

学校编码: 10384  
学号: 15420141152004

分类号 \_\_\_\_\_ 密级 \_\_\_\_\_  
UDC \_\_\_\_\_

厦门大学

硕士 学位 论文

# 复杂关联数据的张量模型与应用研究

Tensor Model of Complex Correlated Data and Application  
Research

王丹

指导教师姓名: 钱争鸣 教授

专业名称: 统计学

论文提交日期: 2017 年 4 月

论文答辩时间: 2017 年 4 月

学位授予日期: 2017 年 6 月

答辩委员会主席: \_\_\_\_\_

评 阅 人: \_\_\_\_\_

2017 年 4 月

## 厦门大学学位论文原创性声明

本人呈交的学位论文是本人在导师指导下，独立完成的研究成果。  
本人在论文写作中参考其他个人或集体已经发表的研究成果，均在文  
中以适当方式明确标明，并符合法律规范和《厦门大学研究生学术活  
动规范（试行）》。

另外，该学位论文为( )课题(组)  
的研究成果，获得( )课题(组)经费或实验室的  
资助，在( )实验室完成。（请在以上括号内填写  
课题或课题组负责人或实验室名称，未有此项声明内容的，可以不作  
特别声明。）

声明人（签名）：

2017 年 月 日

## 厦门大学学位论文著作权使用声明

本人同意厦门大学根据《中华人民共和国学位条例暂行实施办法》等规定保留和使用此学位论文，并向主管部门或其指定机构送交学位论文（包括纸质版和电子版），允许学位论文进入厦门大学图书馆及其数据库被查阅、借阅。本人同意厦门大学将学位论文加入全国博士、硕士学位论文共建单位数据库进行检索，将学位论文的标题和摘要汇编出版，采用影印、缩印或者其它方式合理复制学位论文。

本学位论文属于：

- ( ) 1.经厦门大学保密委员会审查核定的保密学位论文，于 年 月 日解密，解密后适用上述授权。  
( ) 2.不保密，适用上述授权。

(请在以上相应括号内打“√”或填上相应内容。保密学位论文应是已经厦门大学保密委员会审定过的学位论文，未经厦门大学保密委员会审定的学位论文均为公开学位论文。此声明栏不填写的，默认为公开学位论文，均适用上述授权。)

声明人（签名）：

2017 年 月

## 摘要

在许多研究领域，数据呈现出多模态的结构特性，张量可以清晰完整地表示这类数据。而张量数据的向量化处理丢失了大量的数据结构信息，甚至造成维数灾和过拟合现象。张量分析的方法直接输入张量数据，能够有效保持数据的结构信息；此外，以张量数据为基础的模型和算法，在减少参变量数目的同时，缓解向量方法在模型学习时极易出现的过拟合现象，这意味着张量模型处理高维小样本问题更加有效，为分析高维向量数据提供了一种新的思路。

基于张量的数据分析方法具有更加广泛的应用前景，所以本文探讨两类复杂关联数据的张量模型及其应用，一类具有天然的张量结构，一类可以转化为张量进行处理。本文的研究内容主要包括以下几个方面：

1. 对象关联数据的张量模型及其应用研究。传统方法破坏了张量数据的结构特性，张量可以有效表示该类复杂关联数据。以社会标注系统为背景，用张量及张量分解模型研究该系统的高阶结构特性和统计特性。以系统中的用户、资源、标签三大对象作为三个维度，引入权重值区分“用户-标签-资源”三元关系的关联强度，建立三阶权值张量模型，张量分解后得到最优核张量和三个维度的特征矩阵，以及新的三元关系的评分值，根据最优的评分值产生推荐列表，向用户推荐资源或标签。

2. 属性关联数据（向量数据）的张量模型及其应用研究。在实际问题中，存在很多具有关联关系的属性数据，通常采用向量方法进行处理，向量维数过高时容易产生过拟合现象，张量分析方法在减少模型参变量的同时缓解或避免了过拟合现象。本文基于向量数据的学习方法，主要是支持向量机的分类、回归和特征选择方法，自然地推广到张量空间，得到支持张量机的分类模型、回归模型和特征选择方法，实验对比证明，支持张量机模型不仅可以分析向量数据，还可以有效缓解高维小样本问题。

3. 实证研究：基于张量空间模型的网络舆情分析。通过张量方法向用户个性化推荐网络舆情热点话题，并实现网络舆情文本的有效识别和自动分类。在实现个性化推荐时，用张量对用户和网络舆情建模，通过张量分解分析用户的兴趣倾向，进而向用户个性化推荐网络舆情热点话题，实验结果表明，引入权重值的张

量分解模型将进一步提高推荐资源的准确率，使个性化推荐结果更加精确。在实现网络舆情文本的有效识别时，将文本表示为 $20 \times 20$ 的二阶张量，构造张量分类器对网络舆情文本进行分类，实验结果表明，支持张量机模型在解决网络舆情文本的高维小样本问题和数据偏斜问题时具有更好的泛化性能。可见，将张量空间模型应用于网络舆情分析领域具有广泛的应用价值。

**关键词：**多模态数据；高维数据；张量；张量分析

厦门大学博硕士论文摘要库

## Abstract

In many areas of research, data presents multimodal structure property, this type of data can be expressed clearly and completely by tensor. Vectorization of tensor data lose many data structure information, even result in curse of dimensionality and over-fitting. While tensor analysis methods enter tensor data directly, keeping data structure information effectively; models and algorithms based on tensor data, reduce parameter variables, ease overfitting phenomenon of vector methods in model learning process, this means tensor model is more effective for high dimension small sample problem, provide a new idea for analyzing high dimension vector data.

Data analysis methods based on tensor have wider applications, so, this paper will discuss tensor models of two types of complex correlated data and their applications, one has natural tensor structure, another can be converted to tensor for processing. Research mainly includes following aspects:

1. Tensor model of objects correlated data and application research. Traditional methods destroy structural property of tensor data; tensor can effectively express this type of complex correlated data. Taking social tagging system as background, using tensor and tensor decomposition model to study higher order structural and statistical characteristics of the system. Taking user, resource, tag the three objects as three dimensions, introducing weights to differentiate “user-label-resource” the ternary relation intensity, establishing the three-order weights tensor model, we can get optimal kernel tensor、feature matrix of three dimensions and rating scores of new ternary relations by tensor decomposition, get recommendation list by rating scores, and recommend resources or tags to users.

2. Tensor models of attributes correlated data (vector data) and application research. In realistic problems, there are many attribute data which has associated relations, often handled by vector methods, it will appear overfitting phenomenon when vector dimension is too high, tensor analysis methods ease or avoid occurrence of overfitting by reducing model parameter variables. The paper naturally generalizes to tensor space based on learning methods of vector data, primarily the classification,

regression, and feature selection of support vector machine, finally obtain tensor classification model, tensor regression model, and tensor feature selection methods, example comparisons prove that support tensor machine model not only can be used to analyze vector data, but also can alleviate the high dimensional small sample problems effectively.

3. Empirical study: network public opinion analysis based on tensor space model. By tensor methods, we personalized recommend hot topics of network public opinion to users, and realize the effective identification and automatic classification of network public opinion texts. In the realization of personalized recommendation, using tensor to establish model for user and network public opinion, analyzing interest tendency of user by tensor decomposition, and then recommending hot topics of network public opinion to users individually, experiment results show that introducing weights will increase accuracy of resource recommendation and make personalized recommendation more precise. In the realization of effective identification of network public opinion texts, text is expressed by two order tensor of  $20 \times 20$ , constructing tensor classifier to classify network public opinion texts, experiment results show that support tensor machine model has better performance in solving high-dimensional small sample problem and data skew problem of network public opinion texts. It can be seen that tensor space model has wide application value in the field of network public opinion analysis.

**Keywords:** Multimodal data; High dimensional data; Tensor; Tensor analysis.

# 目录

<b>第一章 绪论.....</b>	<b>1</b>
<b>1.1 研究背景及意义 .....</b>	<b>1</b>
<b>1.2 张量表示与分析 .....</b>	<b>2</b>
1.2.1 从向量表示到张量表示.....	2
1.2.2 从矩阵分解到张量分解.....	4
1.2.3 从向量学习到张量学习.....	6
<b>1.3 国内外文献综述 .....</b>	<b>7</b>
1.3.1 张量的起源与发展.....	7
1.3.2 张量分解的研究现状.....	7
1.3.3 张量学习的研究现状.....	9
<b>1.4 本文的内容框架、个人贡献与不足 .....</b>	<b>10</b>
1.4.1 本文的内容框架.....	10
1.4.2 个人贡献与不足.....	11
<b>第二章 相关张量基本理论.....</b>	<b>13</b>
<b>2.1 张量.....</b>	<b>13</b>
2.1.1 张量的定义.....	13
2.1.2 张量的表示.....	13
<b>2.2 相关张量基本运算 .....</b>	<b>15</b>
2.2.1 张量的矩阵化.....	15
2.2.2 张量的模乘.....	16
2.2.3 张量内积与张量范数.....	17
2.2.4 张量外积.....	17
<b>2.3 张量分解.....</b>	<b>18</b>
2.3.1 Tucker 分解 .....	18
2.3.2 CP 分解.....	20
2.3.3 高阶奇异值分解（HOSVD） .....	21
<b>第三章 复杂对象关联数据的张量模型应用研究.....</b>	<b>23</b>

<b>3.1 社会标注系统概述 .....</b>	<b>23</b>
3.1.1 社会标注系统介绍.....	23
3.1.2 社会标注系统的张量模型.....	24
<b>3.2 传统分析方法的局限性 .....</b>	<b>25</b>
<b>3.3 “用户-资源-标签” 三元关系的加权评分.....</b>	<b>26</b>
3.3.1 用户、资源、标签的权值计算.....	26
3.3.2 “用户-资源-标签” 三元元组的加权评分.....	27
<b>3.4 基于加权元组的社会标注系统张量模型 .....</b>	<b>28</b>
3.4.1 张量模型的定义.....	28
3.4.2 三元组初始权重计算.....	30
3.4.3 张量分解.....	31
3.4.4 生成推荐结果.....	33
<b>3.5 本章小节.....</b>	<b>36</b>
<b>第四章 复杂属性关联数据的张量模型应用研究.....</b>	<b>37</b>
<b>4.1 支持向量机理论 .....</b>	<b>37</b>
4.1.1 支持向量机分类.....	37
4.1.2 支持向量机回归.....	40
4.1.3 支持向量机特征选择.....	42
<b>4.2 传统向量方法的局限性 .....</b>	<b>42</b>
<b>4.3 支持张量机理论 .....</b>	<b>43</b>
4.3.1 支持张量机分类.....	43
4.3.2 支持张量机回归.....	45
4.3.3 支持张量机特征选择.....	48
<b>4.4 向量数据的张量表示 .....</b>	<b>48</b>
<b>4.5 实验对比.....</b>	<b>49</b>
4.5.1 分类实验对比.....	49
4.5.2 回归实验对比.....	54
4.5.3 特征选择实验对比.....	58
<b>4.6 本章小节.....</b>	<b>59</b>

<b>第五章 实证研究：基于张量空间模型的网络舆情分析 .....</b>	<b>61</b>
<b>5.1 研究背景及意义 .....</b>	<b>61</b>
<b>5.2 网络舆情文本的预处理 .....</b>	<b>62</b>
5.2.1 舆情文本的特征表示.....	62
5.2.2 舆情文本的特征选择.....	63
<b>5.3 网络舆情文本的张量表示 .....</b>	<b>64</b>
<b>5.4 多类别网络舆情文本分类方法简介 .....</b>	<b>65</b>
5.4.1 “一对多”多类分类方法.....	65
5.4.2 “一对一”多类分类方法.....	65
<b>5.5 实验一：基于张量分解的网络舆情话题推荐 .....</b>	<b>66</b>
5.5.1 数据准备.....	66
5.5.2 评价指标.....	66
5.5.3 实验结果与分析.....	67
<b>5.6 实验二：基于支持张量机的网络舆情文本分类 .....</b>	<b>68</b>
5.6.1 数据准备.....	68
5.6.2 评价指标.....	68
5.6.3 实验结果与分析.....	69
<b>5.7 本章小节.....</b>	<b>75</b>
<b>第六章 总结与展望.....</b>	<b>77</b>
<b>6.1 论文工作总结.....</b>	<b>77</b>
<b>6.2 未来研究展望.....</b>	<b>78</b>
<b>参考文献.....</b>	<b>79</b>
<b>致 谢 .....</b>	<b>83</b>

厦门大学博硕士论文摘要库

## Contents

<b>Chapter 1 Introduction .....</b>	<b>1</b>
<b>1.1 Research Background and Significance .....</b>	<b>1</b>
<b>1.2 Tensor Representation and Analysis.....</b>	<b>2</b>
1.2.1 From Vector To Tensor Representation.....	2
1.2.2 From Matrix Decomposition To Tensor Decomposition .....	4
1.2.3 From Vector Learning To Tensor Learning .....	6
<b>1.3 Literature Review .....</b>	<b>7</b>
1.3.1 The Origin and Development of Tensor .....	7
1.3.2 Research Status of Tensor Decomposition.....	7
1.3.3 Research Status of Tensor Learning.....	9
<b>1.4 Framework、 Personal Contributions and Deficiencies.....</b>	<b>10</b>
1.4.1 Framework .....	10
1.4.2 Personal Contributions and Deficiencies .....	11
<b>Chapter 2 Related Basic Theory of Tensor.....</b>	<b>13</b>
<b>2.1 Tensor.....</b>	<b>13</b>
2.1.1 Definition of Tensor .....	13
2.1.2 Representation of Tensor.....	13
<b>2.2 Related Basic Operations of Tensor .....</b>	<b>15</b>
2.2.1 Tensor Matricization .....	15
2.2.2 Tensor Mode Multiplication.....	16
2.2.3 Tensor Inner Product and Tensor Norm .....	17
2.2.4 Tensor Outer Product .....	17
<b>2.3 Tensor Decomposition.....</b>	<b>18</b>
2.3.2 Tucker Decomposition .....	18
2.3.3 CP Decomposition .....	20
2.3.4 Higher Order Singular Value Decomposition(HOSVD).....	21

## **Chapter 3 Application Research on Tensor Model of Complex Objects Correlated Data ..... 23**

<b>3.1 Overview of Social Tagging System.....</b>	<b>23</b>
3.1.1 Introduction of Social Tagging System.....	23
3.1.2 Tensor Model of Social Tagging System .....	24
<b>3.2 Limitation of Traditional Analysis Methods .....</b>	<b>25</b>
<b>3.3 Weighted Score of "User-Item-Tag" Ternary Relation .....</b>	<b>26</b>
3.3.1 Weight Calculation of User、 Item、 Tag.....	26
3.3.2 Weighted Rating Score of "User-Item-Tag" ternary Tuple .....	27
<b>3.4 Tensor Model of Social Tagging System Based On Weighted Tuples .....</b>	<b>28</b>
3.4.1 Definition of Tensor Model.....	28
3.4.2 Original Weight Calculation of Ternary Tuples .....	30
3.4.3 Tensor Decomposition .....	31
3.4.4 Recommendation Results.....	33
<b>3.5 Conclusion .....</b>	<b>36</b>

## **Chapter 4 Application Research on Tensor Model of Complex Attributes Correlated Data ..... 37**

<b>4.1 Support Vector Machine Theory .....</b>	<b>37</b>
4.1.1 Classification of Support Vector Machine .....	37
4.1.2 Regression of Support Vector Machine .....	40
4.1.3 Feature Selection of Support Vector Machine .....	42
<b>4.2 Limitation of Traditional Vector Methods .....</b>	<b>42</b>
<b>4.3 Support Tensor Machine Theory .....</b>	<b>43</b>
4.3.1 Classification of Support Tensor Machine .....	43
4.3.2 Regression of Support Tensor Machine .....	45
4.3.3 Feature Selection of Support Tensor Machine .....	48
<b>4.4 Tensor Representation of Vector Data.....</b>	<b>48</b>
<b>4.5 Experiment Comparison .....</b>	<b>49</b>

4.5.1 Classification Experiment Comparison .....	49
4.5.2 Regression Experiment Comparison.....	54
4.5.3 Feature Selection Experiment Comparison .....	58
<b>4.6 Conclusion .....</b>	<b>59</b>
<b>Chapter 5 Empirical Research: Network Public Opinion Analysis Based On Tensor Space Model.....</b>	<b>61</b>
<b>5.1 Research Background and Significance .....</b>	<b>61</b>
<b>5.2 Preprocessing of Network Public Opinion Text .....</b>	<b>62</b>
5.2.1 Feature Representation of Public Opinion Text .....	62
5.2.2 Feature Selection of Public Opinion Text .....	63
<b>5.3 Tensor Representation of Network Public Opinion Text.....</b>	<b>64</b>
<b>5.4 Methods Introduction of Multi-class Network Public Opinion Text Classification.....</b>	<b>65</b>
5.4.1 "One-Versus-Rest" Multi-class Classification Method .....	65
5.4.2 "One-Versus-One" Multi-class Classification Method .....	65
<b>5.5 Experiment one: Network Public Opinion Topic Recommendation Based on Tensor Decomposition.....</b>	<b>66</b>
5.5.1 Data Preparation.....	66
5.5.2 Evaluation Index .....	66
5.5.3 Experiment Result and Analysis .....	67
<b>5.6 Experiment Two: Text Categorization of Network Public Opinion Based on Support Tensor Machine .....</b>	<b>68</b>
5.6.1 Data Preparation.....	68
5.6.2 Evaluation Index .....	68
5.6.3 Experiment Result and Analysis .....	69
<b>5.7 Conclusion .....</b>	<b>75</b>
<b>Chapter 6 Conclusions and Future Studies .....</b>	<b>77</b>
<b>6.1 Research Conclusions .....</b>	<b>77</b>

<b>6.2 Future Studies .....</b>	<b>78</b>
<b>References .....</b>	<b>79</b>
<b>Acknowledge .....</b>	<b>83</b>

厦门大学博硕士论文摘要库

# 第一章 绪论

## 1.1 研究背景及意义

随着信息科技和互联网技术的快速发展，我们进入大数据时代。大数据时代的到来使传统的数据处理方式开始面临严峻挑战。大数据不仅指数据的数据量庞大，还包括多模态、多属性的复杂关联数据。科学、网络及社会领域中许多数据呈现出多模态的结构特性，如网络流量、社会关系网、彩色图像等，它们的数据结构特性是高阶的，即张量特性，或者可以被组织成张量结构，张量是对这类数据自然而本质的表达方式，可以保留数据元素相互间的内在关联和空间拓扑结构。此外，张量作为向量和矩阵的推广，具有良好的计算特性和分析表达能力。所以，张量是处理多模态高维数据的有力工具。

现实世界中这些具有高阶特性的数据用张量这种数据结构可以直观精确地表达，如果想要进一步分析这类数据，传统方法是将其转化为向量形式，用基于向量的机器学习方法进行分析。对于张量数据的直接转化将导致高维向量的产生，在学习过程中容易出现维数灾和过拟合现象，而且当转换为向量时，原始数据的高阶结构特性遭到破坏，使数据包含的一些重要信息也不复存在。因此，为了应对处理这类高阶张量数据时遇到的问题，许多研究学者深入探索了基于张量数据的机器学习方法，这些方法对张量数据并不作处理，而是将其直接输入，这样，这些数据的高阶结构特性得到了保持，除此之外，张量算法还具有一定的特殊性，直接输入张量数据训练张量模型时需要确定的参变量数目会减少，缓解或者避免了过拟合现象。

除了多模态结构的张量数据外，现实中更多的是包含多个属性的向量数据。在实际问题中，一个对象往往包含多个属性，举一个简单的例子，人的属性包括其名字、年龄、性别、人际关系等，商品的属性包括其名称、类别、材质、特殊说明等。处理这类数据的方法一般是基于向量的统计学习方法，但是大数据时代数据量爆炸式的增长，使向量数据的维数呈几何级数式增长，而在现实研究中，存在很多我们人为无法控制和解决的因素，如样本的获取方式、成本以及一些不

Degree papers are in the "[Xiamen University Electronic Theses and Dissertations Database](#)". Full texts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to [etd@xmu.edu.cn](mailto:etd@xmu.edu.cn) for delivery details.

厦门大学博硕士论文全文数据库