

SpringerBriefs in Computer Science

For further volumes:
<http://www.springer.com/series/10028>

David B. Skillicorn

Understanding High-Dimensional Spaces

David B. Skillicorn
School of Computing
Queen's University
Kingston, ON
Canada

ISSN 2191-5768 ISSN 2191-5776 (electronic)
ISBN 978-3-642-33397-2 ISBN 978-3-642-33398-9 (eBook)
DOI 10.1007/978-3-642-33398-9
Springer Heidelberg New York Dordrecht London

Library of Congress Control Number: 2012948185

ACM Computing Classification (1998): H.2, E.1, I.2

© The Author 2012

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Preface

High-dimensional spaces arise naturally as a way of modelling datasets with many attributes. Such a dataset can be directly represented in a space spanned by the attributes, with each record of the dataset represented as a point in the space with its position depending on its attribute values. Such spaces are not easy to work with because of their high dimensionality: our intuition about space is not reliable, and measures such as distance do not provide as clear information as we might expect.

High-dimensional spaces have not received as much attention as their applications deserve, partly for these reasons. Some areas where there has been substantial research are: images and video, with high-dimensional representations based on one attribute per pixel; and spaces with highly non-convex clusters. For images and video, the high dimensionality is an artifact of a direct representation, but the inherent dimensionality is usually much lower, and easily discoverable. Spaces with a few highly non-convex clusters do occur, but are not typical of the kind of datasets that arise in practice.

There are at least three main areas where complex high dimensionality and large datasets arise naturally. The first is data collected by online retailers (e.g. Amazon), preference sites (e.g. Pandora), social media sites (e.g. Facebook), and the customer relationship data of all large businesses. In these applications, the amount of data available about any individual is large but also sparse. For example, a site like Pandora has preference information for every song that a user has listened to, but this is still a tiny fraction of all of the songs that the site cares about. A site like Amazon has information about which items any customer has bought, but this is a small fraction of what is available.

The second is data derived from text (and speech). The word usage in a set of documents produces data about the frequency with which each word is used. As in the first case, all of the words used in a given document are visible, but there are always many words that are not used at all in it. So such datasets are large (because easy to construct), wide (because languages contain many words), and sparse (because any document uses a small fraction of the possible words).

The third is data collected for a security, defence, law enforcement or intelligence purpose; or collected about computer networks for cybersecurity. Such

datasets are large and wide because of the need to enable as good solutions as possible by throwing the data collection net wide. This third domain differs from the previous two because of greater emphasis on the anomalous or outlying parts of the data rather than the more central and common place.

High-dimensional datasets are usually analyzed in two ways: by finding the set of clusters they contain; or by looking for the outliers—almost two sides of the same coin. However, these simple strategies conceal subtleties that are often ignored. A cluster cannot really be understood without seeing its relationships to other clusters “around” it; and outliers cannot be understood without understanding both the clusters that they are nearest to, and what other outliers are “around” them. The development of the idea of local outliers has helped with this latter issue, but is still weak because a local outlier is defined only with respect to its nearest non-outlying cluster.

In this book we introduce two ideas that are not completely new, but which have not received as much attention as they should have, and for which the research results are partial and scattered. In essence, we suggest a new way of thinking about how to understand high-dimensional spaces using two models: the *skeleton* which relates the clusters to one another, and *boundaries in empty space* which provides a new perspective on outliers, and on outlying regions.

This book should be useful to those who are analyzing high-dimensional spaces using existing tools, and who feel that they are not getting as much out of the data as they could; also their managers who are trying to understand the path forward in terms of what is possible, and how they might get there. The book assumes either that the reader has a reasonable grasp of mainstream data mining tools and techniques, or does not need to get into the weeds of the technology but needs a sense of the landscape. The book may also be useful for graduate students and other researchers who are looking for open problems, or new ways to think about and apply older techniques.

Acknowledgments

My greatest debt is to Mike Bourassa. Discussions in the course of his doctoral research first surfaced many of the ideas described here, and we had many long conversations about what it meant to be interesting, before we converged on the meaning that is expanded here. I am also grateful to all my students who, by providing an audience for me to explain both simple and complex ideas, help me to understand them better.

Kingston, April 2012

David Skillicorn

Contents

1	Introduction	1
1.1	A Natural Representation of Data Similarity	3
1.2	Goals	8
1.3	Outline	10
2	Basic Structure of High-Dimensional Spaces	13
2.1	Comparing Attributes	13
2.2	Comparing Records	14
2.3	Similarity	14
2.4	High-Dimensional Spaces	16
2.5	Summary	18
3	Algorithms	19
3.1	Improving the Natural Geometry	19
3.1.1	Projection	20
3.1.2	Singular Value Decompositions	20
3.1.3	Random Projections	22
3.2	Algorithms that Find Standalone Clusters	23
3.2.1	Clusters Based on Density	23
3.2.2	Parallel Coordinates	24
3.2.3	Independent Component Analysis	24
3.2.4	Latent Dirichlet Allocation	25
3.3	Algorithms that Find Clusters and Their Relationships	25
3.3.1	Clusters Based on Distance	25
3.3.2	Clusters Based on Distribution	26
3.3.3	Semidiscrete Decomposition	27
3.3.4	Hierarchical Clustering	29
3.3.5	Minimum Spanning Tree with Collapsing	29
3.4	Overall Process for Constructing a Skeleton	30

3.5	Algorithms that Wrap Clusters	31
3.5.1	Distance-Based	32
3.5.2	Distribution-Based.	32
3.5.3	1-Class Support Vector Machines	32
3.5.4	Autoassociative Neural Networks	33
3.5.5	Covers	34
3.6	Algorithms to Place Boundaries Between Clusters.	34
3.6.1	Support Vector Machines	35
3.6.2	Random Forests	35
3.7	Overall Process for Constructing Empty Space	36
3.8	Summary	37
4	Spaces with a Single Center	39
4.1	Using Distance	39
4.2	Using Density	40
4.3	Understanding the Skeleton	42
4.4	Understanding Empty Space	43
4.5	Summary	45
5	Spaces with Multiple Centers	47
5.1	What is a Cluster?	48
5.2	Identifying Clusters	50
5.2.1	Clusters Known Already	50
5.3	Finding Clusters	50
5.4	Finding the Skeleton	55
5.5	Empty Space.	58
5.5.1	An Outer Boundary and Novel Data	58
5.5.2	Interesting Data	60
5.5.3	One-Cluster Boundaries	63
5.5.4	One-Cluster-Against-the-Rest Boundaries	63
5.6	Summary	64
6	Representation by Graphs	67
6.1	Building a Graph from Records.	68
6.2	Local Similarities	68
6.3	Embedding Choices	69
6.4	Using the Embedding for Clustering	70
6.5	Summary	71
7	Using Models of High-Dimensional Spaces	73
7.1	Understanding Clusters.	73
7.2	Structure in the Set of Clusters	76
7.2.1	Semantic Stratified Sampling	77
7.3	Ranking Using the Skeleton	78

Contents	ix
7.4 Ranking Using Empty Space.	87
7.4.1 Applications to Streaming Data	89
7.4.2 Concealment.	90
7.5 Summary	91
8 Including Contextual Information	93
8.1 What is Context?	93
8.1.1 Changing Data	93
8.1.2 Changing Analyst and Organizational Properties.	94
8.1.3 Changing Algorithmic Properties	95
8.2 Letting Context Change the Models.	95
8.2.1 Recomputing the View	95
8.2.2 Recomputing Derived Structures	96
8.2.3 Recomputing the Clustering	97
8.3 Summary	98
9 Conclusions	99
References	103
Index	107