

《智能信息处理》课程作业

## 基于形式概念分析的图像检索

董雪松

作业	分数[20]
得分	

2021 年 11 月 29 日

# 基于形式概念分析的图像检索

董雪松

(大连海事大学 信息科学技术学院, 大连 116026)

**摘 要** 形式概念分析是由 Wille 教授于 1982 年首先提出的, 它提供了一种支持数据分析的有效工具。形式概念包括外延和内涵两部分, 它本质上描述了形式对象与形式属性之间的关系。概念格通过 Hasse 图生动和简洁地体现了这些概念之间的泛化和特化关系。从形式背景中生成概念格的过程实质上是一种概念聚类过程。概念格可以用于完成许多机器学习的任务, 并且概念格在信息检索、数字图书馆、软件工程和知识发现等方面也有很广泛的应用。

**关键词** 形式概念分析; 形式背景; 概念格; 形式概念; 图像检索

## Image retrieval based on formal concept analysis

DongXuesong

( School of Information Science and Technology, Dalian Maritime University, Dalian 116026)

**Abstract** Formal concept analysis was first proposed by Professor wille in 1982. It provides an effective tool to support data analysis. The concept of form includes denotation and connotation, which essentially describes the relationship between formal objects and formal attributes. Concept lattice vividly and concisely reflects the generalization and specialization relationship between these concepts through Hasse diagram. The process of generating concept lattice from formal context is essentially a concept clustering process. Concept lattice can be used to complete many machine learning tasks, and it is also widely used in information retrieval, digital library, software engineering and knowledge discovery.

**Key words** Formal concept analysis; Formal background; Concept lattice; search image

# 1 基本概念与定义

## 1.1 引言

随着科学技术的进步,特别是多媒体数字化技术和网络技术的发展和推广,现代信息检索所处理的对象和规模都有了很大的变化,并对人们的生活和社会的发展产生了重要的影响。<sup>[1]</sup>飞速膨胀的数字图像逐渐成为信息的主流,大量的组织和机构都有非常庞大的图像收藏,这些收藏往往都是以数字化的形式被存储的,并且可以通过计算机网络存取。数码摄影和数码摄像等数字化技术的出现也使得越来越多的普通人能够拥有个人的数字多媒体集合。这些大小各异、形式多样的数字化视听集不仅会长期存在,<sup>[2]</sup>而且会越来越多地产生于各种场合。因特网中的令人感兴趣的图像的数量正以极高的速度增长。然而,由于这些数字图像是无序地分布在世界各地,其中所包含的信息无法被有效地访问和利用。因此,人们迫切需要一种能够快速而且准确地查找访问图像的技术,这就是图像检索技术<sup>[3]</sup>。

早期图像数据库的检索方法主要是基于文本方式,通过人工对每一幅图像建立关键词等描述信息作为图像索引,采用传统数据库方式来满足图像检索的要求。文本标注可以简单清楚地描述图像高层抽象的语义,但是,随着因特网的出现和网络技术的发展,这种检索方法的不足和局限性日益凸出,不同人对同一幅图像有不同理解,这样文本描述信息就存在着多义性,完全人工标注工作量太大,并且关键词无法完全概括图像内容。

为此,人们提出了基于内容(Content Based Image Retrieval,CBIR)的图像检索,从一定程度上解决基于文本图像检索的局限性与人们图像检索需求之间的矛盾。CBIR 使用可以直接从图像中获得的客观的视觉内容特征,如颜色、纹理、形状等来判断图像之间的相似性<sup>[4,5]</sup>。

在基于内容的图像检索中,以图像的底层视觉和形象特征为索引对图像进行检

索具有计算简单、性能稳定的特点,但目前这些特征都有一定的局限性<sup>[5]</sup>。然而,人们判断图像的相似性并非仅仅建立在图像视觉特征的相似性上,语义图像检索也就应运而生,成为解决图像简单视觉特征和用户检索丰富语义之间存在的语义鸿沟的关键。为了实现更为贴近用户理解能力的自然而简洁的查询方式,并提高图像检索的精度,在 CBIR 技术领域进行包含语义的检索方法的研究是十分重要的。

图像检索的过程可以理解为从图像数据中形成概念的过程,反映的是从图像数据中抽取出特征,形成概念,来研究分析概念之间的关系,从而检索出最相似图像。形式概念分析理论是由德国的数学家 Wille 于 1982 年提出来的<sup>[6]</sup>。形式概念分析中的“形式”一词表示我们正在处理领域的工作,通过与这些工作相联系的结构化的概念的联系,发现可理解的,有意义的知识;“概念”是对哲学中概念的一种数学表示,是对人们认知的知识的一种数学化描述。形式概念分析理论的反映概念形成的过程与数据挖掘中从数据中产生知识的过程相似,因此,本文试图将形式概念分析理论引入图像检索,并探讨其可行性。

## 1.2 形式概念的定义

形式概念(Formal Concept Analysis, FCA)由 Wille 博士于 1982 年提出,是应用数学的一个分支、是信息处理的一种理论、是知识处理的一种理论。形式概念是用来构建自然概念的层次联通结构的。形式概念的定义如下所示:

设形式对象集  $U \quad X \in U$

形式属性集  $A \quad B \in A$

二元关系  $R \in U \times A$

若  $X = \{x | x \in U, \forall a \in B, xRa\}$

$B = \{a | a \in A, \forall x \in X, xRa\}$

则二元组  $(X, B)$  被称为形式概念

$X$  中  $x$  每个  $x$  都有全部属性

$B$  中  $a$  每个  $x$  都有的属性

我们可以把形式概念理解成为数学上的概念，因为形式概念等于对象集的属性集，其中对象集和属性集都是在数学上成立的。与概念的表示方法类似，形式概念也有三种表示方法，分别是表达式法，二维表法和图示法。形式概念的作用就是构建自然概念的层次连通结构，为了更好的解释此作用，根据上述描述的内容，在自然概念的基础之上可以建立一个新的形式概念内容，从而对现实世界有一个好的理解。下面继续以人狗的例子来说明形式概念。

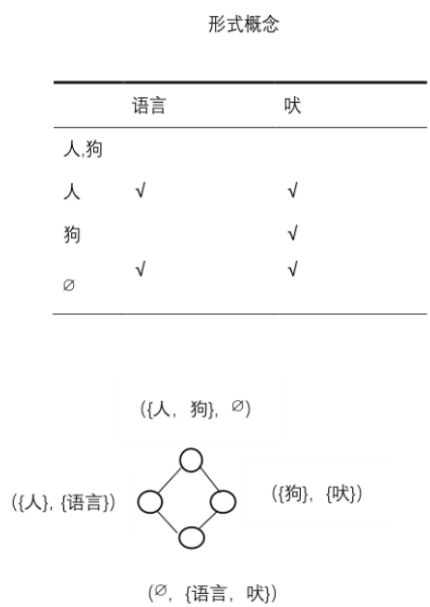


图 1 形式概念的表示

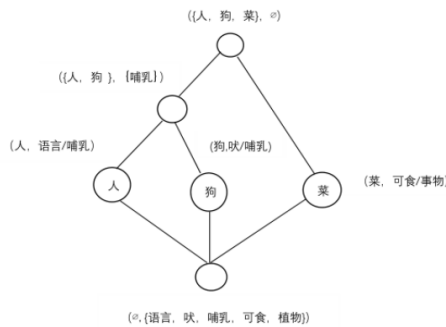


图 2 形式概念的实例

如图 2 所示，为形式概念分析的实例，运用形式概念可以极大限度地对集合中具有某种关系或者含有某些共同属性的元素进行分类，发现由属性和对象构成的概念和概念之间的关系。例如人和狗都具有哺乳这一属性，则将人和狗、哺乳分别

作为对象和属性提取出来，由此类推，则可以发现不同对象之间的联系和不同对象属性之间的联系。事实证明，应用形式概念，在信息检索，智能信息处理等方面起着巨大的作用。

1.3 形式背景和概念格

定义 1：形式背景是一个知识学科结构的三元组  $KS=(A, K, R)$ ，其中  $A$  是作者(对象)的集合， $K$  是主题关键词(属性)的集合， $R$  是  $A$  和  $K$  之间的一个二元关系，即  $R \subseteq A \times K$ 。 $aRk$  表示  $a \in A$  与  $k \in K$  之间存在关系  $R$ ，读作作者(对象) $a$  具有关键词(属性) $k$ 。表 1 是一个由 8 个作者和他们在论文中使用过的 9 个关键词所构成的形式背景。其中，“ $\times$ ”表示作者  $a_i$  标注了关键词  $k_j$ ，空格表示作者  $a_i$  未使用关键词  $k_j$ (下同)。

表 1 形式背景举例

	k <sub>1</sub>	k <sub>2</sub>	k <sub>3</sub>	k <sub>4</sub>	k <sub>5</sub>	k <sub>6</sub>	k <sub>7</sub>	k <sub>8</sub>	k <sub>9</sub>
a <sub>1</sub>	×	×					×		
a <sub>2</sub>	×	×					×	×	
a <sub>3</sub>	×	×	×				×	×	
a <sub>4</sub>	×		×				×	×	×
a <sub>5</sub>	×	×		×		×			
a <sub>6</sub>	×	×	×	×		×			
a <sub>7</sub>	×		×	×	×				
a <sub>8</sub>	×		×	×		×			

定义 2：设  $P$  是作者集合  $A$  的一个子集，定义  $f(P)=\{k \in K | \forall a \in P, aRk\}$ ，表示  $P$  中作者共同关键词的集合；相应地，设  $T$  是主题关键词集合  $K$  的一个子集，定义  $g(T)=\{a \in P | \forall k \in T, aRk\}$ ，表示具有  $T$  中所有关键词的作者集合。以表 1 为例，设  $p1=\{a1, a2\}$ ，则  $f(p1)=\{k1, k2, k7\}$ ；设  $T1=\{k1, k2, k7\}$ ，则  $g(T1)=\{a1, a2, a3\}$ 。

定义 3：形式背景  $(A, K, R)$  上的一个形式概念 (formal concept) 是二元组  $(P, T)$ ，其中  $P \subseteq A, T \subseteq K$ ，且满足  $f(P)=T$  和  $g(T)=P$ 。我们称  $P$  是形式概念  $(P, T)$  的外延， $T$  是形式概念  $(P, T)$  的内涵。在定义 2 的例子中， $f(p1)=T1$ ，但  $g(T1) \neq p1$ ，所以  $(p1, T1)$  不是形式概念。但若设  $p2=\{a1, a2, a3\}, T2=\{k1, k2, k7\}$ ，

则  $f(p_1)=T_2$  ,  $g(T_2)=p_2$  , 所以  $(p_2, T_2)$  是形式概念。

定义 4: 若  $(p_1, T_1)$ 、 $(p_2, T_2)$  是某个形式背景  $(A, K, R)$  上的两个概念, 如果  $p_1 \subseteq p_2$  (或  $T_2 \subseteq T_1$ ), 那么  $(p_1, T_1)$  被称为  $(p_2, T_2)$  的子概念,  $(p_2, T_2)$  被称为  $(p_1, T_1)$  的超概念, 并将其记作  $(p_1, T_1) \leq (p_2, T_2)$ 。关系  $\leq$  被称为形式概念之间的偏序关系。超概念与子概念的关系是所有形式概念集合上的偏序关系。例如在表 1 中, 以概念  $C_1 = (\{a_2, a_3\}, \{k_1, k_2, k_7, k_8\})$  和概念  $C_2 = (\{a_1, a_2, a_3\}, \{k_1, k_2, k_7\})$  为例, 因为概念  $C_1$  的外延  $\{a_2, a_3\} \subseteq$  概念  $C_2$  的外延  $\{a_1, a_2, a_3\}$ , 而同时概念  $C_2$  的内涵  $\{k_1, k_2, k_7\} \subseteq$  概念  $C_1$  的内涵  $\{k_1, k_2, k_7, k_8\}$ , 所以概念  $C_1$  是概念  $C_2$  的子概念, 概念  $C_2$  是概念  $C_1$  的超概念。

## 2 基于形式概念分析的图像

### 检索算法

当用户输入一幅待查询图像时, 图像检索过程便开始进行。待检索的图像经过处理, 生成形式背景, 并建立成一个概念格。图像库中的图像只有当与待检索图像进行比较时, 才建立概念格。因此, 在检索的任一时刻, 计算机内存中只有一幅待检索的图像和一幅图像库中的图像。

该方法是自动提取待检索图像的合适特征, 所以在检索过程中不需要用户参与, 因此不用提供与用户交互功能。

基于形式概念分析的图像检索算法

(1) 根据文献<sup>[9]</sup>中算法求各类图像  $C_i$  的特征向量  $\{X_j\}_{j=1}^N$  及待查询图像  $D$  的特征向量  $Y$ 。

(2) 计算待检索图像  $D$  对各类图像  $C_i$  的图像隶属度  $\mu_{C_i}(Y)$ 。

(3) 若  $\mu_{C_i}(Y) \geq \frac{\max\{\mu_{C_j}(Y)\}}{2}$ , 其中,  $1 \leq i \leq m$ , 则  $C_i$  为待检索图像的候选类别。

(4) 建立待查询图像  $D$  的概念格, 生成概念集合  $S = \{S_1, S_2, \dots, S_m\}$ , 其中  $S_i$  是查询概念格中的一个概念,  $m$  是待查询概念格的概念总数。

(5) 建立候选类别  $C_1$  中的一副图像  $C_{11}$  的概念格, 生成概念集合  $T = \{T_1, T_2, \dots, T_n\}$ , 其中  $T_i$  是概念格中的一个概念,  $n$  是概念格的概念总数。

(6) 将待查询的概念集合中  $S_1$  的与图像中的每一个节点  $T_i$  进行概念匹配。

(7) 迭代(6), 直到  $S$  中所有  $S_i$  均匹配完。

(8) 根据公式计算待查询图像与此幅图像的相似度。

(9) 迭代(5)-(8), 直至候选类别  $C_1$  中所有图像均被匹配完。

(10) 迭代(5)-(9), 直至候选类别  $C_i$  中所有类别的图像均被匹配完。

(11) 按相似度从大到小的顺序排列输出。

## 3 总结

由 Wille R 于 1982 年首先提出的形式概念分析提供了一种支持数据分析的有效工具。形式概念包括外延和内涵两部分, 它本质上描述了形式对象与形式属性之间的关系。形式背景则以二维表的形式描述了不同形式对象以及他们的形式属性之间的关系。在此基础上的概念格则是以 Hasse 图的形式, 描述了不同形式对象之间的父子关系。概念格结构模型来源于形式概念分析理论, 是形式概念分析理论中的核心数据分析工具。

概念格体现了一种概念层次结构, 实现了对数据的可视化。因此, 概念格受到了人们的广泛关注。概念格理论经过几十年的发展, 如今已被广泛运用于软件工程、知识工程、人工智能等领域。本文从概念一词的含义到形式概念的定义及表示, 再到形式背景, 再到格的解释以及概念格的概念以及应用等方面总结了概念格的研究进展。当然, 概念格仍是一个高速发展的领域, 对于粗糙模糊概念格的研究、基于

概念格的数据挖掘模型的实现等等都是以后的重点研究方向。

### 参 考 文 献

- [1] 胡可云, 陆玉昌, 石纯一. 概念格及其应用进展[J]. 清华大学学报(自然科学版), 2000, 40(9):77-81.
- [2] 张云中. 基于形式概念分析的领域本体构建方法研究[D]. 吉林大学, 吉林, 2009.
- [3] Rui Y , Huang T S, Chang S F. Image retrieval: current techniques, promising directions, and open issues. Journal of Visual Communication and Image Representation, 1999,10(4): 39-62
- [4] 庄越挺,潘云鹤.基于内容的图像检索综述.模式识别与人工智能,1999.6
- [5] 朵琳, 杨丙. 一种基于用户兴趣概念格的推荐评分预测方法[J]. 小型微型计算机系统, 2020(10).
- [6]Ganter B, Wille R.. Formal Concept Analysis. Springer-Verlag, Berlin, Heidelberg, New York,1996.
- [7] Gorkani M, Picard R W. Texture Orientation for Sorting Photos at a Glance. In:Proc of Int. Conf. Part Rec.,1994,459-464