3852 расстояний. При этом первый кластер состоит из образов 3, 4, 5, 6, 8 экземпляров. Второй кластер состоит из образов 1, 2, 7 экземпляров.

## ПРИЛОЖЕНИЕ Б

### ПРИМЕР РЕШЕНИЯ ЗАДАЧИ КЛАССИФИКАЦИИ МЕТОДОМ k-СРЕДНИХ КЛАСТЕРНОГО АНАЛИЗА

**Задача:** дана партия диодов, для описания свойств которой осуществили выборку в объеме 8 экземпляров. Информативные параметры экземпляров выборки представлены в табл. Б1. Необходимо определить состав двух кластеров так, чтобы эти кластеры были в наибольшей степени различны.

Для кластеризации использовать метод кластеризации k-средних, за меру расстояния принять Евклидово расстояние, за начальные центры кластеров – образы экземпляров 1 и 2.

Таблица Б1. Значения информативных параметров  
экземпляров выборки

№ экз.	1	2	3	4	5	6	7	8
X1	9	10	2	6	12	4	12	13
X2	6	5	7	4	10	5	3	9

Решение задачи включает в себя выполнение следующих этапов.

1. Изобразим расположение образов экземпляров на двумерной координатной плоскости (рис. Б1). По вертикальной оси отложим значения параметра X1, по горизонтальной – значения параметра X2. Число возле изображения образа объекта – номер экземпляра.

2. Для подготовки данных и приведения значения параметров к единому масштабу измерений проведем стандартизацию выборки по следующим формулам:

$$X_{\text{cp}} = \frac{\sum_{i=1}^n X_i}{n};$$

$$Z_i = \frac{X - X_{\text{cp}}}{\sigma};$$

$$\sigma = \sqrt{D} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - X_{\text{cp}})^2},$$

где  $X_{\text{cp}}$  – среднее значение параметра  $X$ ;

$X_i$  – значение параметра  $i$ -го экземпляра;

$n$  – объем выборки;

$\sigma$  – среднее квадратичное отклонение;

$Z_i$  – преобразованное значение параметра  $X$   $i$ -го экземпляра;

$D$  – дисперсия.

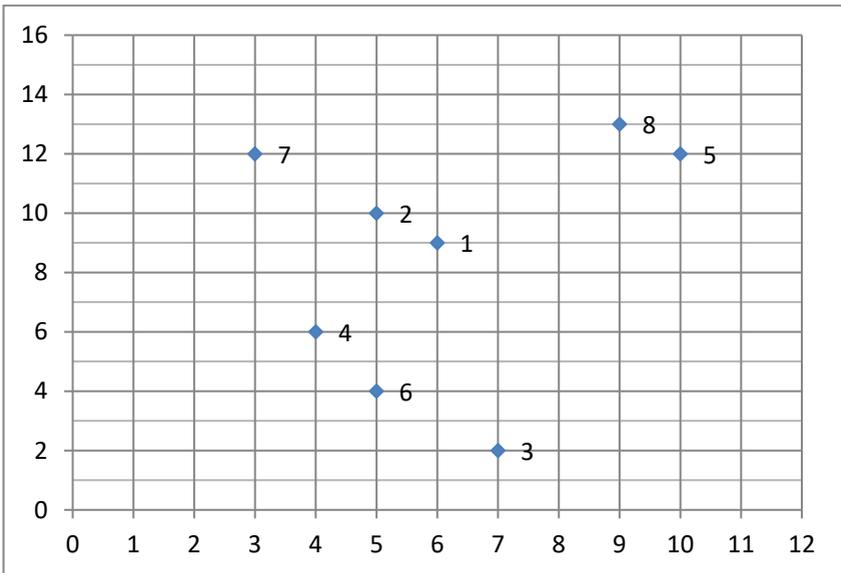


Рис. Б1. Изображение образов объектов на двухмерной координатной плоскости

Таким образом, получим следующие таблицы для параметров  $X_1$  и  $X_2$  (табл. Б2, табл. Б3). Округление значений произведем до 4 знака после запятой.

По результатам вычислений сформируем сводную таблицу преобразованных данных (табл. Б4) и изобразим образы объектов на двумерной координатной плоскости, по вертикальной оси которой отложим значения преобразованного параметра  $Z_1$ , по горизонтальной – значения параметра  $Z_2$  (рис. Б2). Число возле изображения образа объекта – номер экземпляра.

Таблица Б2. Таблица преобразования данных параметра  $X_1$

№ экз.	$X_1$	$X_1 - X_{cp}$	$(X_1 - X_{cp})^2$	$Z_1$
1	9	0,5	0,25	0,1228
2	10	1,5	2,25	0,3685
3	2	-6,5	42,25	-1,5967
4	6	-2,5	6,25	-0,6141
5	12	3,5	12,25	0,8598
6	4	-4,5	20,25	-1,1054
7	12	3,5	12,25	0,8598
8	13	-4,5	20,25	1,1054
<hr/>				
$X_{cp}$	8,5	$\sigma$	4,0708	

Таблица Б3. Таблица преобразования данных параметра  $X_2$

№ экз.	$X_2$	$X_2 - X_{cp}$	$(X_2 - X_{cp})^2$	$Z_2$
1	6	-0,125	0,0156	-0,0517
2	5	-1,125	1,2656	-0,4656
3	7	0,875	0,7656	0,3621
4	4	-2,125	4,5156	-0,8794
5	10	3,875	15,0156	1,6036
6	5	-1,125	1,2656	-0,4656
7	3	-3,125	9,7656	-1,2932

Окончание табл. Б3

<b>8</b>	9	2,875	8,2656	1,1898
$X_{cp}$	6,125	$\sigma$	2,4165	

Таблица Б4. Сводная таблица значений преобразованных данных

№ экз.	1	2	3	4	5	6	7	8
<b>Z1</b>	0,1228	0,3685	-1,5967	-0,6141	0,8598	-1,1054	0,8598	1,1054
<b>Z2</b>	-0,0517	-0,4656	0,3621	-0,8794	1,6036	-0,4656	-1,2932	1,1898

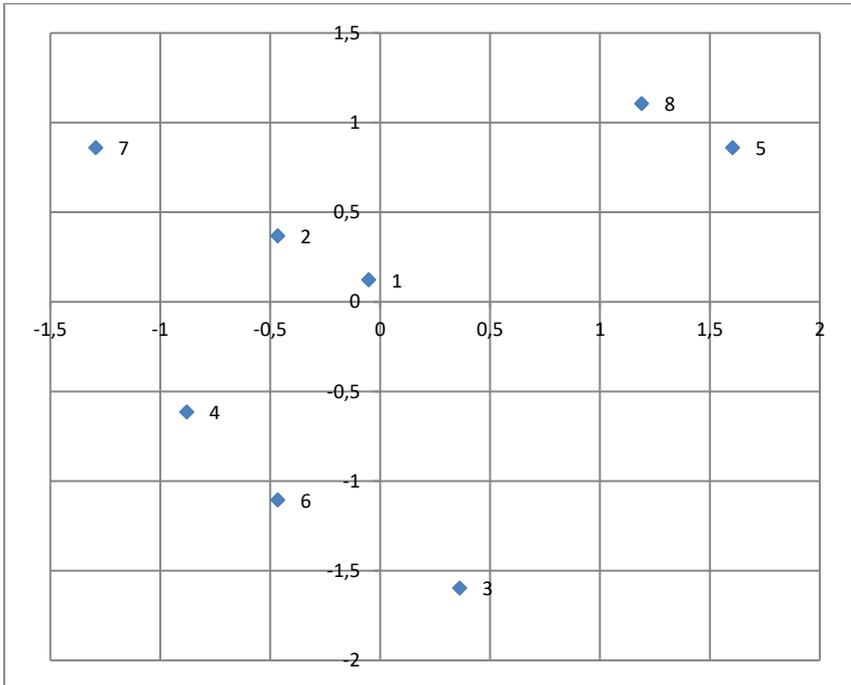


Рис. Б2. Изображение образов объектов на двухмерной координатной плоскости преобразованных данных

На двумерной координатной плоскости преобразованных данных взаимное расположение образов объектов осталось прежним. Это означает, что связи между объектами сохранились.

3. Построим матрицу расстояний с использованием заданной меры расстояний (Евклидово расстояние). Матрица представляет собой таблицу, в которой наименование строк и столбцов – номер экземпляра, а пересечение строк и столбцов – расстояние между соответствующими экземплярами в соответствии с выбранной мерой расстояния. Очевидно, что на пересечении строк и столбцов, соответствующих одинаковым экземплярам, расстояние будет равно 0. Матрица расстояний представлена табл. Б5.

4. Согласно заданию необходимо найти состав двух кластеров, за начальные центры кластеров приняты образы экземпляров 1 и 2. Распределим образы экземпляров 3, 4, 5, 6, 7, 8 по кластерам по критерию минимального расстояния до центров начальных центров, используя матрицу расстояний. Таким образом, первый кластер состоит из 1, 3, 5, 6, 8 экземпляров, второй кластер – из 2, 4, 7.

Теперь можно пересчитать центры вновь сформированных кластеров как покоординатные средние значения:

$$Z_1^{(1)} = \frac{0,123 - 1,597 + 0,860 - 1,105 + 1,105}{5} = -0,123;$$

$$Z_2^{(1)} = \frac{-0,052 + 0,362 + 1,604 - 0,466 + 1,190}{5} = 0,528;$$

$$Z_1^{(2)} = \frac{0,368 - 0,614 + 0,860}{3} = 0,205;$$

$$Z_2^{(2)} = \frac{-0,466 - 0,879 - 1,293}{3} = -0,879.$$

Таблица Б5. Табличное представление матрицы расстояний

№ ЭКз.	1	2	3	4	5	6	7	8
1	0	0,481	1,769	1,108	1,812	1,296	1,444	1,583
2	0,481	0	2,132	1,066	2,127	1,474	0,962	1,812
3	1,769	2,132	0	1,583	2,752	0,962	2,962	2,826
4	1,108	1,066	1,583	0	2,887	0,642	1,531	2,96
5	1,812	2,127	2,752	2,887	0	2,854	2,897	0,481
6	1,296	1,474	0,962	0,642	2,854	0	1,814	2,493
7	1,444	0,962	2,962	1,531	2,897	1,814	0	2,495
8	1,583	1,812	2,826	2,96	0,481	2,493	2,495	0

Таблица Б6. Координаты центров кластеров

Координаты центров	Кластер № 1	Кластер № 2
Z1	-0,123	0,205
Z2	0,528	-0,879

5. Пересчитаем расстояния от каждого образа до центров кластеров (табл. Б7) согласно формуле Евклидова расстояния:

$$d_{AB} = \sqrt{(Z_1^{(A)} - Z_1^{(B)})^2 + (Z_2^{(A)} - Z_2^{(B)})^2}.$$

Таблица Б7. Таблица расстояний от образа каждого экземпляра до центра каждого кластера

№ Экз	1	2	3	4	5	6	7	8
№ Класт								
1	0,630	1,109	1,483	1,490	1,457	1,397	2,069	1,395
2	0,831	0,44	2,188	0,819	2,568	1,374	0,775	2,256

6. Проанализировав табл. Б7, можно сделать вывод, что составы кластеров изменились – образ 6 экземпляра до кластера № 2 находится на меньшем расстоянии, чем до кластера № 1. Тогда составы кластеров меняются. Теперь первый кластер состоит из 1, 3, 5, 8 экземпляров, второй кластер – из 2, 4, 6, 7.

Теперь можно пересчитать центры вновь сформированных кластеров как покоординатные средние значения:

$$Z_1^{(1)} = \frac{0,123 - 1,597 + 0,860 + 1,105}{4} = 0,123;$$

$$Z_2^{(1)} = \frac{-0,052 + 0,362 + 1,604 + 1,190}{4} = 0,776;$$

$$Z_1^{(2)} = \frac{0,368 - 0,614 - 1,105 + 0,860}{4} = -0,123;$$

$$Z_2^{(2)} = \frac{-0,466 - 0,879 - 0,466 - 1,293}{4} = -0,776.$$

Таблица Б8. Координаты центров кластеров

Координаты центров	Кластер № 1	Кластер № 2
Z1	0,123	-0,123
Z2	0,776	-0,776

7. Пересчитаем расстояния от каждого образа до центров кластеров (табл. Б9).

Таблица Б9. Таблица расстояний от образа каждого экземпляра до центра каждого кластера

№ Экз № Класт	1	2	3	4	5	6	7	8
	1	0,828	1,266	1,769	1,812	1,108	1,747	2,196
2	0,765	0,581	1,862	0,502	2,575	1,030	1,111	2,318

8. Проанализировав табл. Б9, можно сделать вывод, что составы кластеров изменились – образ 1 экземпляра до кластера № 2 находится на меньшем расстоянии, чем до кластера № 1. Тогда составы кластеров меняются. Теперь первый кластер состоит из 3, 5, 8 экземпляров, второй кластер – из 1, 2, 4, 6, 7.

Теперь можно пересчитать центры вновь сформированных кластеров как покоординатные средние значения:

$$Z_1^{(1)} = \frac{-1,597 + 0,860 + 1,105}{3} = 0,123;$$

$$Z_2^{(1)} = \frac{0,362 + 1,604 + 1,190}{3} = 1,052;$$

$$Z_1^{(2)} = \frac{0,123 + 0,368 - 0,614 - 1,105 + 0,860}{5} = -0,074;$$

$$Z_2^{(2)} = \frac{-0,052 - 0,466 - 0,879 - 0,466 - 1,293}{5} = -0,631.$$

Таблица Б10. Координаты центров кластеров

Координаты центров	Кластер № 1	Кластер № 2
Z1	0,123	-0,074
Z2	1,052	-0,631

9. Пересчитаем расстояния от каждого образа до центров кластеров (табл. Б11).

Таблица Б11. Таблица расстояний от образа каждого экземпляра до центра каждого кластера

№ Экз	1	2	3	4	5	6	7	8
№ Класт								
1	1,104	1,538	1,853	2,067	0,921	1,953	2,458	0,992
2	0,612	0,472	1,818	0,594	2,422	1,044	1,145	2,170

10. Проанализировав табл. Б11, можно сделать вывод, что составы кластеров изменились – образ 3 экземпляра до кластера № 2 находится на меньшем расстоянии, чем до кластера № 1. Тогда составы кластеров меняются. Теперь первый кластер состоит из 5, 8 экземпляров, второй кластер – из 1, 2, 3, 4, 6, 7.

Теперь можно пересчитать центры вновь сформированных кластеров как по координатным средним значениям:

$$Z_1^{(1)} = \frac{0,860 + 1,105}{2} = 0,983;$$

$$Z_2^{(1)} = \frac{1,604 + 1,190}{2} = 1,397;$$

$$Z_1^{(2)} = \frac{0,123 + 0,368 - 1,597 - 0,614 - 1,105 + 0,860}{6} = -0,328;$$

$$Z_2^{(2)} = \frac{-0,052 - 0,466 + 0,362 - 0,879 - 0,466 - 1,293}{6} = -0,466.$$

Таблица Б12. Координаты центров кластеров

Координаты центров	Кластер № 1	Кластер № 2
Z1	0,983	-0,328
Z2	1,397	-0,466

11. Пересчитаем расстояния от каждого образа до центров кластеров (табл. Б13).

Таблица Б13. Таблица расстояний от образа каждого экземпляра до центра каждого кластера

№ Экз	1	2	3	4	5	6	7	8
№ Класт								
1	1,685	1,962	2,780	2,780	0,241	2,798	2,693	0,240
2	0,612	0,696	1,515	0,502	2,387	0,777	1,448	2,190

12. По табл. Б13 видно, что у всех образов экземпляров наименьшие расстояния наблюдаются с теми кластерами, которым эти образы принадлежат, пересчет центров кластеров из-за изменения состава кластеров не требуется. Следовательно, кластеры стабилизировались.

**Вывод:** при использовании кластеризации методом k-средних были найдены два в наибольшей степени различных кластера. Первый кластер состоит из образов 5, 8 экземпляров, второй кластер состоит из образов 1, 2, 3, 4, 6, 7 экземпляров.

Учебное издание

*Мишанов Роман Олегович*

**ИСПОЛЬЗОВАНИЕ МЕТОДОВ  
СТАТИСТИЧЕСКОЙ КЛАССИФИКАЦИИ  
ДЛЯ ОПРЕДЕЛЕНИЯ КАЧЕСТВА  
КОМПОНЕНТОВ БОРТОВЫХ РЭС**

*Учебное пособие*

Редактор Л. Р. Дмитриенко  
Компьютерная верстка Л. Р. Дмитриенко

Подписано в печать 09.12.2021. Формат 60x84 1/16.  
Бумага офсетная. Печ. л. 5,0.  
Тираж 25 экз. Заказ . Арт. – (РЗУ)/2021.

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«САМАРСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ АКАДЕМИКА С.П. КОРОЛЕВА»  
(САМАРСКИЙ УНИВЕРСИТЕТ)  
443086, Самара, Московское шоссе, 34.

---

Издательство Самарского университета.  
443086, Самара, Московское шоссе.

