

LUCAS FERNANDES BRUNIALTI

Resolução do problema de coagrupamento  
em matrizes de dados esparsas usando  
fatoração de matrizes

São Paulo

2016

**LUCAS FERNANDES BRUNIALTI**

**Resolução do problema de coagrupamento em  
matrizes de dados esparsas usando fatoração de  
matrizes**

Versão original

Dissertação apresentada à Escola de Artes, Ciências e Humanidades da Universidade de São Paulo para obtenção do título de Mestre em Ciências pelo Programa de Pós-graduação em Sistemas de Informação.

Área de concentração: Metodologia e Técnicas da Computação

Orientador: Profa. Dra. Sarajane Marques Peres

São Paulo

2016

Ficha catalográfica

Dissertação de autoria de Lucas Fernandes Brunialti, sob o título “**Resolução do problema de coagrupamento em matrizes de dados esparsas usando fatoração de matrizes**”, apresentada à Escola de Artes, Ciências e Humanidades da Universidade de São Paulo, para obtenção do título de Mestre em Ciências pelo Programa de Pós-graduação em Sistemas de Informação, na área de concentração Sistemas de Informação, aprovada em \_\_\_ de \_\_\_\_\_ de \_\_\_\_\_ pela comissão julgadora constituída pelos doutores:

**Prof. Dr.** \_\_\_\_\_

Presidente

Instituição: \_\_\_\_\_

**Prof. Dr.** \_\_\_\_\_

Instituição: \_\_\_\_\_

**Prof. Dr.** \_\_\_\_\_

Instituição: \_\_\_\_\_

*Escreva aqui sua dedicatória, se desejar, ou remova esta página...*



*“Escreva aqui uma epígrafe, se desejar, ou remova esta página...”*

*(Autor da epígrafe)*

## Resumo

BRUNIALTI, Lucas Fernandes. **Resolução do problema de coagrupamento em matrizes de dados esparsas usando fatoração de matrizes**. 2016. 105 f.

Dissertação (Mestrado em Ciências) – Escola de Artes, Ciências e Humanidades, Universidade de São Paulo, São Paulo, 2016.

Coagrupamento é uma estratégia para análise de dados capaz de encontrar grupos de objetos, então denominados cogrupos, que são similares entre si de acordo com um subconjunto dos seus atributos descritivos, e assim, um objeto pode pertencer a mais de um grupo se subconjuntos diferentes de atributos forem considerados. Essa característica pode ser particularmente útil para aplicações nas quais a similaridade parcial entre objetos faz sentido, conferindo ao resultado da análise de dados algumas características interessantes como **serendipidade** ou flexibilidade no modelo de agrupamento. Contextos de aplicação caracterizados por apresentar subjetividade, como mineração de textos, são candidatos a serem submetidos à estratégia de coagrupamento; a flexibilidade em associar textos de acordo com características parciais representa um tratamento mais adequado à tal subjetividade. Entretanto, análise de grupos considerando dados textuais representa um contexto no qual existe o **problema de esparsidade de dados, que precisa ser adequadamente tratado para que os bons resultados sejam obtidos**. Um método para implementação de coagrupamento capaz de lidar com esse tipo de dados é a fatoração de matrizes. Nesta dissertação de mestrado são propostas duas estratégias para coagrupamento baseadas em fatoração de matrizes, capazes de encontrar cogrupos organizados com sobreposição em uma matriz esparsa de valores reais positivos. As estratégias são apresentadas em termos de suas definições formais e seus algoritmos para implementação. Resultados experimentais são fornecidos a partir de problemas baseados em conjuntos de dados sintéticos e em conjuntos de dados reais contextualizados na área de mineração de textos. Os resultados confirmam a hipótese de que as estratégias propostas são capazes de descobrir cogrupos com sobreposição, e que tal organização de cogrupos fornece informação detalhada, e portanto de valor diferenciado, para a mineração de textos.

Palavras-chaves: Coagrupamento. Fatoração de Matrizes. **Esparsidade**. Análise de Agrupamento. Mineração de Texto.



## Abstract

BRUNIALTI, Lucas Fernandes. **Matrix Factorization for cogruping in sparse data matrices**. 2016. 105 p. Dissertation (Master of Science) – School of Arts, Sciences and Humanities, University of São Paulo, São Paulo, DefenseYear.

Coagrupamento é uma estratégia para análise de dados capaz de encontrar grupos de objetos, então denominados cogrupos, que são similares entre si de acordo com um subconjunto dos seus atributos descritivos, e assim, um objeto pode pertencer a mais de um grupo se subconjuntos diferentes de atributos forem considerados. Essa característica pode ser particularmente útil para aplicações nas quais a similaridade parcial entre objetos faz sentido, conferindo ao resultado da análise de dados algumas características interessantes como **serendipidade** ou flexibilidade no modelo de agrupamento. Contextos de aplicação caracterizados por apresentar subjetividade, como mineração de textos, são candidatos a serem submetidos à estratégia de coagrupamento; a flexibilidade em associar textos de acordo com características parciais representa um tratamento mais adequado à tal subjetividade. Entretanto, análise de grupos considerando dados textuais representa um contexto no qual existe o **problema de esparsidade de dados, que precisa ser adequadamente tratado para que os bons resultados sejam obtidos**. Um método para implementação de coagrupamento capaz de lidar com esse tipo de dados é a fatoração de matrizes. Nesta dissertação de mestrado são propostas duas estratégias para coagrupamento baseadas em fatoração de matrizes, capazes de encontrar cogrupos organizados com sobreposição em uma matriz esparsa de valores reais positivos. As estratégias são apresentadas em termos de suas definições formais e seus algoritmos para implementação. Resultados experimentais são fornecidos a partir de problemas baseados em conjuntos de dados sintéticos e em conjuntos de dados reais contextualizados na área de mineração de textos. Os resultados confirmam a hipótese de que as estratégias propostas são capazes de descobrir cogrupos com sobreposição, e que tal organização de cogrupos fornece informação detalhada, e portanto de valor diferenciado, para a mineração de textos.

Keywords: Coclustering. Matrix Factorization. Sparsity. Clustering Analysis. Text Mining.

## Lista de figuras

Figura 1 – Fatoração da matriz original de dados $X$ em três outras matrizes: $U$ , $S$ e $V$ . . . . .	33
Figura 2 – A reconstrução da primeira linha $\mathbf{x}_1$ de $X$ , através da multiplicação da matriz indicadora de grupos de linhas $U$ pela matriz dos protótipos de linhas ( $SV^T$ ). . . . .	34
Figura 3 – Base de protótipos obtidas com FMN sem restrições (a) e com restrições de ortogonalidade nas matrizes . . . . .	37
Figura 4 – Representação gráfica do problema de coagruamento com sobreposição unidimensional e contextualização do domínio de documentos (notícias) . . . . .	49
Figura 5 – Fatoração da matriz original de dados $X$ em cinco outras matrizes: $U$ , $S$ , $V_1$ , $V_2$ e $V_3$ . . . . .	49
Figura 6 – Dados sintéticos gerados a partir das diferentes estruturas de cogrupos. (a) Um único cogrupo. (b) Cogrupos com linhas e colunas sem intersecção. (c) Cogrupos com estrutura em xadrez. (d) Cogrupos sem intersecção nas linhas e com intersecção nas colunas. (e) Cogrupos com intersecção nas linhas e sem intersecção nas colunas. . . . .	61
Figura 7 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo <i>k-means</i> . . . . .	65
Figura 8 – Resultado da reconstrução da base de dados (e) utilizando o algoritmo <i>k-means</i> com $k = 5$ . . . . .	66
Figura 9 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo <i>fuzzy k-means</i> . . . . .	67
Figura 10 – Resultado da reconstrução da base de dados (e) utilizando o algoritmo <i>fuzzy k-means</i> com $k = 5$ . . . . .	67
Figura 11 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo <i>ONMTF</i> . . . . .	68
Figura 12 – Resultado da reconstrução da base de dados (d) com $k = 5$ e (e) com $l = 5$ , respectivamente, utilizando o algoritmo <i>ONMTF</i> . . . . .	69

Figura 13 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo <i>FNMTF</i> . . . . .	70
Figura 14 – Resultado da reconstrução da base de dados (d) com $k = 5$ e (e) com $l = 5$ , respectivamente, utilizando o algoritmo <i>FNMTF</i> . . . . .	71
Figura 15 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo <i>OvNMTF</i> . . . . .	72
Figura 16 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo <i>BinOvNMTF</i> . . . . .	73
Figura 17 – Resultado da reconstrução da base de dados (e) utilizando o algoritmo <i>BinOvNMTF</i> com $k = 5$ . . . . .	73
Figura 18 – Distribuições em formato de diagrama de caixa do Índice de Rand para as melhores configurações de $l$ e representação para cada algoritmo na base de dados <i>IG toy</i> . . . . .	85
Figura 19 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada para as melhores configurações de $l$ e representação para cada algoritmo na base de dados <i>IG toy</i> . . . . .	85
Figura 20 – Distribuições em formato de diagrama de caixa do Índice de Rand de todas as configurações para cada algoritmo na base de dados <i>IG toy</i> . . . . .	86
Figura 21 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada de todas as configurações para cada algoritmo na base de dados <i>IG toy</i> . . . . .	86
Figura 22 – Distribuições em formato de diagrama de caixa do Índice de Rand para as melhores configurações de $l$ e representação para cada algoritmo na base de dados <i>NIPS</i> . . . . .	88
Figura 23 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada para as melhores configurações de $l$ e representação para cada algoritmo na base de dados <i>NIPS</i> . . . . .	89
Figura 24 – Distribuições em formato de diagrama de caixa do Índice de Rand de todas as configurações para cada algoritmo na base de dados <i>NIPS</i> . . . . .	89

Figura 25 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada de todas as configurações para cada algoritmo na base de dados <i>NIPS</i> . . . . .	90
Figura 26 – Exemplo de uma notícia do canal arena e sua disposição quanto aos cogrupos de notícias e palavras. . . . .	93
Figura 27 – Visualização <i>word cloud</i> das top-100 palavras, dispostas em uma área quadrada, para cada cogrupos gerados pelo algoritmo <i>ONMTF</i> . . . . .	95
Figura 28 – Visualização <i>word cloud</i> de palavras para cada cogrupos do cogrupos de notícias “arena”, gerados pelo algoritmo <i>OvNMTF</i> . . . . .	96
Figura 29 – Visualização <i>word cloud</i> de palavras para cada cogrupos do cogrupos de notícias “jovem”, gerados pelo algoritmo <i>OvNMTF</i> . . . . .	96
Figura 30 – Visualização <i>word cloud</i> de palavras para cada cogrupos do cogrupos de notícias “esportes”, gerados pelo algoritmo <i>OvNMTF</i> . . . . .	98

## Lista de algoritmos

Algoritmo 1 – Algoritmo baseado em atualização multiplicativa para solução do BVD . .	36
Algoritmo 2 – Algoritmo baseado em atualização multiplicativa para solução do ONMTF	38
Algoritmo 3 – Algoritmo baseado em atualização multiplicativa e na teoria de de derivação na superfície com restrições (Variedade Stiefel) para solução do ONMTF .	40
Algoritmo 4 – Algoritmo FNMTF . . . . .	44
Algoritmo 5 – Algoritmo baseado em atualização multiplicativa para solução do OvNMTF	54
Algoritmo 6 – Algoritmo iterativo para solução do BinOvNMTF . . . . .	58

## Lista de tabelas

Tabela 1 – Resumo de qualidade de reconstrução: <i>ok</i> - permite reconstrução de forma natural; $\times$ - sem informação sobre interseção parcial; + - preserva informação de interseção parcial . . . . .	64
Tabela 2 – Avaliação da capacidade de quantização segundo o erro de quantização com os melhores destacados em negrito. . . . .	74
Tabela 3 – Estatísticas das bases de dados usadas nos experimentos. . . . .	76
Tabela 4 – Distribuição de notícias por ano (base de dados <i>IG</i> ). . . . .	77
Tabela 5 – Distribuição de notícias por canal (base de dados <i>IG</i> ). . . . .	77
Tabela 6 – Distribuição de trabalhos acadêmicos por ano (base de dados <i>NIPS</i> ). . . . .	77
Tabela 7 – Distribuição de trabalhos acadêmicos por áreas técnicas (base de dados <i>NIPS</i> ). . . . .	78
Tabela 8 – Índice de Rand médio por representação do conjunto de dados <i>IG toy</i> com $k = 3$ , destacando os melhores resultados por representação. . . . .	84
Tabela 9 – Informação Mútua Normalizada média por representação do conjunto de dados <i>IG toy</i> com $k = 3$ , destacando os melhores resultados por representação. . . . .	84
Tabela 10 – Melhores resultados obtidos para o conjunto de dados <i>IG toy</i> com $k = 3$ , destacando os melhores resultados por algoritmo. . . . .	84
Tabela 11 – Melhores resultados do Índice de Rand para experimentos com a base de dados <i>IG</i> . * - testes não finalizados . . . . .	87
Tabela 12 – Índice de Rand médio por representação do conjunto de dados <i>NIPS</i> com $k = 9$ , destacando os melhores resultados por representação. . . . .	87
Tabela 13 – Informação Mútua Normalizada média por representação do conjunto de dados <i>NIPS</i> com $k = 9$ , destacando os melhores resultados por representação. . . . .	87
Tabela 14 – Melhores resultados obtidos para o conjunto de dados <i>NIPS</i> com $k = 9$ , destacando os melhores resultados por algoritmo. . . . .	88
Tabela 15 – Matriz <i>S</i> para o algoritmo <i>ONMTF</i> com $k = 3$ e $l = 5$ na base de dados <i>IG toy</i> normalizada para que a soma da linha seja igual à 1 . . . . .	92
Tabela 16 – top-20 palavras para cada cogruppo de palavras, após a realização do particionamento . . . . .	94

Tabela 17 – Matriz $S$ para execução do algoritmo $OvNMTF$ com $k = 3$ e $l = 2$ na base de dados <i>IG toy</i> . . . . .	96
Tabela 18 – top-20 palavras para cada cogrupo de palavras, após a realização do particionamento . . . . .	97

## **Lista de abreviaturas e siglas**

Sigla/abreviatura 1	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 2	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 3	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 4	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 5	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 6	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 7	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 8	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 9	Definição da sigla ou da abreviatura por extenso
Sigla/abreviatura 10	Definição da sigla ou da abreviatura por extenso



## Lista de símbolos

$\Gamma$	Letra grega Gama
$\Lambda$	Lambda
$\zeta$	Letra grega minúscula zeta
$\in$	Pertence

## Sumário

1	Introdução . . . . .	19
1.1	Definição do problema . . . . .	21
1.1.1	Estruturas de coagrupamentos . . . . .	22
1.1.2	Coagrupamento e fatorização de matrizes . . . . .	22
1.2	Hipótese . . . . .	23
1.3	Objetivos . . . . .	23
1.4	Metodologia . . . . .	24
1.5	Organização do documento . . . . .	25
2	Conceitos Fundamentais . . . . .	26
2.1	Tipos de cogrupos . . . . .	26
2.2	Algoritmos para Biclusterização . . . . .	27
2.3	Avaliação de Biclusterização . . . . .	29
3	Fatoração de matrizes não-negativas para coagrupamento . . . . .	31
3.1	Decomposição de Valores em Blocos para Coagrupamento . . . . .	32
3.2	Fatoração Ortogonal Tripla de Matrizes Não-negativas . . . . .	37
3.3	Fatoração Tripla Rápida de Matrizes Não-negativas . . . . .	41
3.4	Considerações Finais . . . . .	45
4	Fatoração de matrizes não-negativas para coagrupamento com sobreposição unidimensional . . . . .	47
4.1	Fatoração Tripla de Matrizes Não-negativas com Sobreposição Unidimensional . . . . .	50
4.2	Fatoração Binária Tripla de Matrizes Não-negativas com Sobreposição Unilateral . . . . .	55
4.3	Considerações Finais . . . . .	57
5	Experimentos e Resultados . . . . .	60
5.1	Experimentos com Bases de Dados sintéticas . . . . .	61

5.1.1	Análise da reconstrução . . . . .	64
5.1.2	Reconstrução a partir dos resultados do algoritmo <i>k-means</i>	64
5.1.3	Reconstrução a partir dos resultados do algoritmo <i>fuzzy</i> <i>k-means</i> . . . . .	66
5.1.4	Reconstrução a partir dos resultados do algoritmo <i>ONMTF</i>	68
5.1.5	Reconstrução a partir dos resultados do algoritmo <i>FNMTF</i>	69
5.1.6	Reconstrução a partir dos resultados do algoritmo <i>OvNMTF</i>	71
5.1.7	Reconstrução a partir dos resultados do algoritmo <i>Bi-</i> <i>nOvNMTF</i> . . . . .	72
5.1.8	Análise da Capacidade de Quantização . . . . .	74
5.2	Experimentos com Bases de Dados Reais . . . . .	75
5.2.1	Descrição das bases de dados . . . . .	75
5.2.2	Pré-processamento . . . . .	78
5.2.3	Experimentos quantitativos . . . . .	80
5.2.3.1	Configuração dos experimentos . . . . .	80
5.2.3.2	Resultados . . . . .	82
5.2.4	Análise qualitativa . . . . .	90
5.2.4.1	Análise de dados utilizando <i>ONMTF</i> . . . . .	91
5.2.4.2	Análise de dados utilizando <i>OvNMTF</i> . . . . .	96
	<b>Referências<sup>1</sup></b> . . . . .	99
	<b>Apêndice A—Mineração de Texto</b> . . . . .	103
A.1	Tarefas de pré-processamento . . . . .	103
A.1.1	Representação textual . . . . .	103
A.1.2	Tokenização . . . . .	104
A.1.3	Filtragem . . . . .	104
A.1.4	Stemming . . . . .	105
A.1.5	Redução de Dimensionalidade . . . . .	105

---

<sup>1</sup> De acordo com a Associação Brasileira de Normas Técnicas. NBR 6023.

# 1 Introdução

Segundo [Jain, Murty e Flynn \(1999\)](#), a análise de agrupamento pode ser vista como uma tarefa exploratória que tem o objetivo de organizar uma coleção de dados em um conjunto de grupos, segundo a similaridade ou dissimilaridade existente entre esses dados. Tradicionalmente, os métodos usados para análise de agrupamento são desenvolvidos para minimizar a similaridade intragrupo e maximizar a similaridade intergrupos; e precisam encontrar uma organização “natural” em grupos que acomode cada dado do conjunto sob análise em um grupo específico.

Estratégias de diferentes naturezas – particional, hierárquica, baseada em densidade, etc ([XU; WUNSCH, 2005](#); [HAN; KAMBER, 2006](#)), podem ser usadas para alcançar o objetivo da análise de agrupamento, e cada uma delas possui características que as fazem mais ou menos suscetíveis para conjuntos de dados de diferentes naturezas. Ainda, sob o contexto clássico da tarefa de agrupamento, os métodos precisam lidar com a similaridade intrínseca entre os dados tomando como base a comparação de todas as suas características descritivas e, de alguma forma, precisam ser capazes de descobrir quais características de fato tornam dados em um grupo de dados mais similares entre si.

Ao longo do tempo, pesquisadores da área de análise de agrupamento vêm propondo flexibilizações na definição da tarefa de agrupamento de forma a adequá-la a contextos mais realísticos nos quais a organização natural dos dados em um conjunto de dados não pressupõe restrições como a pertinência de um dado a um único grupo ou a possibilidade de agrupar dados de acordo com similaridades em subconjuntos de atributos descritivos ([referência, referência, referência](#)). Essa forma de tratar a tarefa de agrupamento permite melhorias no processo de descoberta de agrupamentos sobre dois aspectos: facilita o trabalho do método que busca os grupos, pois flexibiliza a maneira como os atributos descritivos dos dados ou a pertinência do dado aos grupos influencia o processo de agrupamento; fornece um conjunto de informações diferenciado que permite que análises mais refinadas sejam realizadas quando da interpretação dos grupos apresentados como resultado.

Esse diferencial pode ser especialmente útil quando o contexto da aplicação da análise de agrupamento apresenta alguma subjetividade em termos de interpretação de resultados, um fato bastante comum em tarefas de mineração de textos, por exemplo. Considere o contexto de um sistema de recomendação (SR) de notícias baseado em conteúdo (um conteúdo textual). Um SR de notícias simples e hipotético poderia apresentar a seguinte

estratégia para elaborar suas recomendações: dado um conjunto de notícias organizadas em grupos por um método de análise de agrupamento com base na similaridade de seus conteúdos; se um usuário visitar uma notícia, o SR verifica quais são as demais notícias pertencentes ao grupo daquela que foi lida pelo usuário e as recomenda para ele (Figura ??a). Embora essas recomendações pareçam ser ideais, e sob algum aspecto de análise elas são, é factível assumir que esse usuário está recebendo um serviço de recomendação prático, mas talvez menos útil e interessante do que poderia ser e com baixa serendipidade (um resultado de alta serendipidade é aquele que é diferente do que é esperado e comumente praticado, e é mais ou igualmente útil ao contexto).

Uma possibilidade de melhoria nesse sistema hipotético seria usar um algoritmo de análise de agrupamento que permitisse descobrir uma organização de grupo de notícias baseada em **similaridades parciais** ou baseada **em partes** (FRANCA, 2010; HO, 2008). Assim o sistema seria dotado da capacidade de perceber, por exemplo, que algumas notícias podem trazer conteúdo referente a diferentes contextos se forem analisadas apenas sob determinados aspectos. Nesse caso, os grupos formados durante a análise de agrupamento seriam capazes de refletir a diversidade de contextos abordados em uma notícia, fazendo-a pertencer a diferentes grupos, por diferentes motivos. A recomendação, nesse caso, seria potencialmente mais serendípita. Por exemplo, é sabido que eventos de beisebol – o *superbowl* – possuem uma abertura cultural na qual grandes artistas da música fazem apresentações; ou eventos de esportes radicais, como tirolesa, acontecem em eventos de música contemporâneos – *rock in rio*. Tais notícias deveriam aparecer em grupos caracterizados por notícias referentes a esporte, notícias referentes a música, ou notícias referentes a esporte e música (Figura ??(b, c, d)).

Na figura ??(b,c,d) é introduzido graficamente o conceito de coagrupamento. Nesse contexto, o problema de mineração de textos é modelado como o problema de encontrar uma organização dos textos em grupos que considerem similaridades parciais. Assim, um texto pode pertencer a um ou mais cogrupos, a depender dos atributos descritivos sendo considerados. A nomenclatura coagrupamento deriva da estratégia de análise de dados executada durante o processo de descoberta de grupos. Nesse caso, tanto os dados (linhas) quanto os seus atributos (descritivos) são mutuamente submetidos a uma análise de similaridade, e portanto, grupos de dados (linhas) são estabelecidos com respeito a grupos de atributos (colunas).

A associação da análise de coagrupamentos a mineração de textos é interessante por diferentes aspectos. A mineração de textos constitui-se como um problema no qual é preciso lidar com a necessidade de apresentação de resultados com boa interpretabilidade e com um espaço dos dados de alta-dimensionalidade. O primeiro problema é bem resolvido com a estratégia de coagrupamento pois os grupos de atributos que são gerados por ela podem revelar informação antes escondida nos dados (TJHI; CHEN, 2009), e que em um processo de agrupamento tradicional não poderiam ser, pelo menos diretamente, descobertas. Ainda segundo (TJHI; CHEN, 2009), análise de coagrupamento pode apresentar bom desempenho em espaços de alta-dimensionalidade porque seu processo de agrupar atributos (características) pode ser visto com uma redução de dimensionalidade dinâmica para o espaço dos dados.

A despeito da capacidade intrínseca do processo de coagrupamento em lidar de forma diferenciada com o problema de alta-dimensionalidade, ainda se faz necessário notar que no contexto de mineração de dados, ocorre também o problema de esparsidade na representação dos dados. Assim, para implementar a estratégia de coagrupamento com alguma eficiência, é necessário adotar métodos que tenham a capacidade de lidar com esparsidade.

Dentre os diferentes métodos existentes na literatura referentes à implementação de análise de coagrupamento (citar artigos dos vários algoritmos que seguem outras linhas), métodos que usam fatoração de matrizes não negativas (LEE; SEUNG, 2000; LEE; SEUNG, 1999) têm sido vistos como uma boa alternativa a ser aplicada no contexto de mineração de textos (XU; LIU; GONG, 2003; SHAHNAZ et al., 2006; YOO; CHOI, 2010).

## 1.1 Definição do problema

Eu estou entendendo que temos duas facetas do problema. Um deles é o mais obvio que é dar um jeito de descobrir os grupos com sobreposição e mostrar que é útil para recomendação/textos etc. O outro é fazer a fatoração de matriz funcionar pra isso. Então acho que temos que dividir essa definição em duas partes: sobreposição nos cogrupos e fatoração funcionando nisso. Por isso dividi em duas partes, para ver se conseguimos mostrar isso.

### 1.1.1 Estruturas de coagrupamentos

Então aqui entra a parte de mostrar as estruturas de cogrupos possíveis e destacar aquela que fatoração já resolve e depois a que não resolve e a que queremos resolver. Também contextualizar no problema de recomendação ou análise de textos. Figuras precisam entrar aqui para mostra as estruturas.

### 1.1.2 Coagrupamento e fatorização de matrizes

A estratégia de coagrupamento pode ser apresentada como o processo de agrupamento simultâneo de linhas e colunas em uma matriz de dados, de forma que seja possível encontrar **cogrupos** nos quais um **grupo de objetos** (linhas) associado a um deles diz respeito a objetos que são similares entre si considerando um **grupo de atributos** (colunas), também associado ao cogrupo.

Com maior formalidade, dada uma matriz  $X(N, M)$  (prefiro dizer  $X \in \mathbb{R}^{N \times M}$ ) em que  $N$  é o número de linhas,  $M$  o número de colunas e  $x(m, n)$  é, geralmente, um número real representando a relação entre a linha  $x_n$  (linha  $n$ ?) e a coluna  $y_m$  (coluna  $m$ ?), o problema de **coagrupamento** consiste em encontrar um conjunto  $\mathcal{C}$  de submatrizes  $G(I, J)$ , onde  $I = \{i_1, \dots, i_r\}$  com  $r \leq L$  e  $J = \{j_1, \dots, j_s\}$  com  $s \leq C$ , que maximize a similaridade entre os elementos  $g\{i, j\}$ . (Quem é  $L$  e  $C$ ?)

A **esparsidade** em uma matriz é caracterizada pela existência de poucos elementos diferentes de zero (0). Em termos gerais, a esparsidade de uma matriz pode ser medida como a proporção de elementos iguais a zero (0) que ela contém. Problemas de otimização que envolvem matrizes esparsas são caracterizados por apresentarem alta complexidade combinatorial para os quais algoritmos eficientes em matrizes não esparsas tem seu desempenho bastante prejudicado.

**Fatorar uma matriz** consiste em encontrar duas, ou mais, novas matrizes que ao serem multiplicadas, reconstroem a matriz original. Considere uma matriz  $R(N, M)$  em que  $N$  é o número de linhas,  $M$  o número de colunas. A fatoração desta matriz em duas novas matrizes consiste em encontrar duas matrizes  $U(N, K)$  e  $D(M, K)$ , tal que  $R = U \times D^t = \hat{R}$  (Por que colocar  $\hat{R}$ ? Só se for para dizer  $R \approx U \times D^t = \hat{R}$ ). Se  $K$  é escolhido tal que seja menor do que  $N$  e  $M$ , então é dito que  $U$  e  $D$  são representações compactas de  $R$  (Estranho, pois se  $K = M - 1$  parece que precisaremos de mais números

e não haverá nenhuma compactação de dados). Se a matriz  $R$ , e as suas decomposições, são não negativas, tem-se o caso de fatoração de matriz não-negativa.

O problema de coagrupamento pode ser modelado de tal forma que a fatoração de matriz é capaz de fornecer uma aproximação da organização em cogrupos presente no conjunto de dados sob análise.

Considere que o conjunto de dados sob análise é representado pela matriz  $X(N, M)$ , a fatoração dessa matriz em duas (ou mais) novas matrizes  $U(N, K)$  e  $D(M, K)$  significa que  $K$  grupos de linhas foram descobertos, de acordo com  $K$  grupos de colunas.

Se três matrizes são geradas na fatoração,  $U(N, L)$ ,  $S(L, K)$  e  $D(M, K)$ , a interpretação pode incluir uma noção de pesos (matriz  $S$ ) que relacionam grupos de linhas e grupos de colunas, e dimensões diferentes para as matrizes  $U$  e  $D$  podem ser admitidas de modo que o número de grupos de linhas pode ser diferente do número de grupo de colunas.

Imagino que agora tem que entrar uma apresentação rápida do algoritmo novo.

## 1.2 Hipótese

Fatoração de matrizes considerando a decomposição da matriz original em ... como descrever em algo nível aqui?? ... possibilita a descoberta de cogrupos com sobreposição (de colunas); a partir das novas matrizes é possível extrair informação detalhada sobre a relação dos grupos de linhas em relação ao grupo de colunas que pode agregar valor à solução de um problema real de recomendação.

## 1.3 Objetivos

O objetivo geral desse trabalho é o desenvolvimento de novas estratégias de coagrupamento baseadas em fatoração de matrizes, que sejam capazes de descobrir cogrupos com sobreposição em uma dimensão da matriz, isto é, ou sobreposição de colunas ou sobreposição de linhas, considerando uma matriz de valores reais positivos. Com a proposição dessas novas estratégias, este trabalho cobre uma lacuna presente na área de coagrupamento baseado em algoritmos de fatoração de matrizes.

Este trabalho tem como objetivos específicos a aplicação das novas estratégias em um contexto de aplicação real, de forma a ilustrar que elas



Como objetivos específicos, este trabalho mostra que a aplicação das novas estratégias em ambientes controlados (matrizes de dados com cogrupos sintéticas) e em um contexto de aplicação real:

- alcançam resultados tão bons quanto, ou melhores que, as estratégias correlatas que não permitem a sobreposição de dimensões nos cogrupos, quando medidas de qualidade quantitativas são consideradas;
- alcança resultados tão bons quanto, ou melhores que, estratégias de agrupamento quando análise de particionamento clássico e medidas de qualidade quantitativas são consideradas;
- é capaz de melhorar a interpretabilidade qualitativa dos resultados quando comparada aos resultados fornecidos por estratégias de agrupamento clássico.

## 1.4 Metodologia

A análise exploratória da literatura especializada foi escolhida como estratégia para a aquisição de conhecimento sobre a área de coagrupamento e fatoração de Matrizes aplicada à coagrupamento.

E não estou conseguindo encontrar uma forma de descrever a parte referente à concepção das estratégias propostas e também não sei como definir as estratégias referente às derivações.

A fim de permitir a validação das estratégias propostas e, portanto, a verificação da hipótese, fez-se necessário a definição de: (a) um ambiente de teste controlado, representado por uma coleção de conjuntos de dados sintéticos, contendo cada um dos conjuntos situações diferentes referentes à estrutura de coagrupamento e variações em relação à esparsidade; (b) um contexto para realização de uma prova de conceito, no qual um conjunto de dados real foi construído.

Para a prova de conceito foi escolhido usar o conteúdo referente à notícias publicadas no portal iG<sup>1</sup>. Trata-se de um portal de notícias brasileiro muito conhecido, com um volume de notícias bastante grande e com notícias categorizadas em canais, que representam os assuntos dessas notícias. Essas características conferem liberdade para a configuração

---

<sup>1</sup> <http://ig.com.br/>

de experimentos de diferentes naturezas, como experimentos considerando determinadas categorias de notícias, tipos de notícias ou datas de publicação das notícias.

A partir do conteúdo de notícias do portal iG foi construído um corpus de dados textuais, categorizados de acordo com as categorias já usadas no referido portal. Todo o conteúdo do corpus passou por rotinas de pré-processamento comuns na área de Mineração de Texto: *tokenização*, filtragem de *stopwords*, remoção de sufixos (*stemming*), representação da relação “termos  $\times$  documentos” usando estratégias de frequência de termos, como TF-IDF e *n-grams*.

Os resultados da aplicação das estratégias de coagrupamento foram validados utilizando técnicas de avaliação interna, para a verificação da consistência dos biclusters encontrados (SANTAMARÍA; MIGUEL; THERÓN, 2007), e externas (HOCHREITER et al., 2010), avaliando o quanto os biclusters encontrados estão em consenso com as classes de notícias (HOCHREITER et al., 2010).

Então precisaremos dizer aqui como fizemos a avaliação qualitativa.

## 1.5 Organização do documento

Esta dissertação é composta por XXX capítulos incluindo esta introdução. Os demais capítulos estão divididos em duas partes: a primeira é dedicada a explorar a estratégia de coagrupamento implementada com algoritmos baseadas em fatoração de matrizes; a segunda é dedicada a explorar o contexto de sistemas de recomendação baseados em conteúdo textual a partir da aplicação das estratégias de coagrupamento estudadas.

No capítulo 2 são apresentados os principais conceitos referentes à área de coagrupamentos. Especificar mais detalhes ... ..

Estratégias de fatorização de matrizes aplicadas à coagrupamentos são discutidas no capítulo .....

A principal contribuição deste trabalho, as estratégias ... .., é apresentada em detalhes no capítulo ... ..

## 2 Conceitos Fundamentais

Esses conceitos fundamentais ainda vem na qualificação. Não estão ajustados ao que está sendo discutido agora, para a fase final.

Técnicas e algoritmos de Biclusterização são usados, principalmente, no contexto de expressão genética. No entanto, algoritmos de Biclusterização se fazem úteis quando se deseja encontrar *modelos locais*. Ou seja, enquanto algoritmos de clusterização têm o intuito de encontrar *modelos globais*, que geram grupos de dados levando em consideração todas as características, algoritmos de Biclusterização geram grupos de dados em que as características tem alta correlação (FRANCA, 2010; MADEIRA; OLIVEIRA, 2004).

Para a descrição do problema formal de Biclusterização usa-se a seguinte definição (MADEIRA; OLIVEIRA, 2004): seja uma matriz  $A$ , de dimensão  $N \times M$ , um conjunto de linhas  $X = \{x_1, \dots, x_n, \dots, x_N\}$  (aqui deveria ser  $X = \{1, 2, \dots, n, \dots, N\}$ ) e um conjunto de colunas  $Y = \{y_1, \dots, y_m, \dots, y_M\}$  (mesmo comentário anterior), em que  $a_{nm}$  geralmente é um número real e representa a relação entre a linha  $x_n$  e a coluna  $y_m$ ; o problema de Biclusterização é encontrar biclusters, que são submatrizes de  $A$ , denotados por  $A_{IJ}$ , em que  $I \subseteq X$  e  $J \subseteq Y$ . Assim, o bicluster  $A_{IJ}$  é um grupo dos objetos em  $I$ , perante as características com alta correlação  $J$ .

### 2.1 Tipos de cogrupos

Como a definição de bicluster não inclui uma prévia estrutura da matriz  $A$  e dos biclusters  $A_{IJ}$ , diversos algoritmos propostos na literatura diferem quanto ao tipo de bicluster que são capazes de encontrar. Uma taxonomia dos tipos de biclusters é proposta por Madeira e Oliveira (2004):

- *Biclusters com valores constantes*, se trata de biclusters em que todos os valores de  $A_{IJ}$  são constantes:  $a_{ij} = \mu, \forall i, j \in I, J$ , (aqui o valor de  $\mu$  não deveria ser indexado por  $IJ$ , isto é,  $\mu_{IJ}$ ? O mesmo para os outros  $\mu$  que aparecem abaixo.) onde  $\mu$  é um valor constante dentro de  $A_{IJ}$ . Porém, em conjuntos de dados reais, esses biclusters estão presentes com algum tipo de ruído  $\mu + \eta_{ij}$ , onde  $\eta_{ij}$  é o ruído associado com os valores de  $\mu$  e  $a_{ij}$  (MADEIRA; OLIVEIRA, 2004).

- *Biclusters com valores constantes nas linhas ou colunas*, se trata de biclusters com valores constantes nas linhas:  $a_{ij} = \mu + \alpha_i, \forall i, j \in I, J$  ou  $a_{ij} = \mu \cdot \alpha_i, \forall i, j \in I, J$ , onde  $\alpha_i$  é um fator aditivo ou multiplicativo para cada linha; ou ainda biclusters com valores constantes nas colunas:  $a_{ij} = \mu + \beta_j, \forall i, j \in I, J$  ou  $a_{ij} = \mu \cdot \beta_j, \forall i, j \in I, J$ , onde  $\beta_j$  é um fator aditivo ou multiplicativo para cada coluna (MADEIRA; OLIVEIRA, 2004).
- *Biclusters com valores coerentes*, em que são considerados valores próximos entre si (coerentes) para definição de um bicluster:  $a_{ij} = \mu + \alpha_i + \beta_j, \forall i, j \in I, J$ , ou  $a_{ij} = \mu' \cdot \alpha'_i \cdot \beta'_j, \forall i, j \in I, J$ , sendo que se  $\mu = \log \mu' \implies \alpha_i = \alpha'_i, \beta_j = \beta'_j$  (MADEIRA; OLIVEIRA, 2004).
- *Biclusters com evoluções coerentes* têm seus valores com evoluções coerentes, por exemplo, um bicluster com  $a_{i4} \leq a_{i3} \leq a_{i2} \leq a_{i1}$  tem valores com evolução coerente na coluna (MADEIRA; OLIVEIRA, 2004) (estranho considerar a ordem das colunas ou das linhas, já que na maioria dos problemas pode ser bem arbitrário.). Seus valores podem ser gerados por uma função geradora de valores com evolução coerente  $a_{ij} = g(a_{ij}), \forall i, j \in I, J$ , sendo  $g(\cdot)$  não linear e não constante, para que o tipo de bicluster não seja classificado nos casos anteriores. (Muito estranho  $a = g(a)$ , pois  $g()$  deveria ser a identidade)

Os biclusters também diferem quanto as suas estruturas. Cada algoritmo usado para implementar Biclusterização faz uma suposição da estrutura de biclusters que é capaz de encontrar. A Figura ?? sumariza as diferentes estruturas de biclusters, com as linhas e colunas ordenadas para permitir a visualização dos biclusters por meio do mapa de calor dos valores de  $A$ , sendo os biclusters  $A_{IJ}$  representados por cores sólidas e o fundo da matriz ruído.

## 2.2 Algoritmos para Biclusterização

Diversos algoritmos para encontrar biclusters, de diferentes tipos e estruturas, foram propostos na literatura (TANAY; SHARAN; SHAMIR, 2005; MADEIRA; OLIVEIRA, 2004).

Um dos algoritmos de Biclusterização mais comum e simples. que encontra biclusters com valores coerentes, em estrutura com sobreposição e arbitrariamente posicionados, é o *Coupled Two-way Clustering* (CTWC) (GETZ; LEVINE; DOMANY, 2000). O algoritmo

CTWC é capaz de encontrar biclusters através da clusterização de objetos e atributos (linhas e colunas), separadamente. O algoritmo de clusterização usado por [Getz, Levine e Domany \(2000\)](#) foi o *Superparamagnetic Clustering* (SPC), o qual é capaz de determinar o número de clusters automaticamente, e com uma estratégia de clusterização hierárquica *top-down* é capaz de gerar clusters estáveis ([GETZ; LEVINE; DOMANY, 2000](#)). O SPC tem como entrada uma matriz de similaridade e um parâmetro temperatura, que controla o quão estáveis serão os clusters que o algoritmo gerará. Assim, o CTWC encontra clusters estáveis de linhas e colunas através do SPC, e iterativamente executa o SPC nos clusters de linhas e colunas encontrados, mantendo na memória um par do subconjunto de linhas e do subconjunto de colunas (biclusters), assim como os clusters estáveis de linhas e colunas, separadamente.

Já o algoritmo de [Cheng e Church \(2000\)](#) é capaz de encontrar o mesmo tipo de bicluster que o algoritmo CTWC, porém usando uma estratégia gulosa: biclusters com valores coerentes e estrutura com sobreposição e arbitrariamente posicionados. Este algoritmo está sendo objeto de estudo desse projeto de mestrado para aplicação em dados textuais e por isso segue aqui descrito em mais detalhes. Nesse algoritmo, para encontrar biclusters, ou  $\delta$ -biclusters, na matriz  $A$ , os autores definem o *Resíduo Quadrático Médio* (RQM):

$$\begin{aligned} H_{IJ} &= \frac{1}{|I||J|} \sum_{i,j \in I,J} (a_{ij} - a_{iJ} - a_{IJ} + a_{IJ})^2 \\ H_{iJ} &= \frac{1}{|J|} \sum_{j \in J} (a_{ij} - a_{iJ} - a_{IJ} + a_{IJ})^2 \\ H_{IJ} &= \frac{1}{|I|} \sum_{i \in I} (a_{ij} - a_{iJ} - a_{IJ} + a_{IJ})^2 \end{aligned}$$

em que

$$a_{iJ} = \frac{1}{|J|} \sum_{j \in J} a_{ij}, \quad a_{IJ} = \frac{1}{|I|} \sum_{i \in I} a_{ij}, \quad a_{IJ} = \frac{1}{|I||J|} \sum_{i,j \in I,J} a_{ij}$$

onde  $H_{IJ}$  é o RQM de uma submatriz  $A_{IJ}$ ,  $H_{iJ}$  o RQM da linha  $i$ ,  $H_{IJ}$  o RQM da coluna  $j$ ,  $a_{iJ}$  a média dos valores da linha  $i$ ,  $a_{IJ}$  a média dos valores da coluna  $j$  e  $a_{IJ}$  a média dos valores da submatriz  $A_{IJ}$ , definida pelos subconjuntos  $I$  e  $J$ .

Então, um bicluster perfeito  $A_{IJ}$  teria o RQM  $H_{IJ} = 0$ , pois  $a_{ij} = a_{ij}, \forall i, j \in I, J$ , ( $a_{ij} = a_{ij}$ , hein? não seria  $a_{ij} = a_{iJ} + a_{IJ} - a_{IJ}$ ?) fazendo  $a_{iJ} = a_{IJ} = a_{IJ}$ . No entanto, se apenas minizar o RQM, um bicluster com apenas um elemento seria perfeito, o que pode

não refletir a realidade. Além disso, em conjunto de dados reais existe ruído, podendo esconder o bicluster perfeito.

Para encontrar biclusters, ou  $\delta$ -biclusters, [Cheng e Church \(2000\)](#) usam uma estratégia gulosa que retira linhas e colunas, visando a minimização do RQM, respeitando um parâmetro  $\delta$ , que é calibrado pelo usuário. Então, um bicluster é encontrado quando o RQM de uma submatriz  $A_{IJ}$  é  $H_{IJ} \leq \delta$ , para algum  $\delta \geq 0$ . As etapas de remoções de elementos da matriz são apresentadas nos algoritmos 1 e 2.

O algoritmo 2 é usado para acelerar o processo de busca de um  $\delta$ -bicluster, convergindo mais rapidamente para uma solução quanto maior for o parâmetro  $\alpha$ , em que  $\alpha \geq 0$ . Ainda, para amenização do problema de encontrar  $\delta$ -biclusters perfeitos com apenas um elemento, ou poucos elementos, é utilizado o algoritmo 3, que adiciona nós sem aumentar o RQM do bicluster.

Por fim, o algoritmo 4 é a consolidação dos algoritmos 3, 2 e 1 e a iteração para encontrar  $k$   $\delta$ -biclusters, um a um, sendo  $k$  fornecido pelo usuário.

Além dos algoritmos apresentados, existem outros algoritmos que são capazes de encontrar outros tipos de biclusters (Seção 2.1), além de serem recentes ([FRANÇA; ZUBEN, 2010](#); [YANG; LESKOVEC, 2013](#); [HOCHREITER et al., 2010](#); [CABANES; BENNANI; FRESNEAU, 2012](#)), mostrando que ainda há interesse na área de pesquisa de Biclusterização.

## 2.3 Avaliação de Biclusterização

Para determinar parâmetros, descobrir a qualidade e/ou estabilidade dos biclusters encontrados por algoritmos, é necessário estabelecer métricas de avaliação. Existem duas maneiras de avaliar biclusters ([HOCHREITER et al., 2010](#)): *interna*, usa os dados dos resultados dos algoritmos, juntamente com métricas de qualidade e/ou estabilidade, para avaliar as soluções geradas; *externa*, utiliza os dados reais das soluções de biclusters de um conjunto de dados, usando estratégias para comparação, obtendo assim, maior confiança nas soluções.

A avaliação interna pode não ser tão precisa quanto a avaliação externa, porém é útil para descobrir parâmetros ótimos. Apesar de [Prelić et al. \(2006\)](#) sugerirem não usar avaliações internas, por não estar claro como estender noções de separação e homogeneidade, [Santamaría, Miguel e Therón \(2007\)](#) descreveu métricas de consistência para verificar se

um bicluster é consistente com a sua definição, seja aditiva, multiplicativa e/ou constante, fazendo uma comparação dos elementos do bicluster:

$$C_l(A_{IJ}) = \frac{1}{|I|} \sum_{i=1}^{|I|-1} \sum_{j=i+1}^{|I|} \sqrt{\sum_{k=1}^{|J|} (a_{ik} - a_{jk})^2}$$

$$C_c(A_{IJ}) = \frac{1}{|J|} \sum_{i=1}^{|J|-1} \sum_{j=i+1}^{|J|} \sqrt{\sum_{k=1}^{|I|} (a_{ki} - a_{kj})^2}$$

em que  $C_l(A_{IJ})$  é o índice de consistência das linhas do bicluster  $A_{IJ}$  e  $C_c(A_{IJ})$  é o índice de consistência das colunas do bicluster  $A_{IJ}$ . Ainda, a consistência do bicluster inteiro  $C$  pode ser definida pela média:

$$C(A_{IJ}) = \frac{|I| \cdot C_l + |J| \cdot C_c}{|I| + |J|}$$

Uma das métricas externas que são usadas para comparar biclusters encontrados com biclusters reais em um conjunto de dados, é a métrica *concensus score* ([HOCHREITER et al., 2010](#)). Essa métrica calcula a maximização das similaridades entre biclusters encontrados e reais, usando o *índice de Jaccard* como medida de similaridade e o algoritmo Húngaro para solucionar o problema de maximização. A saída da avaliação é um *score*  $\in [0, 1]$ , em que 0 significa que os biclusters comparados são totalmente diferentes, e 1 o inverso.

### 3 Fatoração de matrizes não-negativas para coagrupamento

Fatoração de matrizes não-negativas (*Non-negative Matrix Factorization* - NMF) foi estudada como um método para análise de dados capaz de extrair conhecimento sobre um objeto a partir do estudo de suas partes (LEE; SEUNG, 1999), como um contraponto a métodos mais populares como Análise de Componentes Principais (*Principal Component Analysis* - PCA) e Quantização Vetorial, porém, considerando a fatoração matrizes positivas ou negativas. Lee e Seung (1999) apresentam tal abordagem a partir de sua aplicação no aprendizado de características de faces (em dados do tipo imagem) e na análise de características semânticas de textos. A análise de textos também foi usada como aplicação na ilustração da aplicação de fatorização de matrizes por Ho (2008) que segue a ideia de aprendizado de partes, por Kuang (2014) que aplica fatoração de matrizes não-negativas sobre um problema formulado como análise de agrupamento (clustering), e por Long, Zhang e Yu (2005), Ding et al. (2006), Yoo e Choi (2010), que formulam o problema como coagrupamento (cogruping). Ainda, outros contextos são submetidos à análise sob a formulação de problemas de coagrupamento, sendo alguns exemplos a clusterização de genes e análise de microarray em bioinformática (KLUGER et al., 2003) e a filtragem colaborativa em sistemas de recomendação (SALAKHUTDINOV; MNIH, 2008). Na aplicação de NMF em filtragem colaborativa em sistemas de recomendação, destaca-se o modelo baseado em NMF de Koren (2009) para predição de preferências de usuários por filmes, que ganhou em primeiro lugar o Prêmio Netflix (*Netflix Prize*)<sup>1</sup>.

A adequação da fatoração de matrizes para tarefas modeladas como agrupamento ou coagrupamento ocorre porque muitas das representações usadas em aplicações nessas áreas se apresentam como uma relação entre pares de elementos pertencentes a conjuntos finitos, como apresentado em (LONG; ZHANG; YU, 2005). Por exemplo, na resolução da tarefa de agrupamento de documentos, em mineração de textos, usa-se, comumente, dois conjuntos finitos, documentos e palavras, sendo que a relação entre eles é representada pela ocorrência de uma determinada palavra em um determinado documento. Note ainda que a relação expressa entre os elementos, como no contexto de palavras e documentos, apresenta-se como uma matriz de dados positiva, característica que ilustra a aplicabilidade de NMF.

---

<sup>1</sup> <http://www.netflixprize.com/>



Formalmente, algoritmos de coagrupamento baseados em NMF têm como entrada uma matriz de dados  $X \in \mathbb{R}_+^{n \times m}$ , contendo números reais positivos com  $n$  linhas e  $m$  colunas. Esta matriz é formada por um conjunto de vetores de linhas  $\mathcal{N} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  e um conjunto de vetores de colunas  $\mathcal{M} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , e as relações existentes entre cada linha  $x$  e cada coluna  $y$  são representadas por  $x_{ij}$  considerando os índices  $i = \{1, \dots, n\}$  e  $j = \{1, \dots, m\}$ , que é justamente um valor da matriz  $X$ . Cada valor em  $x_{ij}$  representa, então, a relação existente entre pares de elementos em algum contexto de interesse. O objetivo é encontrar  $k$  partições de  $\mathcal{N}$ , denotadas pelos subconjuntos ordenados  $\mathcal{K}_p \subseteq \mathcal{N}$ ,  $l$  partições para  $\mathcal{M}$ , denotadas pelos subconjuntos ordenados  $\mathcal{L}_q \subseteq \mathcal{M}$ , considerando os índices  $p = \{1, \dots, k\}$  e  $q = \{1, \dots, l\}$ . Então, os subconjuntos  $\{\mathcal{K}_1, \dots, \mathcal{K}_k\}$  e  $\{\mathcal{L}_1, \dots, \mathcal{L}_l\}$  são os cogrupos de linhas e colunas, respectivamente.

Para implementação da NMF como uma estratégia para resolução do problema de coagrupamento, diferentes algoritmos foram apresentados na literatura. Cada uma delas considera o problema de NMF com diferentes restrições que permitem propor soluções para o problema de coagrupamento de diferentes naturezas. Este capítulo se destina a apresentar três das implementações existentes que são usadas como base para a proposta desta dissertação: decomposição de valores em blocos (Seção 3.1 introduzida por Long, Zhang e Yu (2005)); fatoração ortogonal tripla de matrizes não-negativas (Seção 3.2) introduzida por Ding et al. (2006); e fatoração tripla rápida de matrizes não-negativas (Seção 3.3) introduzida por Wang et al. (2011). Outras implementações correlatas podem ser encontradas em (LI; DING, 2006).

Todos algoritmos de NMF para resolução das tarefas de agrupamento e coagrupamento tem em comum a natureza recursiva, pois são encontradas partições de linhas a partir de partições de colunas, para então, encontrar melhores partições de colunas a partir das partições de linhas. Esse processo recursivo continua até que a aproximação entre a fatoração e a matriz fatorada ( $X$ ) através de alguma medida atinja um mínimo local ou global, resultando em uma solução para o particionamento de linhas e colunas.

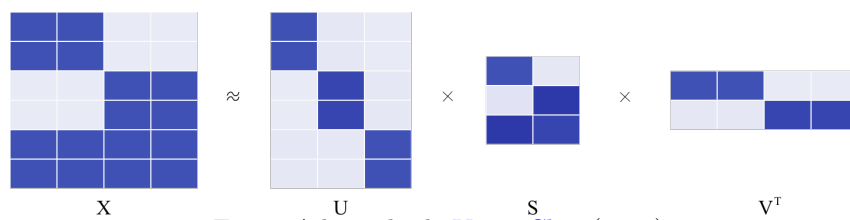
### 3.1 Decomposição de Valores em Blocos para Coagrupamento

A Decomposição de Valores em Blocos (*Block Value Decomposition* - BVD) foi proposto por Long, Zhang e Yu (2005) como uma abordagem para tratar o problema de coagrupamento, com base em fatoração de matrizes não-negativas. Esta decomposição

recebe esse nome por ter a capacidade de encontrar estruturas em blocos escondidas na matriz de dados. Isso é possível porque o algoritmo BVD é capaz de explorar a relação entre linhas e colunas da matriz de dados por meio da decomposição dela em três matrizes:  $U$  uma matriz de coeficientes de linhas,  $S$  uma matriz com estrutura em blocos, e  $V$  uma matriz de coeficientes de colunas. Segundo os autores, tais coeficientes em  $U$  e  $V$  podem ser vistos como um fator que associa linhas às partições encontradas no conjunto de linhas, e que associa as colunas às partições encontradas no conjunto de colunas, respectivamente; e  $S$  pode ser vista como uma representação compacta da matriz original de entrada e permite sua reconstrução aproximada a partir da operação  $USV^T$ . Sob o ponto de vista de resolução do problema de coagrupamento, então, o objetivo no BVD é encontrar grupos de linhas e colunas de forma simultânea, sendo  $k$  grupos de  $\mathcal{N}$  (linhas) e  $l$  grupos de  $\mathcal{M}$  (colunas).

Ainda, os autores proponentes da abordagem BVD defendem que interpretações intuitivas podem ser derivadas da análise das combinações das matrizes geradas na fatoração, quando aplicadas a um contexto específico. Um exemplo fornecido no trabalho de [Long, Zhang e Yu \(2005\)](#) é que, considerando uma matriz de entrada que representa a relação “documentos por palavras” (linhas por colunas) cada coluna de  $US$  captura a ideia de uma base para a representação de grupos de palavras; e cada linha em  $SV^T$  captura a ideia de uma base para a representação de grupos de documentos. Uma representação gráfica para o resultado de uma fatoração de matrizes pode ser visto na figura 1.

Figura 1 – Fatoração da matriz original de dados  $X$  em três outras matrizes:  $U$ ,  $S$  e  $V$



Fonte: Adaptado de [Yoo e Choi \(2010\)](#)

Na figura 1 considere que uma célula com cor escura representa a existência de uma relação entre linha e coluna, e uma célula com cor clara representa a inexistência de uma relação entre linha e coluna, e que essas relações são estabelecidas adequadamente em cada contexto de aplicação. Transportando o exemplo da figura para o contexto de uma matriz de “documentos por palavras” tem-se um conjunto de seis documentos e quatro palavras, sendo que, por exemplo, o primeiro documento possui uma relação com as duas

primeiras palavras, e não possui relação com as terceira e a quarta palavras. A matriz  $U$  pode ser interpretada como uma matriz “documentos por grupos de documentos”, sendo portanto uma situação em que seis documentos estão agrupados em três grupos ( $k = 3$ ): os dois primeiros documentos no primeiro grupo, os dois próximos no segundo grupo e os dois últimos no terceiro grupo. A matriz  $V^T$  pode ser interpretada como uma matriz de “grupos de palavras por palavras”, sendo portanto uma situação em que há dois grupos de palavras ( $l = 2$ ) no contexto das quatro palavras existentes. Finalmente, a matriz  $S$  representa uma relação entre “grupos de documentos” e “grupos de palavras”. A primeira linha da matriz  $S$  indica que há uma relação entre o primeiro grupo de linhas e o primeiro grupos de palavras, e que não há uma relação entre o primeiro grupo de linhas e o segundo grupo de palavras. Seguindo a interpretação intuitiva apresentada por (LONG; ZHANG; YU, 2005), uma das combinações possíveis,  $SV^T$ , está ilustrada na figura 2. Observe que a matriz  $SV^T$  representa “três grupos de documentos por quatro palavras”, constituindo-se como uma base de representação para grupos de documentos.

Figura 2 – A reconstrução da primeira linha  $\mathbf{x}_1$ . de  $X$ , através da multiplicação da matriz indicadora de grupos de linhas  $U$  pela matriz dos protótipos de linhas ( $SV^T$ ).

The diagram shows the reconstruction of the first row of matrix  $X$  through matrix multiplication. It is structured as follows:

- Top Row:** Matrix  $X$  (a 4x4 grid) is approximately equal to matrix  $U$  (a 4x4 grid) multiplied by the product of matrix  $S$  (a 2x2 grid) and matrix  $V^T$  (a 2x4 grid). The entire product  $S \times V^T$  is enclosed in large parentheses.
- Middle Row:** Matrix  $X$  is equal to matrix  $U$  multiplied by matrix  $SV^T$  (a 2x4 grid).
- Bottom Row:** Matrix  $X$  is approximately equal to matrix  $U(SV^T)$  (a 4x4 grid).

The matrices are represented as grids of blue and light blue squares, indicating binary or weighted values. The labels  $X$ ,  $U$ ,  $S$ ,  $V^T$ , and  $U(SV^T)$  are placed below their respective grids.

Fonte: Lucas Fernandes Brunialti, 2016

Note ainda que nessa estrutura de matrizes fatoradas é possível identificar protótipos responsáveis por cada parte da reconstrução da matriz original  $X$ , visto que cada base de

representação pode ser vista como um conjunto de vetores protótipos: as linhas de  $(SV^T)$  são vetores base (protótipos de grupos de linhas). O raciocínio apresentado para interpretação da figura 2 pode também ser feito para a interpretação da combinação  $US$ . Importante salientar que, assim como ocorre na resolução da tarefa de agrupamento, diferentes grupos de linhas e de colunas podem ser obtidos, representando diferentes soluções para o problema. E, neste caso ainda, diferentes matrizes  $S$  podem ser obtidas para uma mesma organização de grupos de linhas e colunas. Portanto, qualquer interpretação derivada dessa análise deve ser considerada como apenas uma das formas possíveis de análise dos dados provenientes do contexto de aplicação.<sup>2</sup>

Semelhante ao *fuzzy k-means*, é possível observar esta fatoração como uma ótica de compactação. O BVD compacta a matriz de dados em uma matriz com fatores que correlacionam cada linha com cada grupo de linhas ( $U$ ). Além disso, esta compactação adiciona a idéia de uma matriz de fatores  $V$  para compactar as colunas de  $X$ , e  $S$  que é uma visão compactada de  $X$  em  $kl$  elementos. Portanto, a compactação transforma  $nm$  elementos em  $nk + kl + ml$  elementos, através das matrizes  $U$ ,  $S$  e  $V$ .

O problema de coagrupamento implementado sob a abordagem BVD é formalmente apresentado como (LONG; ZHANG; YU, 2005):

**Problema 1** (Problema de Decomposição de Valores em Blocos).

$$\begin{aligned} \mathcal{F}_1(U, S, V) = \min_{U, S, V} \quad & \|X - USV^T\|_F^2 \\ \text{sujeito a} \quad & U \geq 0, \\ & S \geq 0, \\ & V \geq 0 \end{aligned}$$

em que  $U \in \mathbb{R}_+^{n \times k}$ ,  $S \in \mathbb{R}_+^{k \times l}$ ,  $V \in \mathbb{R}^{m \times l}$ , e  $\|\cdot\|_F$  denota a norma de Frobenius para matrizes.

A implementação para o processo de minimização do problema 1 é descrito no algoritmo 1, o qual é baseado em atualizações multiplicativas e tem sua convergência demonstrada via teoria de otimização não linear (veja (LONG; ZHANG; YU, 2005)). Nesse algoritmo considere  $t$  o contador de iterações,  $U^{(t)}$ ,  $S^{(t)}$  e  $V^{(t)}$ , as matrizes  $U$ ,  $S$  e  $V$ , na iteração  $t$ , respectivamente, e  $\odot$  é o produto de Hadamard.

<sup>2</sup> Medidas de avaliação de agrupamento ou coagrupamento internas podem ser aplicadas às diferentes soluções encontradas de maneira a guiar um processo decisório no contexto de aplicação.

---

**Algoritmo 1** Algoritmo baseado em atualização multiplicativa para solução do BVD
 

---

```

1: function BVD( $X, t_{max}, k, l$ )
2:   Inicialize:  $U^{(0)} \leftarrow \mathcal{U}(0, 1), V^{(0)} \leftarrow \mathcal{U}(0, 1), S^{(0)} \leftarrow \frac{1}{nm} \sum_{i,j} x_{ij}$  e  $t \leftarrow 0$ .
3:   while (não convergiu) ou ( $t \leq t_{max}$ ) do
4:
5:     
$$U^{(t+1)} \leftarrow U^{(t)} \odot \frac{XV^{(t)}S^{(t)T}}{U^{(t)}S^{(t)}V^{(t)T}V^{(t)}S^{(t)T}}$$

6:
7:     
$$V^{(t+1)} \leftarrow V^{(t)} \odot \frac{X^T U^{(t+1)} S^{(t)}}{V^{(t)} S^{(t)T} U^{(t+1)T} U^{(t+1)} S^{(t)}}$$

8:
9:     
$$S^{(t+1)} \leftarrow S^{(t)} \odot \frac{U^{(t+1)T} X V^{(t+1)}}{U^{(t+1)T} U^{(t+1)} S^{(t)} V^{(t+1)T} V^{(t+1)}}$$

10:     $t \leftarrow t + 1$ 
11:  end while
12:  return  $U^{(t)}, S^{(t)}, V^{(t)}$ 
13: end function

```

---

A inicialização dos elementos das matrizes  $U$ ,  $S$  e  $V$  são gerados através de uma distribuição uniforme que ignora zeros ( $\mathcal{U}(0, 1) \in ]0, 1]$ ). Como condições para assumir a convergência, neste trabalho, considera-se a diferença do erro de aproximação em duas iterações consecutivas menor ou igual a um  $\epsilon$ :

$$\left\| X - U^{(t)} S^{(t)} V^{(t)T} \right\|_F^2 - \left\| X - U^{(t+1)} S^{(t+1)} V^{(t+1)T} \right\|_F^2 \leq \epsilon$$

O algoritmo também pára caso a  $t$ -ésima iteração seja igual ao número máximo de iterações ( $t_{max}$ ).

Note que como  $U$  e  $V$  possuem valores no domínio dos reais, então, não é possível obter as partições diretamente, sem um processo de pós-processamento. Um modo simples de obter o particionamento para as linhas, é o seguinte:

$$\mathcal{K}_p = \mathcal{K}_p + \{x_i\} \mid p = \arg \max_{p'} \mathbf{u}_i. \quad \forall i = \{1, \dots, n\}, \forall p, p' = \{1, \dots, k\}$$

Isso significa que uma linha  $i$  pertencerá a um cogruppo  $p$  (ou partição) se para todos os  $k$  cogrupos, o fator  $u_{ip}$  for maior que todos os outros fatores para os outros cogrupos, presentes no vetor  $\mathbf{u}_i$ .

A complexidade de tempo do algoritmo é possível ser calculada fixando as condições  $k \simeq l$ ,  $n \simeq m$ ,  $k \ll n$ ,  $l \ll m$  e usando um algoritmo para otimizar a ordem das multiplicações (Cormen et al. (2001) discute algoritmos para tal otimização):  $\mathcal{O}\left(t_{max}(nl(m+k) + ml(n+k) + mk(n+l) + k^2(n+l) + l^2(m+k))\right)$ .

### 3.2 Fatoração Ortogonal Tripla de Matrizes Não-negativas

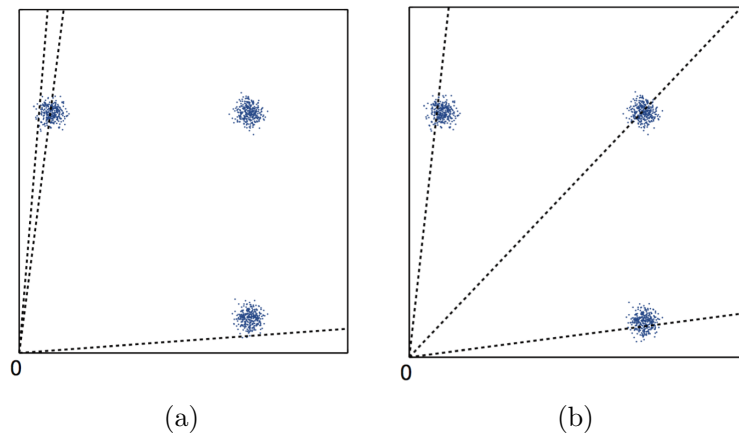
Baseado no problema 1, [Ding et al. \(2006\)](#) propõem o problema 2, e o chama de Fatoração Ortogonal Tripla de Matrizes Não-negativas (*Orthogonal Non-negative Matrix Tri-factorization* - ONMTF).

**Problema 2** (Problema de Fatoração Ortogonal tripla de Matrizes Não-negativas).

$$\begin{aligned} \mathcal{F}_2(U, S, V) &= \min_{U, S, V} \|X - USV^T\|_F^2 \\ \text{sujeito a } &U \geq 0, S \geq 0, V \geq 0, \\ &U^T U = I, \\ &V^T V = I \end{aligned}$$

em que  $U \in \mathbb{R}_+^{n \times k}$ ,  $S \in \mathbb{R}_+^{k \times l}$ ,  $V \in \mathbb{R}^{m \times l}$  e  $\|\cdot\|_F$  denota a norma de Frobenius. Na formulação desse problemas, duas restrições de ortonormalidade são acrescentadas,  $U^T U = I$  e  $V^T V = I$ , em que  $I$  é a matriz identidade, para as matrizes de grupos de linhas e grupos de colunas, respectivamente. Tais restrições restringem o problema da fatoração  $X \approx USV^T$  para um número menor de possíveis soluções, buscando a unicidade, como mostrado na figura 3. No entanto, o algoritmo não garante que as restrições de ortonormalidade são respeitadas, apesar das restrições colocadas encontrarem soluções de particionamento com maior ortogonalidade.

Figura 3 – Base de protótipos obtidas com FMN sem restrições (a) e com restrições de ortogonalidade nas matrizes



Fonte: [Yoo e Choi \(2010\)](#)

A figura 3 representa um problema com três grupos (três nuvens de pontos). As linhas pontilhadas representam os protótipos obtidos por FMN (a) sem restrição de

ortogonalidade nas matrizes, como o BVD (a), e (b) com restrição de ortogonalidade nas matrizes, como o ONMTF. Note que a base obtida com FMN ortogonal tende a encontrar protótipos mais centralizados nos grupos. Já a base obtida com FMN sem restrições tende a encontrar uma região convexa que contém os pontos dos grupos (incluindo regiões que abrangem pontos de grupos originalmente projetados como sendo diferentes). A solução obtida em (a), embora correta, não é desejável.

Em termos de compactação, igualmente ao BVD, o algoritmo para solução do ONMTF transforma os  $nm$  elementos da matriz  $X$  em  $nk + kl + ml$  elementos, através da fatoração em  $U$ ,  $S$  e  $V$ .

Ding et al. (2006) propõem uma solução para implementação do processo de minimização para o problema 2 semelhante ao que foi apresentado na seção 3.1, também baseado em atualizações multiplicativas e com convergência demonstrada com base na teoria de otimização não linear. O algoritmo para tal processo de minimização é apresentado no algoritmo 2, no qual  $t$  o contador de iterações,  $U^{(t)}$ ,  $S^{(t)}$  e  $V^{(t)}$ , as matrizes  $U$ ,  $S$  e  $V$ , na iteração  $t$ , respectivamente,  $\odot$  é o produto de Hadamard e  $\mathcal{U}(0, 1) \in ]0, 1]$  uma função que gera valores de uma distribuição uniforme que ignora zeros. Também, a mesma condição de convergência aplicada no algoritmo 1 pode ser aplicada nesse caso.

---

**Algoritmo 2** Algoritmo baseado em atualização multiplicativa para solução do ONMTF

---

```

1: function ONM3F( $X, t_{max}, k, l$ )
2:   Initialize:  $U^{(0)} \leftarrow \mathcal{U}(0, 1), V^{(0)} \leftarrow \mathcal{U}(0, 1), S^{(0)} \leftarrow \mathcal{U}(0, 1)$  e  $t \leftarrow 0$ .
3:   while (não convergiu) ou ( $t \leq t_{max}$ ) do
4:

```

$$U^{(t+1)} \leftarrow U^{(t)} \odot \sqrt{\frac{XV^{(t)}S^{(t)T}}{U^{(t)}U^{(t)T}XV^{(t)}S^{(t)T}}} \quad (1)$$

$$V^{(t+1)} \leftarrow V^{(t)} \odot \sqrt{\frac{X^T U^{(t+1)} S}{V^{(t)} V^{(t)T} X^T U^{(t+1)} S}} \quad (2)$$

$$S^{(t+1)} \leftarrow S^{(t)} \odot \sqrt{\frac{U^{(t+1)T} X V^{(t+1)}}{U^{(t+1)T} U^{(t+1)} S^{(t)} V^{(t+1)T} V^{(t+1)}}} \quad (3)$$

```

7:      $t \leftarrow t + 1$ 
8:   end while
9:   return  $U^{(t)}, S^{(t)}, V^{(t)}$ 
10: end function

```

---

Note que é possível obter o particionamento de linhas e colunas da mesma forma que foi descrita com o BVD. Calculando a complexidade de tempo, fixando as mesmas condições  $k \simeq l$ ,  $n \simeq m$ ,  $k \ll n$ ,  $l \ll m$  e usando um algoritmo para otimizar a ordem

das multiplicações, é possível encontrar a mesma complexidade antes calculada para o algoritmo 1.

No artigo de Yoo e Choi (2010), é proposta uma abordagem mais simples para a derivação das regras de atualização multiplicativas, considere uma função de otimização qualquer  $\mathcal{J}$  e seu respectivo gradiente  $\nabla \mathcal{J}$ :

$$\nabla \mathcal{J} = [\nabla \mathcal{J}]^+ - [\nabla \mathcal{J}]^-$$

onde  $[\nabla \mathcal{J}]^+$  é a parte positiva do gradiente,  $[\nabla \mathcal{J}]^-$  a parte negativa do gradiente. Se  $[\nabla \mathcal{J}]^+ \geq 0$  e  $[\nabla \mathcal{J}]^- \geq 0$ , então, é possível definir uma regra de atualização multiplicativa, para otimizar os parâmetros  $\Theta$  da função  $\mathcal{J}$ :

$$\Theta \leftarrow \Theta \odot \left( \frac{[\nabla \mathcal{J}]^-}{[\nabla \mathcal{J}]^+} \right)^\eta \quad (4)$$

onde  $\odot$  representa o produto Hadamard,  $(\cdot)^\eta$  representa a potência para cada elemento, e  $\eta$  uma taxa de aprendizado ( $0 < \eta \leq 1$ ). Então, se  $\Theta$  for inicializado com elementos positivos, é possível verificar que a regra de atualização multiplicativa da equação 4 mantém a não-negatividade de  $\Theta$ .

Também, é utilizada uma abordagem diferente para a derivação de regras de atualização multiplicativas, visando um algoritmo para a solução do problema 2. Neste caso, o gradiente é calculado com base em uma superfície com restrições que preserva a ortogonalidade. Essa superfície com restrições é chamada de Variedade de Stiefel (*Stiefel Manifold*).

Assim, Yoo e Choi (2010) faz o uso da estratégia da equação 4 e da teoria de derivação na superfície com restrições (Variedade Stiefel), para propor uma solução para o problema 2 (ONMTF) alternativa às atualizações das equações 1, 2 e 3, através da atualização multiplicativa. Essas atualizações são apresentadas no algoritmo 3, considere  $t$  o contador de iterações,  $U^{(t)}$ ,  $S^{(t)}$  e  $V^{(t)}$ , as matrizes  $U$ ,  $S$  e  $V$ , na iteração  $t$ , respectivamente,  $\odot$  o produto de Hadamard,  $\mathcal{U}(0, 1) \in ]0, 1]$  uma função que gera valores de uma distribuição uniforme que ignora zeros,  $\text{diag}(\cdot)$  uma função que extrai a diagonal principal de uma matriz e a transforma em um vetor, e  $\mathbf{1}$  um vetor de uns com dimensão que torne a multiplicação possível.

Considerando a preservação da ortonormalidade, o algoritmo apresentado é capaz de preservar melhor a ortonormalidade. Para tal afirmação, Yoo e Choi (2010) realizou



um experimento que consistiu em medir as restrições de ortonormalidade em  $U$  e  $V$ , através das restrições  $\|U^T U - I\|$  e  $\|V^T V - I\|$ , respectivamente, durante cada iteração do algoritmo proposto, comparando-o com o algoritmo proposto por [Ding et al. \(2006\)](#).

---

**Algoritmo 3** Algoritmo baseado em atualização multiplicativa e na teoria de de derivação na superfície com restrições (Variedade Stiefel) para solução do ONMTF

---

```

1: function ONMTF( $X, t_{max}, k, l$ )
2:   Initialize:  $U^{(0)} \leftarrow \mathcal{U}(0, 1), V^{(0)} \leftarrow \mathcal{U}(0, 1), S^{(0)} \leftarrow \mathcal{U}(0, 1)$  e  $t \leftarrow 0$ .
3:   while (não convergiu) ou ( $t \leq t_{max}$ ) do
4:
5:     
$$U^{(t+1)} \leftarrow U^{(t)} \odot \frac{X V^{(t)} S^{(t)T}}{U^{(t)} S^{(t)} V^{(t)T} X^T U^{(t)}}$$

6:
7:     
$$V^{(t+1)} \leftarrow U^{(t)} \odot \frac{X^T U^{(t+1)} S^{(t)}}{V^{(t)} S^{(t)T} U^{(t+1)T} X V^{(t)}}$$

8:
9:     
$$S^{(t+1)} \leftarrow S^{(t)} \odot \frac{U^{(t+1)T} X V^{(t+1)}}{U^{(t+1)T} U^{(t+1)} S^{(t)} V^{(t+1)T} V^{(t+1)}}$$

10:
11:     $t \leftarrow t + 1$ 
12:  end while
13:
14:    
$$U^{(t)} \leftarrow U^{(t)} \text{diag}(S^{(t)} \text{diag}(\mathbf{1}^T V^{(t)}) \mathbf{1})$$

15:
16:    
$$V^{(t)} \leftarrow V^{(t)} \text{diag}(\mathbf{1}^T \text{diag}(\mathbf{1}^T U^{(t)}) S^{(t)})$$

17:
18:  return  $U^{(t)}, S^{(t)}, V^{(t)}$ 
19: end function

```

---

Ainda, [Yoo e Choi \(2010\)](#) propõem uma normalização baseada numa interpretação probabilística da fatoração de  $X$  em  $USV^T$ , para então realizar o particionamento de linhas e colunas, como apresentado no algoritmo 3. Note que o particionamento de linhas e colunas é realizado como nos outros algoritmos já apresentados e a mesma condição de convergência aplicada no algoritmo 1 e 2 pode ser aplicada nesse caso.

Calculando a complexidade de tempo, seguindo as mesmas restrições apresentadas anteriormente para os outros algoritmos, é possível verificar que o algoritmo 3, proposto por [Yoo e Choi \(2010\)](#) para solução do problema 2 (ONMTF), é equivalente à complexidade apresentada para o algoritmo 1 e 2.

### 3.3 Fatoração Tripla Rápida de Matrizes Não-negativas

O problema de Fatoração Tripla Rápida de Matrizes Não-negativas (*Fast Non-negative Matrix Tri Factorization* - FNMTF), formalizado no problema 3, foi proposto por Wang et al. (2011) com os seguintes argumentos contra o uso prático dos problemas até então propostos para encontrar cogrupos: eles exigem soluções algorítmicas iterativas, com intensas multiplicações de matrizes em cada passo do algoritmo; eles propõem encontrar coagrupamentos flexíveis (com restrições relaxadas), o que implica em encontrar inúmeras soluções para a tarefa de coagrupamento.

**Problema 3** (Fatoração tripla rápida de Matrizes Não-negativas).

$$\begin{aligned} \mathcal{F}_3(U, S, V) = \min_{U, S, V} & \|X - USV^T\|_F^2 \\ & U \in \Psi^{n \times k}, \\ & V \in \Psi^{m \times l} \end{aligned}$$

em que  $S \in \mathbb{R}_+^{k \times l}$ ,  $\Psi = \{0, 1\}$ ,  $\sum_{p=1}^k \mathbf{u}_p = 1$  e  $\sum_{q=1}^l \mathbf{v}_q = 1$ .

Wang et al. (2011) denominam  $U$  como uma matriz indicadora dos grupos de linhas,  $V$  como uma matriz indicadora dos grupos de colunas, e  $S$  como uma matriz que contém os fatores que relacionam um grupo de linhas aos grupos de colunas, e um grupo de colunas aos grupos de linhas. Note que nesse caso não é necessário uma etapa de pós-processamento para particionamento como nos outros algoritmos apresentados.

Ainda, as restrições  $\sum_{p=1}^k \mathbf{u}_p = 1$  e  $\sum_{q=1}^l \mathbf{v}_q = 1$  indicam que uma linha e uma coluna, respectivamente, têm que pertencer à algum grupo. Então, apesar dessas restrições serem semelhantes à ortonormalidade, elas não são, pois não resolvem o caso em que há uma possibilidade de haver grupos vazios, ou seja, sem nenhum elemento pertencente a este grupo, enquanto a restrição de ortonormalidade, garante que não haverá grupos vazios, seja este grupo de linhas ou de colunas.

Como  $U$  e  $V$  neste caso têm as restrições descritas, e são capazes de fornecer o particionamento de linhas e colunas, respectivamente, de forma direta, a capacidade de compactação nesse caso é semelhante ao algoritmo *k-means*, brevemente descrito no Capítulo 2. Porém, como a fatoração compacta as colunas e as relações entre grupos de linhas e colunas (matriz  $S$ ), a matriz  $X$  é compactada em  $n + kl + m$  elementos.

Como não há restrições em  $S$ , com exceção da positividade que é garantida pela positividade de  $X$ , é possível encontrar uma regra de atualização para  $S$ , e portanto, minimização de  $\mathcal{F}_3$ :

$$\begin{aligned}\nabla_S \mathcal{F}_3 &= U^T X V - U^T U S V^T V &= 0 \\ \implies U^T U S V^T V &= U^T X V \\ \implies (U^T U)^{-1} U^T U S V^T V (V^T V)^{-1} &= (U^T U)^{-1} U^T X V (V^T V)^{-1}\end{aligned}$$

$$\therefore S = (U^T U)^{-1} U^T X V (V^T V)^{-1}$$

Assim, o problema de minimização se transforma nos subproblemas de atualização de  $U$  e  $V$ . Uma estratégia semelhante aquela aplicada no clássico algoritmo *k-means* pode ser aplicada. Primeiramente, fixa-se  $S$  e  $V$  e resolve-se o problema 3 para  $U$  de forma iterativa, verificando quais dos protótipos de linhas (linhas de  $\tilde{V}$ ) mais se aproxima das linhas de  $X$ . Em seguida, fixa-se  $S$  e  $U$  para alcançar uma solução iterativa para  $V$ , através dos protótipos de colunas (colunas de  $\tilde{U}$ ), assim como mostrado no algoritmo 4.

Note que é possível interpretar a regra de atualização para  $S$ :  $(U^T U)$  é uma matriz que na diagonal contém a contagem do número de linhas em cada grupo de linhas e zeros nas demais posições,  $(U^T U)^{-1}$  é o cálculo da média parcial para cada grupo de linhas, e  $U^T X$  seleciona e soma as linhas de  $X$  para cada grupo de linhas. A mesma interpretação pode ser feita para  $(V^T V)$ ,  $(V^T V)^{-1}$  e  $XV$ . O seguinte exemplo mostra uma matriz de dados com 6 linhas e 3 colunas, particionada por 3 grupos de linhas e 2 grupos de colunas, ilustrando a interpretação da atualização de  $S$ :

$$X = \begin{bmatrix} \text{—} & \mathbf{x}_{1.} & \text{—} \\ \text{—} & \mathbf{x}_{2.} & \text{—} \\ \text{—} & \mathbf{x}_{3.} & \text{—} \\ \text{—} & \mathbf{x}_{4.} & \text{—} \\ \text{—} & \mathbf{x}_{5.} & \text{—} \\ \text{—} & \mathbf{x}_{6.} & \text{—} \end{bmatrix} = \begin{bmatrix} | & | & | \\ \mathbf{x}_{.1} & \mathbf{x}_{.2} & \mathbf{x}_{.3} \\ | & | & | \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}, U^T U = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}, V^T V = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

A matriz  $U$  do exemplo, apresenta um particionamento das 6 linhas em 3 grupos, sendo que 3 linhas pertencem ao primeiro grupo, 2 linhas ao segundo grupo, e 1 linha ao terceiro grupo, como mostra a diagonal principal da matriz  $U^T U$ , a mesma interpretação pode ser feita para  $V$  e  $V^T V$ .

$$\begin{aligned} (U^T X)V &= \left( \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \text{---} & \mathbf{x}_{1.} & \text{---} \\ \text{---} & \mathbf{x}_{2.} & \text{---} \\ \text{---} & \mathbf{x}_{3.} & \text{---} \\ \text{---} & \mathbf{x}_{4.} & \text{---} \\ \text{---} & \mathbf{x}_{5.} & \text{---} \\ \text{---} & \mathbf{x}_{6.} & \text{---} \end{bmatrix} \right) \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_{1.} + \mathbf{x}_{2.} + \mathbf{x}_{6.} \\ \mathbf{x}_{3.} + \mathbf{x}_{5.} \\ \mathbf{x}_{4.} \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \\ &= \begin{bmatrix} x_{13} + x_{23} + x_{63} & x_{11} + x_{12} + x_{21} + x_{22} + x_{61} + x_{62} \\ x_{33} + x_{53} & x_{31} + x_{32} + x_{51} + x_{52} \\ x_{43} & x_{41} + x_{42} \end{bmatrix} \end{aligned}$$

A multiplicação em  $X$  por  $U^T$  pela esquerda, representa a soma da seleção de todas as linhas de  $X$  pertencentes à um mesmo grupo de linhas. O mesmo ocorre quando multiplica-se  $V$  pela direita de  $X$ , porém, representando a soma da seleção das colunas de  $X$  pertencentes à um mesmo grupo de colunas. Sendo assim, a operação  $U^T X V$  representa a soma de todos os elementos de  $X$  pertencentes à um mesmo grupo de linhas e colunas, ou seja, por exemplo o elemento da linha 1 e coluna 1 de  $U^T X V$ , que foi calculado no exemplo, é simplesmente a soma de todos os elementos de  $X$  que pertencem ao primeiro grupo de linhas e ao primeiro grupo de colunas.

$$(U^T U)^{-1} = \begin{bmatrix} \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & 1 \end{bmatrix}, (V^T V)^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$$

$$S = (U^T U)^{-1} U^T X V (V^T V)^{-1} = \begin{bmatrix} \frac{x_{13}+x_{23}+x_{63}}{3} & \frac{x_{11}+x_{12}+x_{21}+x_{22}+x_{61}+x_{62}}{3 \times 2} \\ \frac{x_{33}+x_{53}}{2} & \frac{x_{31}+x_{32}+x_{51}+x_{52}}{2 \times 2} \\ x_{43} & \frac{x_{41}+x_{42}}{2} \end{bmatrix}$$

Com o exemplo mostrado, é possível visualizar a atualização de  $S$  de forma mais intuitiva. Por exemplo, para calcular um elemento  $s_{pq}$  de  $S$ , será simplesmente o cálculo da média de todos elementos em  $X$  que pertencem ao grupo de linhas  $p$  e ao grupo de colunas  $q$ . Dessa forma, também é possível propor uma solução iterativa para  $S$ .

O algoritmo 4 ilustra o processo de minimização para o problema 3. Nesse algoritmo considere os índices  $i = \{1, \dots, n\}$ ,  $j = \{1, \dots, m\}$ ,  $p = p' = \{1, \dots, k\}$ , e  $q = q' = \{1, \dots, l\}$ , o contador de iterações  $t$ ,  $U^{(t)}$ ,  $S^{(t)}$  e  $V^{(t)}$  como sendo as matrizes  $U$ ,  $S$  e  $V$  na iteração  $t$ , respectivamente,  $\mathcal{U}(0, 1) \in ]0, 1]$  uma função que gera valores de uma distribuição uniforme que ignora zeros, e  $\|\cdot\|^2$  é a norma frobenius para vetores. Também, as mesmas condições de convergência usada nos algoritmos anteriores pode ser aplicada nesse caso.

---

**Algoritmo 4** Algoritmo FNMTF
 

---

```

1: function FNMTF( $X, t_{max}, k, l$ )
2:   Inicialize:  $U^{(0)}, V^{(0)} \leftarrow 0, 1 \mid \sum_{p=1}^k \mathbf{u}_p = 1, \sum_{q=1}^l \mathbf{v}_q = 1, S^{(0)} \leftarrow \mathcal{U}(0, 1)$  e  $t \leftarrow 0$ .
3:   while (não convergiu) ou ( $t \leq t_{max}$ ) do
4:      $S^{(t+1)} \leftarrow (U^{(t)T} U^{(t)})^{-1} U^{(t)T} X V^{(t)} (V^{(t)T} V^{(t)})^{-1}$ 
5:      $\tilde{V} \leftarrow S^{(t+1)} V^{(t)T}$ 
6:      $(U^{(t+1)})_{ip} \leftarrow \begin{cases} 1 & p = \arg \min_{p' \in \{1, \dots, k\}} \|\mathbf{x}_{i \cdot} - \tilde{\mathbf{v}}_{p'}\|^2 \\ 0 & \text{caso contrário} \end{cases} \quad \forall i, p$ 
7:      $\tilde{U} \leftarrow U^{(t+1)} S^{(t+1)}$ 
8:      $(V^{(t+1)})_{jq} \leftarrow \begin{cases} 1 & q = \arg \min_{q' \in \{1, \dots, l\}} \|\mathbf{x}_{\cdot j} - \tilde{\mathbf{u}}_{q'}\|^2 \\ 0 & \text{caso contrário} \end{cases} \quad \forall j, q$ 
9:      $t \leftarrow t + 1$ 
10:  end while
11:  return  $U^{(t)}, S^{(t)}, V^{(t)}$ 
12: end function

```

---

A análise de complexidade de tempo do algoritmo 4 difere das demais apresentadas. Se for usado um algoritmo iterativo para o cálculo das matrizes inversas, aproveitando o fato que ambas contém elementos diferentes de 0 na diagonal principal, é possível chegar no seguinte resultado:  $\mathcal{O}\left(t_{max}(nl(m+k) + kl(n+m) + nk + ml)\right)$ . Já é possível perceber

que a complexidade do algoritmo é menor das apresentadas nos algoritmos 1, 2 e 3. Isso se explica pois a atualização de  $U$  e  $V$  é feita de forma iterativa. Ainda, há multiplicações de matrizes para o cálculo dos centróides  $\tilde{U}$  e  $\tilde{V}$ , e para atualização de  $S$ , que podem ser calculados de forma iterativa, melhorando ainda mais a complexidade de tempo do algoritmo.

### 3.4 Considerações Finais

É importante ressaltar, que nenhum dos algoritmos apresentados são capazes de garantir convergência para um mínimo global dos problemas apresentados, então, é possível encontrar diversas soluções em execuções diferentes dos algoritmos. Essa é uma limitação também presente nos algoritmos de agrupamento clássicos.

Além disso, a fim de melhor motivar a proposta dos novos algoritmos, apresentados no capítulo 4, é interessante compreender os algoritmos ONMTF e FNMTF em termos de suas capacidades de quantização do espaço dos dados e de geração de informação sobre os dados e em termos do processo de descoberta de cogrupos.

Do ponto de vista de quantização do espaço dos dados, a quantidade de informação que precisa ser armazenada é dependente da organização da matriz  $S$ , ou seja, do número de grupos de linhas  $k$  e do número de colunas  $l$  necessários para explicar a matriz de dados original. Como será discutido mais à frente neste trabalho, para determinados tipos de organização de cogrupos, especificamente aqueles em que há sobreposição de colunas nos grupos de colunas, os algoritmos implementados sob essas estratégias exigem uma quantidade de grupos de colunas ( $l$ ) que pode ser maior do que o número de grupos de colunas desejados. Os algoritmos propostos tem o objetivo de superar essa limitação, permitindo que a quantização do espaço seja mais próxima da desejada em termos de grupos de colunas.

Do ponto de vista de geração de informação sobre os dados, os algoritmos ONMTF e FNMTF são capazes de fornecer como as linhas da matriz se organizam em grupos e como as colunas da matriz se organizam em grupos. Transferindo essa informação para um contexto de aplicação, significa dizer que os algoritmos são capazes de explicar como os dados se organizam no espaço (matriz  $U$ ) e, intuitivamente e sob uma forma de organização, como grupos de atributos desses dados (matriz  $V$ ) podem estar associados (por meio da matriz  $S$ ) à essa organização dos dados. Esse tipo de informação, também presente nos

algoritmos propostos, não é explicitamente fornecida por algoritmos de agrupamento como *k-means* e *fuzzy-k-means*, brevemente discutidos no capítulo 2.

Do ponto de vista do processo de descoberta dos cogrupos, os algoritmos ONMTF e FNMTF são capazes de considerar simultaneamente ambas organizações na resolução do problema de minimização do erro de aproximação da matriz original. Porém, possuem um processo que pode ser caracterizado por um tipo de interdependência entre grupos de linhas, que é na realidade o causador da possibilidade de chegar a uma quantidade de grupos de colunas ( $l$ ) maior do que é realmente necessário para explicar os dados que possuem uma organização de cogrupos com sobreposição de colunas. Esta é a segunda meta de superação obtida nos algoritmos propostos, que por sua formulação, são capazes de resolver o problema de fatoração das matrizes de maneira independente para cada grupo de linhas.

## 4 Fatoração de matrizes não-negativas para coagrupamento com sobreposição unidimensional

Como discutido no capítulo 3, a fatoração de matrizes aplicada ao problema de agrupamento ou coagrupamento pode ser analisada sob, pelo menos, três aspectos: quantização do espaço dos dados; geração de informação sobre os dados; processo de descoberta de cogrupos. Também, para cada uma dessas análises, diferentes estratégias apresentam vantagens e desvantagens.

Com o intuito de apresentar uma alternativa às estratégias de fatoração de matrizes não-negativas para coagrupamento presentes na literatura, objetivando superar algumas das dificuldades apresentadas por eles, propõe-se duas novas estratégias que adicionam  $k$  matrizes  $V$  na fatoração, ao invés de uma única matriz  $V$ . Basicamente, cada uma dessas matrizes representará uma organização de cogrupos de colunas independente, de maneira que aumentar-se-á a flexibilidade para o estabelecimento de relações entre cogrupos de linhas e cogrupos de colunas. As estratégias são:

- OvNMTF: uma estratégia de fatoração tripla de matrizes não-negativas com sobreposição unidimensional, baseada na estratégia BVD;
- BinOvNMTF: uma estratégia de fatoração binária tripla de matrizes não-negativas com sobreposição unidimensional, baseada na estratégia FNMTF.

Sob o ponto de vista de resolução de problemas de coagrupamento, o objetivo dessas estratégias é o mesmo das estratégias já apresentadas neste texto, qual seja, encontrar grupos de linhas e colunas de forma simultânea. Entretanto, visto que se tem a liberdade de organizar um conjunto de matrizes  $V$ , que abstrai grupos e colunas, para ser associado a cada matriz  $U$ , que abstrai os grupos de linhas, é possível que grupos de colunas diferentes sejam formados, a depender do grupo de linhas sendo considerado. Isso significa que problemas em que há interseção de colunas na formação de cogrupos poderão ser adequadamente tratados pelas estratégias propostas.

Formalmente, considere uma matriz de dados  $X \in \mathbb{R}^{n \times m}$  contendo números reais positivos com  $n$  linhas e  $m$  colunas, formada por um conjunto de vetores de linhas  $\mathcal{N} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  e um conjunto de vetores de colunas  $\mathcal{M} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , e as relações existentes entre cada linha  $x$  e cada coluna  $y$  são representadas por  $x_{ij}$  considerando os índices  $i = \{1, \dots, n\}$  e  $j = \{1, \dots, m\}$ . O objetivo é esta matriz é formada por um conjunto de vetores de linhas  $\mathcal{N} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  e um conjunto de vetores de colunas

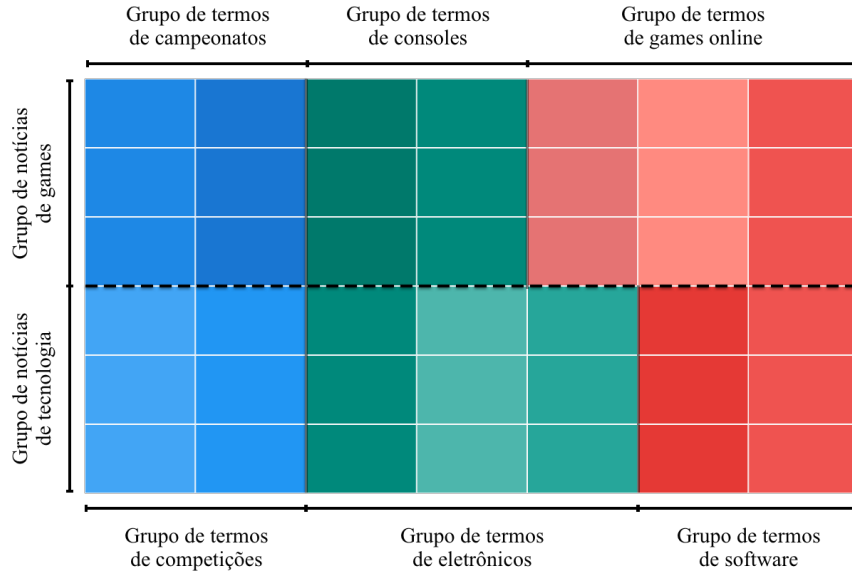


$\mathcal{M} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , e as relações existentes entre cada linha  $x$  e cada coluna  $y$  são representadas por  $x_{ij}$  considerando os índices  $i = \{1, \dots, n\}$  e  $j = \{1, \dots, m\}$ , que é justamente um valor da matriz  $X$ . Cada valor em  $x_{ij}$  representa, então, a relação existente entre pares de elementos em algum contexto de interesse. Modificando o objetivo como foi apresentado na introdução ao capítulo 3, sendo então, encontrar  $k$  partições de  $\mathcal{N}$ , denotadas pelos subconjuntos ordenados  $\mathcal{K}_p \subseteq \mathcal{N}$ ,  $k \times l$  partições para  $\mathcal{M}$  em cada  $\mathcal{K}_p$ , denotadas pelos subconjuntos ordenados  $\mathcal{L}_{pq} \subseteq \mathcal{M}$ , considerando os índices  $p = \{1, \dots, k\}$  e  $q = \{1, \dots, l\}$ . Então, os subconjuntos  $\{\mathcal{K}_1, \dots, \mathcal{K}_k\}$  e  $\{\mathcal{L}_{11}, \dots, \mathcal{L}_{1l}, \dots, \mathcal{L}_{k1}, \dots, \mathcal{L}_{kl}\}$  são os cogrupos de linhas e colunas, respectivamente.

Observe na figura 4 uma representação gráfica da estrutura de cogrupos com sobreposição de colunas, contextualizado em uma aplicação de análise de textos. Note que, na realidade, trata-se de busca por uma solução adequada para um problema com sobreposição unidimensional, ou seja, sobreposição de colunas ou de linhas, já que a sobreposição de linhas é justamente o problema transposto.

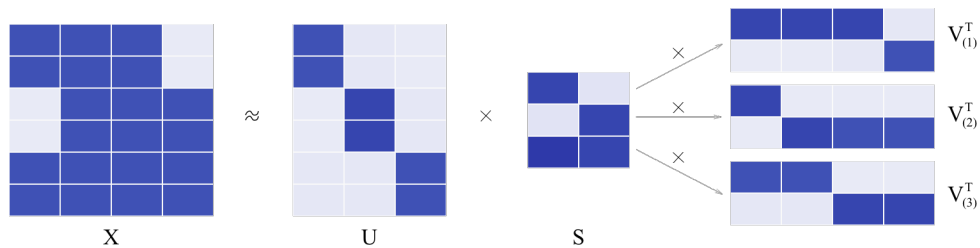
A figura 4 mostra dois grupos de documentos, dos assuntos games e tecnologia, com três documentos em cada. Ainda, este exemplo contém 6 grupos de palavras, divididos igualmente entre os grupos de documentos. Então, pode-se dizer que os grupos de palavras intitulados “campeonatos”, “consoles” e “games online” caracterizam o grupo de documentos sobre games, assim como os grupos de palavras intitulados “competições”, “eletrônicos” e “software” caracterizam o grupo de documentos sobre tecnologia. Note que há sobreposição entre os primeiros três grupos de palavras do grupo de documentos com os três primeiros grupos de palavras do grupo de documentos de tecnologia. Também, é possível observar que apesar das palavras serem as mesmas para todos os documentos, como há independência entre os grupos de palavras para cada grupo de documentos, as palavras caracterizaram grupos de palavras diferentes, como no exemplo, os grupos de palavras intitulados como “games online” e “software”.

Figura 4 – Representação gráfica do problema de coagruamento com sobreposição unidimensional e contextualização do domínio de documentos (notícias)



Da mesma maneira que foi ilustrado no capítulo 3 que é possível derivar interpretações intuitivas da análise das combinações de matrizes geradas na fatoração produzido pelos algoritmos lá discutidos, para os algoritmos discutidos no presente capítulo, interpretações análogas podem ser feitas, porém, consirando a existências das várias matrizes  $V$ . Assumindo novamente que uma matriz de entrada representação a relação “documento por palavras” (linhas por colunas), cada coluna das  $k$  matrizes  $UI_{(p)}S, \forall p = 1, \dots, k$ , captura a ideia de uma base para representação de grupos de palavras, descritos nas matrizes  $V_{(p)}$ ; e cada linha em  $\sum_{p=1}^k SV_{(p)}^T$ , captura a ideia de uma base de representação de grupos de documentos. A representação gráfica para o resultado de uma fatoração de matrizes que permite esse raciocínio é apresentada na figura 5.

Figura 5 – Fatoração da matriz original de dados  $X$  em cinco outras matrizes:  $U$ ,  $S$ ,  $V_1$ ,  $V_2$  e  $V_3$



Todo o universo de interpretações delineado para os algoritmos da literatura podem ser estendidos para esse novo contexto de problema de agrupamento. No entanto, agora

existem três matrizes para determinar os grupos de palavras, cada uma delas responsável por determinar os grupos de palavras para cada grupo de documentos. Na figura 5 considere que uma célula com cor escura representa a existência de uma relação entre linha e coluna, e uma célula com cor clara representa a inexistência de uma relação entre linha e coluna, e que essas relações são estabelecidas adequadamente em cada contexto de aplicação. Tem-se então, um conjunto de seis documentos e quatro palavras. A matriz  $U$  pode ser interpretada como descrito anteriormente, uma matriz “documentos por grupos de documentos”, com seis documentos agrupados em três grupos ( $k = 3$ ): os dois primeiros documentos no primeiro grupo, os dois próximos no segundo grupo e os dois últimos no terceiro grupo. Já as matrizes  $V_{(p)}^T$  têm uma interpretação levemente diferente, ainda podem ser interpretadas como matrizes de “grupos de palavras por palavras”, sendo dois grupos de palavras ( $l = 2$ ), porém, a matriz  $V_{(1)}^T$ , por exemplo, contém os grupos de palavras para o primeiro grupo de documentos apenas (composto pelas linhas 1 e 2 de  $X$ ), no primeiro grupo de palavras estão as colunas 1, 2 e 3, e no segundo, a coluna 4. O mesmo raciocínio pode ser utilizado para as matrizes  $V_{(2)}^T$  e  $V_{(3)}^T$ . Por fim, para a matriz  $S$  pode ser realizada a mesma interpretação, representando uma relação entre grupos de documentos e grupos de palavras, com a ressalva de que irão existir  $k \times l$  grupos de palavras, ou seja, a linha  $p$  da matriz  $S$  irá representar a relação entre o  $p$ -ésimo grupo de documentos e os grupos de palavras encontrados em  $V_{(p)}^T$ .

O restante deste capítulo é destinado a apresentar a formulação dos problemas para cada uma das estratégias, incluindo a apresentação de uma derivação teórica para um algoritmo que implementa o processo de resolução para os problemas. Finalmente, as considerações finais apresentam uma discussão sobre o diferencial dessas estratégias no que diz respeito aos três aspectos citados no início do capítulo.

#### 4.1 Fatoração Tripla de Matrizes Não-negativas com Sobreposição Unidimensional

Baseado no problema 1 (LONG; ZHANG; YU, 2005), o problema 4 é apresentado neste trabalho, e recebe aqui o nome de Fatoração Tripla de Matrizes Não-negativas Sobrepostas (Overlapped Non-negative Matrix Tri-factorization - OvNMTF).

**Problema 4** (Problema de Fatoração Tripla de Matrizes Não-negativas Sobrepostas).

$$\begin{aligned} \mathcal{F}_4(U, S, V_{(1)}, \dots, V_{(k)}) &= \min_{U, S, V_{(1)}, \dots, V_{(k)}} \left\| X - U \sum_{p=1}^k I_{(p)} S V_{(p)}^T \right\|_F^2 \\ \text{sujeito a} \quad &U \geq 0, S \geq 0, \\ &V_{(p)} \geq 0, \quad \forall p \end{aligned}$$

em que  $U \in \mathbb{R}_+^{n \times k}$ ,  $S \in \mathbb{R}_+^{k \times l}$ ,  $V_{(p)} \in \mathbb{R}_+^{m \times l}$ ,  $p = \{1, \dots, k\}$  é o índice para o conjunto de matrizes  $\{V_{(1)}, \dots, V_{(k)}\}$ ,  $I_{(p)} \in \{0, 1\}^{k \times k}$  é uma matriz identidade seletora, e  $\|\cdot\|_F$  denota a norma de Frobenius. Na formulação, o papel das matrizes  $I_k$  seletoras é, justamente, organizar a base de grupos de linhas, de forma que cada um deles seja otimizado de acordo com uma das matrizes  $V_k$ .

É possível perceber, que neste caso, diferente dos apresentados no capítulo 3, têm menor capacidade de compactação, porém, com maior nível de detalhamento. Isso se explica, pois a partir dos  $nm$  elementos da matriz de dados, são gerados  $nk + kl + klm$  para representá-los, diferente.

Desde que o problema 4 é semelhante ao problema 2, é esperado que ele possa também ser resolvido por meio de regras de atualização multiplicativas. Assim, seguindo o exposto em (YOO; CHOI, 2010), a derivação das regras é aqui apresentada por meio de um abordagem baseada no cálculo do gradiente. Contudo, para o cálculo do gradiente de  $\mathcal{F}_4$  a estratégia usada aqui é a apresentada em Yoo e Choi (2010), de forma que  $\nabla \mathcal{F}_4 = [\nabla \mathcal{F}_4]^+ - [\nabla \mathcal{F}_4]^-$ .

Expandindo  $\mathcal{F}_4$ , com base nas propriedade de traço de matrizes, para tornar o cálculo do gradiente mais simples, é possível obter:

$$\begin{aligned} \mathcal{F}_4 &= \text{tr} \left[ (X - U \sum_{p=1}^k I_{(p)} S V_{(p)}^T)^T (U \sum_{p=1}^k I_{(p)} S V_{(p)}^T) \right] \\ &= \text{tr}(X^T X) - 2\text{tr}(X^T U \sum_{p=1}^k I_{(p)} S V_{(p)}^T) + \text{tr}(\sum_{p=1}^k V_{(p)} S^T I_{(p)} U^T U \sum_{p'=1}^k I_{(p')} S V_{(p')}^T) \end{aligned}$$

Note que a matriz seletora tem a seguinte propriedade:  $I_{(p)}^T = I_{(p)}$ .

---

Considere as seguintes igualdades para o cálculo dos gradientes, sendo  $A, Q, B$  e  $C$  matrizes de quaisquer dimensões adequadas para a realização das multiplicações:

$$\begin{aligned} \nabla_Q \text{tr}(AQB) &= A^T B^T \\ \nabla_{Q^T} \text{tr}(AQB) &= BA \end{aligned} \tag{5}$$

$$\begin{aligned} \nabla_Q \text{tr}(AQBQ^T C) &= A^T C^T Q B^T + C A Q B \\ \nabla_{Q^T} \text{tr}(AQBQ^T C) &= B Q^T C A + B^T Q^T A^T C^T \end{aligned} \tag{6}$$

Para o cálculo de  $\nabla_U \mathcal{F}_4$ , considere as partes positiva e negativa desse gradiente,  $[\nabla_U \mathcal{F}_4]^+$  e  $[\nabla_U \mathcal{F}_4]^-$ , respectivamente. Usando a igualdade da equação 5, com  $A = X^T$ ,  $B = \sum_{p=1}^k I_{(p)} S V_{(p)}^T$  e  $Q = U$ , é possível obter  $[\nabla_U \mathcal{F}_4]^-$ , como segue.

$$\begin{aligned} [\nabla_U \mathcal{F}_4]^- &= -2 \nabla_U \left( \text{tr} [X^T U \sum_{p=1}^k I_{(p)} S V_{(p)}^T] \right) \\ &= -2 X \sum_{p=1}^k V_{(p)} S^T I_{(p)} \end{aligned}$$

Para o cálculo de  $[\nabla_U \mathcal{F}_4]^+$ , é utilizado a igualdade da equação 6, com  $A = \sum_{p=1}^k V_{(p)} S^T I_{(p)}$ ,  $B = I$ ,  $Q = U^T$  e  $C = \sum_{p'=1}^k I_{(p')} S V_{(p')}^T$ . Então tem-se

$$\begin{aligned} [\nabla_U \mathcal{F}_4]^+ &= \nabla_U \left( \text{tr} (X^T X) + \text{tr} \left[ \sum_{p=1}^k V_{(p)} S^T I_{(p)} U^T U \sum_{p'=1}^k I_{(p')} S V_{(p')}^T \right] \right) \\ &= U \sum_{p'=1}^k I_{(p')} S V_{(p')}^T \sum_{p=1}^k V_{(p)} S^T I_{(p)} \\ &\quad + U \sum_{p=1}^k I_{(p)} S V_{(p)}^T \sum_{p'=1}^k V_{(p')} S^T I_{(p')} \\ &= 2U \sum_{p=1}^k \sum_{p'=1}^k I_{(p)} S V_{(p)}^T V_{(p')} S^T I_{(p')} \end{aligned}$$

De forma similar, para o cálculo de  $\nabla_S \mathcal{F}_4$ , considere as partes positiva e negativa,  $[\nabla_S \mathcal{F}_4]^+$  e  $[\nabla_S \mathcal{F}_4]^-$ , respectivamente. Usando a igualdade da equação 5 para todas as partes da soma de  $p = 1, \dots, k$ , com  $A = X^T U I_{(p)}$ ,  $Q = S$ , e  $B = V_{(p)}^T$ , é possível obter  $[\nabla_S \mathcal{F}_4]^-$ , como segue:

$$\begin{aligned} [\nabla_S \mathcal{F}_4]^- &= -2 \nabla_S \left( \text{tr} [X^T U I_{(1)} S V_{(1)}^T] + \dots + \text{tr} [X^T U I_{(k)} S V_{(k)}^T] \right) \\ &= -2 \left( I_{(1)} U^T X V_{(1)} + \dots + I_{(k)} U^T X V_{(k)} \right) \\ &= -2 \sum_{p=1}^k I_{(p)} U^T X V_{(p)} \end{aligned}$$

Para o cálculo de  $[\nabla_S \mathcal{F}_4]^+$ , é utilizado a igualdade da equação 6 para todas as partes da soma de  $p = p' = 1, \dots, k$ , com  $A = V_{(p)}$ ,  $Q = S^T$ ,  $B = I_{(p)} U^T U I_{(p')}$  e  $C = V_{(p')}^T$ . Então tem-se

$$\begin{aligned} [\nabla_S \mathcal{F}_4]^+ &= \nabla_S \left( \text{tr} (X^T X) + \text{tr} [V_{(1)} S^T I_{(1)} U^T U I_{(1)} S V_{(1)}^T] + \dots + \text{tr} [V_{(1)} S^T I_{(1)} U^T U I_{(k)} S V_{(k)}^T] \right. \\ &\quad \left. + \dots + \text{tr} [V_{(k)} S^T I_{(k)} U^T U I_{(1)} S V_{(1)}^T] + \dots + \text{tr} [V_{(k)} S^T I_{(k)} U^T U I_{(k)} S V_{(k)}^T] \right) \\ &= (I_{(1)} U^T U I_{(1)} S V_{(1)}^T V_{(1)} + I_{(1)} U^T U I_{(1)} S V_{(1)}^T V_{(1)}) \\ &\quad + \dots + (I_{(1)} U^T U I_{(k)} S V_{(k)}^T V_{(1)} + I_{(k)} U^T U I_{(1)} S V_{(1)}^T V_{(k)}) \\ &\quad + \dots + (I_{(k)} U^T U I_{(1)} S V_{(1)}^T V_{(k)} + I_{(1)} U^T U I_{(k)} S V_{(k)}^T V_{(1)}) \\ &\quad + \dots + (I_{(k)} U^T U I_{(k)} S V_{(k)}^T V_{(k)} + I_{(k)} U^T U I_{(k)} S V_{(k)}^T V_{(k)}) \\ &= 2 \sum_{p=1}^k \sum_{p'=1}^k I_{(p)} U^T U I_{(p')} S V_{(p')}^T V_{(p)} \end{aligned}$$

Finalmente, para o cálculo de  $\nabla_{V_{(p)}} \mathcal{F}_4$ , considere as partes positiva e negativa,  $[\nabla_{V_{(p)}} \mathcal{F}_4]^+$  e  $[\nabla_{V_{(p)}} \mathcal{F}_4]^-$ , respectivamente. Usando a igualdade da equação 5, com  $A = X^T U I_{(p)} S$ ,  $Q = V_{(p)}^T$ , e  $B = I$ ,  $\forall p$ , é possível obter  $[\nabla_{V_{(p)}} \mathcal{F}_4]^-$ , da seguinte forma:

$$\begin{aligned} [\nabla_{V_{(p)}} \mathcal{F}_4]^- &= -2\nabla_{V_{(p)}} \left( \text{tr}[X^T U I_{(1)} S V_{(1)}^T] + \cdots + \text{tr}[X^T U I_{(k)} S V_{(k)}^T] \right) \\ &= -2X^T U I_{(p)} S \end{aligned}$$

Fixando a derivação para  $[\nabla_{V_{(p)}} \mathcal{F}_4]^+$  é nota-se que há dois casos diferentes para a derivação. O caso em que  $p = p'$  e o caso em que  $p \neq p'$ . Para o caso em que  $p = p'$ , é utilizada a igualdade da equação 6, com  $A = I$ ,  $Q = V_{(p)}$ ,  $B = S^T I_{(p)} U^T U I_{(p')} S$  e  $C = I$ ,  $\forall p, p'$ , de forma que

$$\begin{aligned} [\nabla_{V_{(p)}} \mathcal{F}_4]_{p=p'}^+ &= \nabla_{V_{(p)}} \left( \text{tr}[V_{(1)} S^T I_{(1)} U^T U I_{(1)} S V_{(1)}^T] + \cdots + \text{tr}[V_{(k)} S^T I_{(k)} U^T U I_{(k)} S V_{(k)}^T] \right) \\ &= \nabla_{V_{(p)}} \left( \text{tr}[V_{(p)} S^T I_{(p)} U^T U I_{(p)} S V_{(p)}^T] \right) \\ &= 2V_{(p)} S^T I_{(p)} U^T U I_{(p)} S \end{aligned}$$

Para o caso em que  $p \neq p'$ , é usada a igualdade da equação 5 em duas outras situações, aquela em que  $V_{(p)}$  esta localizado à esquerda, então,  $Q = V_{(p)}$ ,  $A = I$  e  $B = S^T I_{(p)} U^T U I_{(p')} S$ ; e aquela em que  $V_{(p)}$  esta localizado à direita, então,  $Q = V_{(p)}^T$ ,  $A = V_{(p')} S^T I_{(p')} U^T U I_{(p)} S V_{(p)}^T$ ,  $\forall p, p'$ . Assim,

$$\begin{aligned} [\nabla_{V_{(p)}} \mathcal{F}_4]_{p \neq p'}^+ &= \nabla_{V_{(p)}} \left( \text{tr}[V_{(1)} S^T I_{(1)} U^T U I_{(2)} S V_{(2)}^T] + \cdots + \text{tr}[V_{(1)} S^T I_{(1)} U^T U I_{(k)} S V_{(k)}^T] \right. \\ &\quad \left. + \cdots + \text{tr}[V_{(k)} S^T I_{(k)} U^T U I_{(1)} S V_{(1)}^T] + \cdots + \text{tr}[V_{(k)} S^T I_{(k)} U^T U I_{(k-1)} S V_{(k-1)}^T] \right) \\ &= \sum_{p' \neq p} \left[ \nabla_{V_{(p)}} \left( \text{tr}[V_{(p)} S^T I_{(p)} U^T U I_{(p')} S V_{(p')}^T] \right) \right. \\ &\quad \left. + \nabla_{V_{(p)}} \left( \text{tr}[V_{(p')} S^T I_{(p')} U^T U I_{(p)} S V_{(p)}^T] \right) \right] \\ &= \sum_{p' \neq p} 2(V_{(p')} S^T I_{(p')} U^T U I_{(p)} S) \end{aligned}$$

Então, é possível calcular  $[\nabla_{V_{(p)}} \mathcal{F}_4]^+$ ,  $\forall p$ , fazendo:

$$\begin{aligned} [\nabla_{V_{(p)}} \mathcal{F}_4]^+ &= [\nabla_{V_{(p)}} \mathcal{F}_4]_{p=p'}^+ + [\nabla_{V_{(p)}} \mathcal{F}_4]_{p \neq p'}^+ \\ &= 2(V_{(p)} S^T I_{(p)} U^T U I_{(p)} S) + \sum_{p' \neq p} V_{(p')} S^T I_{(p')} U^T U I_{(p)} S \\ &= 2\left(\sum_{p'=1}^k V_{(p')} S^T I_{(p')} U^T U I_{(p)} S\right) \end{aligned}$$

O resultado final dos gradientes para  $U, S, V_{(p)}, \forall p \in \{1, \dots, k\}$  são apresentados nas equações 7, 8 e 9, respectivamente.

$$\nabla_U \mathcal{F}_4 = 2\left(-X \sum_{p=1}^k V_{(p)} S^T I_{(p)} + U \sum_{p=1}^k \sum_{p'=1}^k I_{(p)} S V_{(p)}^T V_{(p')} S^T I_{(p')}\right) \quad (7)$$

$$\nabla_S \mathcal{F}_4 = 2 \left( - \sum_{p=1}^k I_{(p)} U^T X V_{(p)} + \sum_{p=1}^k \sum_{p'=1}^k I_{(p)} U^T U I_{(p')} S V_{(p')}^T V_{(p)} \right) \quad (8)$$

$$\nabla_{V_{(p)}} \mathcal{F}_4 = 2 \left( - X^T U I_{(p)} S + \sum_{p'=1}^k V_{(p')} S^T I_{(p')} U^T U I_{(p)} S \right) \quad (9)$$

Sendo assim, as regras de atualização para as matrizes  $U, V_{(p)}, \forall p \in \{1, \dots, k\}$ , são mostradas nas equações 10, 11 e 12, apresentadas no algoritmo 5, o qual implementa o processo de minimização para o problema 4. Nesse algoritmo  $t$  é o contador de iterações,  $U^{(t)}, S^{(t)}$  e  $V_{(p)}^{(t)}$  são as matrizes  $U, S$  e  $V_{(p)}$ , na iteração  $t$ , respectivamente,  $\mathcal{U}(0, 1) \in ]0, 1]$  uma função que gera valores de uma distribuição uniforme que ignora zeros, e  $\odot$  é o produto de Hadamard.

---

**Algoritmo 5** Algoritmo baseado em atualização multiplicativa para solução do OvNMTF

---

```

1: function OvNMTF( $X, t_{max}$ )
2:   Initialize:  $U^{(0)} \geq \mathcal{U}(0, 1), S^{(0)} \geq \mathcal{U}(0, 1), V_{(p)}^{(0)} \geq \mathcal{U}(0, 1), \forall p$  e  $t \leftarrow 0$ .
3:   while (não convergiu) ou ( $t \leq t_{max}$ ) do
4:
```

$$U^{(t+1)} \leftarrow U^{(t)} \odot \frac{\sum_{p=1}^k X V_{(p)}^{(t)} S^{(t)T} I_{(p)}}{\sum_{p=1}^k \sum_{p'=1}^k U^{(t)} I_{(p)} S^{(t)} V_{(p)}^{(t)T} V_{(p')}^{(t)} S^{(t)T} I_{(p')}} \quad (10)$$

```

5:   for  $p \leftarrow 1, k$  do
```

$$V_{(p)}^{(t+1)} \leftarrow V_{(p)}^{(t)} \odot \frac{X^T U^{(t+1)} I_{(p)} S^{(t)}}{\sum_{p'=1}^k V_{(p')} S^T I_{(p')} U^T U I_{(p)} S} \quad (11)$$

```

7:   end for
```

```

8:
```

$$S^{(t+1)} \leftarrow S^{(t)} \odot \frac{\sum_{p=1}^k I_{(p)} U^{(t+1)T} X V_{(p)}^{(t+1)}}{\sum_{p=1}^k \sum_{p'=1}^k I_{(p)} U^{(t+1)T} U^{(t+1)} I_{(p')} S^{(t)} V_{(p)}^{(t+1)T} V_{(p')}^{(t+1)}} \quad (12)$$

```

9:    $t \leftarrow t + 1$ 
```

```

10: end while
```

```

11: return  $U^{(t)}, S^{(t)}, V_{(1)}^{(t)}, \dots, V_{(k)}^{(t)}$ 
```

```

12: end function
```

---

A inicialização dos elementos das matrizes  $U, S$  e  $V_{(1)}, \dots, V_{(k)}$  são gerados através de uma distribuição uniforme que ignora zeros ( $\mathcal{U}(0, 1) \in ]0, 1]$ ). Como condições para assumir a convergência, assim como os algoritmos apresentados no capítulo 3, considera-se a diferença do erro de aproximação em duas iterações consecutivas menor ou igual a um  $\epsilon$ :

$$\left\| X - U^{(t)} \sum_{p=1}^k I_{(p)} S^{(t)} V_{(p)}^{(t)T} \right\|_F^2 - \left\| X - U^{(t+1)} \sum_{p=1}^k I_{(p)} S^{(t+1)} V_{(p)}^{(t+1)T} \right\|_F^2 \leq \epsilon$$

O algoritmo também pára caso a  $t$ -ésima iteração seja igual ao número máximo de iterações ( $t_{max}$ ).

O particionamento, neste caso, não é direto. Então é necessário um processo de pós-processamento para a matriz  $U$  e para as matrizes  $V_{(1)}, \dots, V_{(k)}$ . Um modo simples de obter o particionamento para as linhas, já descrito no capítulo 3, é o seguinte:

$$\mathcal{K}_p = \mathcal{K}_p + \{x_{i.}\} \mid p = \arg \max_{p'} \mathbf{u}_i. \quad \forall i = \{1, \dots, n\}, \forall p, p' = \{1, \dots, k\}$$

Pode-se usar a mesma estratégia para o particionamento das colunas:

$$\mathcal{L}_{pq} = \mathcal{L}_{pq} + \{x_{.j}\} \mid q = \arg \max_{q'} \mathbf{v}_{(p)j}. \quad \forall j = \{1, \dots, m\}, \forall q, q' = \{1, \dots, l\}, \forall p = \{1, \dots, k\}$$

## 4.2 Fatoração Binária Tripla de Matrizes Não-negativas com Sobreposição Unilateral

A segunda estratégia proposta nesta dissertação, segue os pressupostos estabelecidos por Wang et al. (2011) no problema 3. O problema 5, então denominado Fatoração Binária Tripla de Matrizes Não-negativas com Sobreposição Unilateral (Unilateral Overlapped Binary Non-negative Matrix Tri-factorization - BinOvNMTF) também está associado a diminuição da flexibilidade da solução de coagrupamento apresentada no que diz respeito à relação de associação entre grupos de linhas e grupos de colunas, a qual aqui é binária. Porém, mantém a independência do estabelecimento de diferentes bases para os grupos de linhas, uma vez que também assume a existência de  $k$  matrizes  $V$ .

**Problema 5** (Problema de Fatoração Binária Tripla de Matrizes Não-negativas Sobrepostas).

$$\begin{aligned} \mathcal{F}_5(U, S, V_{(1)}, \dots, V_{(k)}) &= \min_{U, S, V_{(1)}, \dots, V_{(k)}} \left\| X - U \sum_{p=1}^k I_{(p)} S V_{(p)}^T \right\|_F^2 \\ \text{sujeito a} \quad &U \in \Psi^{n \times k}, \\ &V_{(p)} \in \Psi^{m \times l}, \quad \forall p \end{aligned}$$

em que  $\Psi = \{0, 1\}$ ,  $\sum_{p=1}^k \mathbf{u}_{p.} = 1$  e  $\sum_{q=1}^l \mathbf{v}_{(p)q.} = 1, \forall p, p = \{1, \dots, k\}$  e  $q = \{1, \dots, l\}$  são os índices que iteram no número de linhas e colunas, respectivamente,  $I_{(p)} \in \{0, 1\}^{k \times k}$  é uma matriz identidade seletora, e  $\|\cdot\|_F$  denota a norma de Frobenius para matrizes. Ainda, as restrições  $\sum_{p=1}^k \mathbf{u}_{p.} = 1$  e  $\sum_{q=1}^l \mathbf{v}_{(p)q.} = 1, \forall p$  garantem que uma linha e uma coluna de um grupo, respectivamente, têm que pertencer à algum grupo.



Em termos de compactação, a interpretação é semelhante à interpretação feita para o problema 4, porém, por causa das restrições serem binárias, os  $nm$  elementos da matriz de dados serão compactados em  $n + kl + km$  elementos.

A mesma estratégia de derivação utilizada no algoritmo 4, também pode ser feita para propor uma solução para o problema 5: solucionar o problema para  $S$  para obter subproblemas que podem ser solucionados através de uma estratégia iterativa. Então, recapitulando o gradiente  $[\nabla_S \mathcal{F}_4]^+$ , já calculado para solução do problema 4, será o mesmo para  $\mathcal{F}_5$ , porém passível de simplificação, como  $(U^T U)$ , para o problema 5 que tem restrições binárias em  $U$ , é uma matriz que contém zeros em todos elementos, com exceção da diagonal principal:

$$\begin{aligned} [\nabla_S \mathcal{F}_5]^+ &= \sum_{p=1}^k \sum_{p'=1}^k I_{(p)} U^T U I_{(p')} S V_{(p')}^T V_{(p)} \\ &= \sum_{p=1}^k \sum_{p'=1}^k \left( I_{(p)} \begin{bmatrix} (U^T U)_{11} & & 0 \\ & \ddots & \\ 0 & & (U^T U)_{kk} \end{bmatrix} I_{(p')} \right) S V_{(p')}^T V_{(p)} \end{aligned}$$

Note que se  $p \neq p'$ ,  $I_{(p)} U^T U I_{(p')}$  e toda expressão, será uma matriz de zeros, por exemplo, se  $p = 1$  e  $p' = 2$ :

$$\begin{aligned} I_{(1)} \begin{bmatrix} (U^T U)_{11} & & 0 \\ & \ddots & \\ 0 & & (U^T U)_{kk} \end{bmatrix} I_{(2)} &= \begin{bmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} (U^T U)_{11} & & 0 \\ & \ddots & \\ 0 & & (U^T U)_{kk} \end{bmatrix} \begin{bmatrix} 0 & \dots & 0 \\ 0 & 1 & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \\ &= (0)_{kk} \end{aligned}$$

$$\therefore [\nabla_S \mathcal{F}_5]^+ = \sum_{p=1}^k I_{(p)} U^T U I_{(p)} S V_{(p)}^T V_{(p)}$$

Assim, é possível escrever o gradiente  $\nabla_S \mathcal{F}_5$ , a fim de encontrar uma expressão para atualização de  $S$ .

$$\begin{aligned} [\nabla_S \mathcal{F}_5]^+ &= -2 \sum_{p=1}^k I_{(p)} U^T X V_{(p)} + 2 \sum_{p=1}^k I_{(p)} U^T U I_{(p)} S V_{(p)}^T V_{(p)} = 0 \\ \implies \sum_{p=1}^k I_{(p)} U^T U I_{(p)} S V_{(p)}^T V_{(p)} &= \sum_{p=1}^k I_{(p)} U^T X V_{(p)} \end{aligned}$$

Como os termos  $I_{(p)} U^T X V_{(p)}$  e  $I_{(p)} U^T U I_{(p)} S V_{(p)}^T V_{(p)}$ ,  $\forall p$ , por serem multiplicados por  $I_{(p)}$  pela esquerda, correspondem a cada linha da matriz que será a atualização de

$S$ , aqui denotado por  $\tilde{S} = I_{(p)}S$ . Observe que  $I_{(p)}(U^T U)^{-1}I_{(p)} = I_{(p)}(U^T U)^{-1}$ , dado a estrutura de  $(U^T U)$ .

$$\begin{aligned} I_{(p)}U^T U \tilde{S} V_{(p)}^T V_{(p)} &= I_{(p)}U^T X V_{(p)} \\ \tilde{S} &= I_{(p)}(U^T U)^{-1}I_{(p)}U^T X V_{(p)} V_{(p)}^T V_{(p)} \\ &= I_{(p)}(U^T U)^{-1}U^T X V_{(p)} V_{(p)}^T V_{(p)} \end{aligned}$$

$$\therefore S = \sum_{p=1}^k I_{(p)}(U^T U)^{-1}U^T X V_{(p)} V_{(p)}^T V_{(p)}$$

É importante ressaltar que a mesma interpretação realizada para atualização de  $S$  no algoritmo 4, também é possível ser transferida para esse contexto.

Assim, é possível transformar o problema 5 em subproblemas para atualização de  $U$  e  $V_{(1)}, \dots, V_{(k)}$ . Usando a mesma estratégia iterativa adotada no algoritmo 4: obtém-se os protótipos de linhas e verifica-se quais linhas de  $X$  mais se aproximam de cada um dos protótipos, para atualização de  $U$ ; então, é feito o mesmo processo para atualização de  $V_{(1)}, \dots, V_{(k)}$ , obtém-se os protótipos de colunas, relativo à  $V_{(p)}, \forall p$ , e verifica-se quais colunas de  $X$  mais se aproximam de cada um dos protótipos.

O implementação desse processo é apresentado no algoritmo 6). Considere  $t$  o contador de iterações,  $U^{(t)}$ ,  $S^{(t)}$  e  $V_{(p)}^{(t)}$  são as matrizes  $U$ ,  $S$  e  $V_{(p)}$ , na iteração  $t$ , respectivamente,  $\mathcal{U}(0, 1) \in ]0, 1]$  uma função que gera valores de uma distribuição uniforme que ignora zeros, e  $\|\cdot\|^2$  é a norma frobenius para vetores.

Note que nesse tipo de fatoração, semelhante à apresentada no problema 3, não necessita de uma fase de pós-processamento para obter o particionamento de linhas e colunas, pois o mesmo é direto a partir das matrizes  $U$  e  $V_{(1)}, \dots, V_{(k)}$ . Ainda, usa-se a mesma estratégia de convergência que a apresentada no algoritmo 5.

### 4.3 Considerações Finais

A motivação para a apresentação dos novos problemas de fatoração de matrizes (OvNMTF e BinOvNMTF), era a possibilidade de superar algumas das dificuldades apresentadas nas fatorações da literatura (ONMTF e FNMTF), no que diz respeito à solução do problema de coagrupamento.

Do ponto de vista de quantização do espaço dos dados, a quantidade de informação armazenada nas novas fatorações tem o potencial de superar as fatorações da literatura.

---

**Algoritmo 6** Algoritmo iterativo para solução do BinOvNMTF
 

---

```

1: function BINOVNMTF( $X, t_{max}, k, l$ )
2:   Initialize:  $U^{(0)} \geq 0, 1 \mid \sum_{p=1}^k \mathbf{u}_{p\cdot} = 1, V_{(p)}^{(0)} \geq 0, 1 \mid \sum_{q=1}^l \mathbf{v}_{(p)q\cdot} = 1, \forall p, S^{(0)} \leftarrow$ 
    $\mathcal{U}(0, 1)$  e  $t \leftarrow 0$ .
3:   while (não convergiu) ou ( $t \leq t_{max}$ ) do
4:
       
$$S^{(t+1)} \leftarrow \sum_{p=1}^k I_{(p)} (U^{(t)T} U^{(t)})^{-1} U^{(t)T} X V_{(p)}^{(t)} (V_{(p)}^{(t)T} V_{(p)}^{(t)})^{-1} \quad (13)$$

5:
       
$$\tilde{V} \leftarrow \sum_{p=1}^k I_{(p)} S^{(t+1)} V_{(p)}^{(t)T}$$

6:
       
$$(U^{(t+1)})_{ip} \leftarrow \begin{cases} 1 & p = \arg \min_{p' \in \{1, \dots, k\}} \|\mathbf{x}_{i\cdot} - \tilde{\mathbf{v}}_{p'\cdot}\|^2 \\ 0 & \text{caso contrário} \end{cases} \quad \forall i, p \quad (14)$$

7:
       
$$\tilde{U}_{(p)} \leftarrow U^{(t+1)} I_{(p)} S^{(t+1)}, \forall p$$

8:
       
$$(V_{(p)}^{(t+1)})_{jq} \leftarrow \begin{cases} 1 & q = \arg \min_{q' \in \{1, \dots, l\}} \|\mathbf{x}_{\cdot j} - \tilde{\mathbf{u}}_{\cdot q'}\|^2 \\ 0 & \text{caso contrário} \end{cases} \quad \forall j, q, p \quad (15)$$

9:      $t \leftarrow t + 1$ 
10:  end while
11:  return  $U^{(t)}, S^{(t)}, V_{(1)}^{(t)}, \dots, V_{(k)}^{(t)}$ 
12: end function

```

---

Isso ocorre porque nas fatorações propostas, o número de colunas  $l$  necessários para explicar a matriz de dados original deve ser mais próxima daquele necessário para obter o número de grupos de colunas desejado, mesmo para os casos em que há interseção de colunas entre as diferentes bases dos grupos de linhas. Entretanto, as fatorações propostas necessitam criar  $k$  matrizes para a abstração dos grupos de colunas, o que incorre em um custo maior de armazenamento dos protótipos dos grupos de colunas, e portanto, menor capacidade de compactação.

Do ponto de vista de geração de informação sobre os dados, as fatorações propostas são capazes de fornecer o mesmo tipo de informação: como as linhas da matriz de entrada se organizam em grupos e como as colunas da matriz de entrada se organizam em grupo. Porém, desde que se tem  $k$  organizações de grupos de colunas, sabe-se que há várias possibilidades dessa organização ocorrer para cada um dos grupos de linhas. Transferindo essa informação para um contexto de aplicação, entende-se que há um conjunto de organização de atributos que estão associados à organização assumida pelos dados no

espaço dos dados. Intuitivamente, pode-se dizer que há diferentes maneiras de justificar o agrupamento de dados descoberto, com base nas similaridades parciais dos atributos que os descrevem.

Do ponto de vista do processo de descoberta dos cogrupos, as fatorações propostas seguem o mesmo princípio seguido nas fatorações da literatura, qual seja, considerar a informação de similaridade de dados e de atributos simultaneamente para resolver o problema de minimização do erro de quantização da matriz original, e consequentemente apresentar protótipos que expliquem os cogrupos. Porém, as fatorações propostas liberam os algoritmos de minimização da necessidade de considerar todos os grupos de linhas na otimização dos grupos de colunas. Desta forma, implementa-se um processo no qual não existe mais a interdependência entre grupos de linhas. Esse fato é que, na realidade, possibilita aproximar a quantidade de cogrupos utilizada para explicar os dados daquela que é a desejada, ainda que a sobreposição de colunas ocorra nos dados.

Essas considerações, assim como aquelas delineadas nos capítulos 2 e 3, são ilustradas no próximo capítulo, no qual resultados experimentais são apresentados.

## 5 Experimentos e Resultados

Para fins de validação dos algoritmos propostos foram realizados experimentos utilizando bases de dados sintéticas e bases de dados textuais reais, sendo que sobre as bases de dados reais, foram criadas versões simplificadas para permitir análises mais detalhadas.

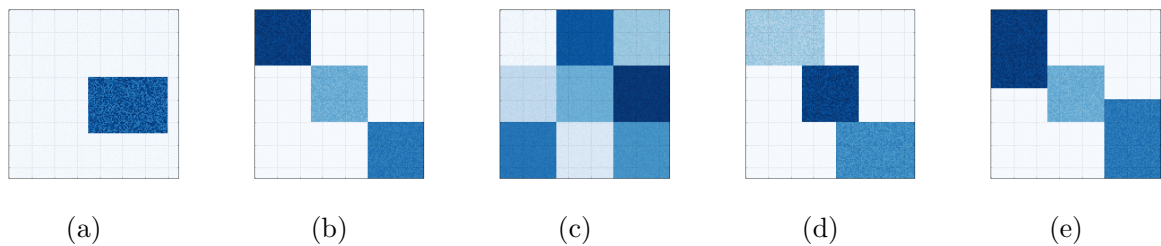
Esses experimentos foram projetados e executados com o fim de ilustrar as capacidades e limitações dos algoritmos de coagrupamento baseados em fatoração de matrizes presentes na literatura e dos algoritmos propostos neste trabalho, todos já apresentados nos capítulos 3 e 4, respectivamente. Tais capacidades e limitações são discutidas neste capítulo em termos de resultados obtidos sobre ambientes controlados (bases de dados sintéticas), ambientes semi-controlados (bases de dados textuais simplificadas) e ambientes aqui denominados não controlados (bases de dados textuais originais). O intuito com a experimentação desses algoritmos em bases de dados textuais é ilustrar seu desempenho como um método de resolução da tarefa de agrupamento de textos, da área de mineração de textos. Também, como os algoritmos clássicos de agrupamento, *k-means* e *fuzzy k-means*, são equivalentes à fatoração de matrizes (veja mais em (DING; HE; SIMON, 2005)), foram aplicados sobre as mesmas bases de dados, para servir como base de comparação dos algoritmos *ONMTF* (algoritmo 3 de Yoo e Choi (2010)), *FNMTF*, *OvNMTF* e *BinOvNMTF*.

Para proporcionar uma visão organizada das capacidades e limitações dos algoritmos, primeiramente são apresentados os experimentos e os resultados obtidos com as bases de dados sintéticas. Tais resultados são apresentados em termos de qualidade de *reconstrução*, a qual é analisada com apoio de visualização gráfica, e de capacidade de *quantização*, a qual é avaliada em termos do erro de quantização dos algoritmos. Então, a qualidade de particionamento é validada nos resultados obtidos com as bases de dados textuais reais originais e simplificadas, fazendo uso de medidas de qualidade de agrupamento. Para essas últimas, uma análise qualitativa é delineada de forma a ilustrar o valor agregado que a flexibilidade de modelos de coagrupamento pode trazer ao contexto de mineração de textos.

## 5.1 Experimentos com Bases de Dados sintéticas

As bases de dados sintéticas foram criadas com inspiração nas diferentes estruturas de cogrupos, propostas por [Madeira e Oliveira \(2004\)](#), descritas com maior detalhamento no capítulo 2. Dentre as nove possíveis estruturas de cogrupos apresentadas por esses autores, foram escolhidas para uso nesses experimentos as três mais simples que, comprovadamente na literatura, os algoritmos *ONMTF* e *FNMTF* são capazes de tratar, e duas estruturas que apresentam interseção parcial de linhas ou colunas, que representam os casos que tais algoritmos não são capazes de tratar de forma natural, ou seja, usando os parâmetros de  $k$  e  $l$ , número de grupos de linhas e número de grupos de colunas, respectivamente, quando escolhidos manualmente analisando as matrizes da figura 6. Portanto, os casos que apresentam interseção parcial de colunas são o alvo dos algoritmos propostos neste trabalho (*OvNMTF* e *BinOvNMTF*). Na figura 6 são apresentadas as visualizações gráficas de cada uma das estruturas de cogrupos em estudo.

Figura 6 – Dados sintéticos gerados a partir das diferentes estruturas de cogrupos. (a) Um único cogrupo. (b) Cogrupos com linhas e colunas sem intersecção. (c) Cogrupos com estrutura em xadrez. (d) Cogrupos sem intersecção nas linhas e com intersecção nas colunas. (e) Cogrupos com intersecção nas linhas e sem intersecção nas colunas.



Fonte: Lucas Fernandes Brunialti, 2016

Para compreender a representação gráfica apresentada na figura 6 considere que cada um dos cinco quadrados maiores representa uma base de dados sintética, que a quantidade de dados da base está representada pela altura do quadrado, que a quantidade de atributos na base está representada pela largura do quadrado. Todas as bases de dados possuem 150 dados (linhas) e 150 atributos (colunas). Considere ainda que cada quadrado ou retângulo, representados com diferentes tons da cor azul, representam um cogrupo, também com suas alturas e larguras representando, respectivamente, a quantidade de linhas e colunas que compõem cada cogrupo. As diferentes tonalidades da cor azul revelam

a similaridade entre os valores assumidos pelos dados em subconjuntos de atributos, o que também revela a existência dos cogrupos. A intensidade da cor é proporcional à intensidade dos valores associados a cada atributo em cada dado, então valores maiores são representados por tonalidades mais escuras e vice-versa. Por exemplo, a seguinte matriz poderia representar o caso da figura 6b, porém em dimensões menores que 150 linhas por 150 colunas:

$$\begin{bmatrix} \mathbf{30,3} & \mathbf{31,6} & \mathbf{30,9} & 1,0 & 0,8 & 0,9 & 0,7 & 0,7 & 0,1 \\ \mathbf{29,1} & \mathbf{30,7} & \mathbf{30,0} & 0,7 & 0,0 & 0,6 & 0,8 & 0,1 & 0,9 \\ \mathbf{30,8} & \mathbf{29,5} & \mathbf{31,5} & 0,2 & 0,7 & 0,2 & 0,9 & 0,7 & 0,5 \\ 0,4 & 0,9 & 1,0 & \mathbf{10,5} & \mathbf{9,2} & \mathbf{10,8} & 0,8 & 0,8 & 0,8 \\ 0,5 & 0,7 & 0,5 & \mathbf{11,1} & \mathbf{10,0} & \mathbf{9,2} & 0,9 & 0,7 & 0,6 \\ 0,3 & 0,4 & 0,5 & \mathbf{10,8} & \mathbf{11,2} & \mathbf{10,9} & 0,5 & 0,3 & 0,1 \\ 0,0 & 0,7 & 0,4 & 0,4 & 0,1 & 0,4 & \mathbf{20,2} & \mathbf{19,6} & \mathbf{20,4} \\ 0,8 & 0,5 & 1,0 & 0,4 & 0,7 & 0,3 & \mathbf{21,2} & \mathbf{20,7} & \mathbf{19,4} \\ 0,0 & 0,6 & 0,4 & 0,6 & 0,1 & 0,1 & \mathbf{19,9} & \mathbf{20,2} & \mathbf{20,9} \end{bmatrix}$$

Todas as bases de dados da figura 6 são matrizes de valores reais positivos geradas artificialmente. Primeiramente, uma matriz de tamanho  $150 \times 150$  é preenchida com valores ponto flutuante, gerados aleatoriamente a partir de uma função  $\mathcal{U}(0, 1)$ , que gera números de uma distribuição uniforme no intervalo  $]0, 1]$ . Em seguida, o tamanho em termo de linhas e colunas e disposição dos cogrupos na base de dados foram determinados de acordo com cada estrutura de cogrupos desejada. Para instanciar cada cogrupos, um conjunto de valores foi criado e distribuído pelas células do cogrupos da seguinte forma:

- um valor central  $c \in \mathcal{C}$ , sendo  $\mathcal{C} = \{20, 40, 60, 80, 100, 120, 140, 160, 180\}$ , foi aleatoriamente escolhido, e retirado de  $\mathcal{C}$  para não haver cogrupos com o mesmo  $c$  em uma mesma matriz;
- o conjunto de valores usado para instanciar as células do cogrupos foi estabelecido por meio da adição de  $c$  a valores reais, gerados a partir da função  $\mathcal{U}(0, 10)$ , que gera números de uma distribuição uniforme no intervalo  $]0, 10]$ ;
- cada um dos valores nesse conjunto foi atribuído às células previamente definidas como pertencentes ao cogrupos.

Assim, considerando que um cogruppo é uma submatriz da matriz original  $X$ , ele pode ser gerado pela equação:

$$x_{ij} = c + \mathcal{U}(0, 10)$$

sendo  $i$  e  $j$  os índices das linhas e colunas de  $X$  escolhidos para compor o cogruppo.

Ainda sobre a interpretação da figura 6, alguns detalhes precisam ser observados. A depender do objetivo da análise de coagrupamento, na figura 6a, por exemplo, mais de um cogruppo pode ser observado, além daquele que é de interesse de análise neste trabalho. Para essa interpretação, considera-se que todo agrupamento de linhas e de colunas, independente de serem úteis ou não, ou de serem interesse para análise ou não, é um cogruppo. Assim, tem-se os seguintes cogruppos na base de dados sintética (a):

- O mais evidente, destacado em azul, é formado pelas linhas [60 .. 109] e pelas colunas [70 .. 139]. Esse é o cogruppo de interesse, nesse projeto, para a resolução da tarefa de coagrupamento aplicada a dados textuais, e é representado na figura 6a pelo quadrado em cor azul.
- O segundo, que não está destacado na visualização, é formado por todas as linhas e as colunas [0 .. 69] e [140 .. 149].
- O terceiro, também não destacado, é formado por todas as colunas e as linhas [0 .. 59] e [110 .. 149].

Sob essa ótica de interpretação, as demais bases de dados possuem:

- (b): três cogruppos de principal interesse e seis cogruppos não destacados na figura;
- (c): seis cogruppos de principal interesse;
- (d): três cogruppos de principal interesse e oito cogruppos não destacados na figura;
- (e): três cogruppos de principal interesse e oito cogruppos não destacados na figura.

Para cada uma das bases de dados sintéticas foram executados os seguintes algoritmos: *k-means*<sup>1</sup>, *fuzzy k-means*<sup>2</sup>, *ONMTF*, *FNMTF*, *OvNMTF* e *BinOvNMTF*<sup>3</sup>. Os resultados foram analisados em termos de qualidade de reconstrução e capacidade de quantização.

<sup>1</sup> Para os experimentos com *k-means* foi usada a implementação da biblioteca *scikit-learn* (PEDREGOSA et al., 2011) da linguagem Python.

<sup>2</sup> Para experimentos com *fuzzy k-means* foi usado a implementação do algoritmo *fuzzy k-means* da biblioteca *scikit-fuzzy* da linguagem Python.

<sup>3</sup> Os algoritmos baseadas em fatoração de matrizes foram implementados pelo autor deste trabalho usando a linguagem Python.



### 5.1.1 Análise da reconstrução

Uma forma de analisar o resultado dos algoritmos estudados neste trabalho é analisá-los quanto a sua capacidade de reconstrução. A reconstrução é feita tomando as matrizes resultantes da fatoração e combinando-as da mesma forma que o seu problema foi proposto, de forma a reconstruir a matriz original.

Essa análise se torna importante quando se tem como objetivo avaliar o comportamento dos algoritmos em diferentes estruturas de organização de cogrupos, com destaque para a análise de maior interesse neste trabalho – coagrupamento com interseção parcial de linhas ou colunas. [A capacidade de reconstrução .... vamos ver se colocamos alguma coisa a mais para justificar esse análise.](#)

Um resumo sobre os resultados obtidos nessa análise é apresentado na tabela 1. O restante dessa seção se destina a detalhar as análise de qualidade de reconstrução.

Tabela 1 – Resumo de qualidade de reconstrução: *ok* - permite reconstrução de forma natural;  $\times$  - sem informação sobre interseção parcial; + - preserva informação de interseção parcial

	base (a)	base (b)	base (c)	base (d)	base (e)
<i>k-means</i>	<i>ok</i>	<i>ok</i>	<i>ok</i>	<i>ok</i> , $\times$	
<i>fuzzy k-means</i>	<i>ok</i>	<i>ok</i>	<i>ok</i>	<i>ok</i> , $\times$	+
<i>ONMTF</i>	<i>ok</i>	<i>ok</i>	<i>ok</i>	+	+
<i>FNMTF</i>	<i>ok</i>	<i>ok</i>	<i>ok</i>		
<i>OvNMTF</i>	<i>ok</i>	<i>ok</i>	<i>ok</i>	<i>ok</i> , +	<i>ok</i> , +
<i>BinOvNMTF</i>	<i>ok</i>	<i>ok</i>	<i>ok</i>	<i>ok</i> , +	

Fonte: Lucas Fernandes Brunialti, 2016

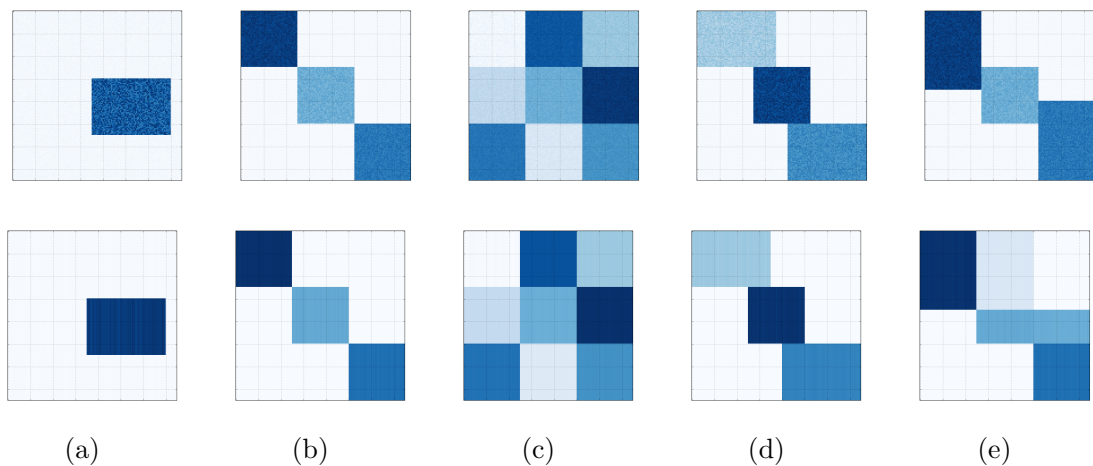
### 5.1.2 Reconstrução a partir dos resultados do algoritmo *k-means*

Para execução dos experimentos com o *k-means* os seguintes parâmetros foram estabelecidos: número máximo de iterações (300) ou a diferença do erro de quantização em duas iterações consecutivas ( $\leq 1 \times 10^{-4}$ ), o que ocorrer primeiro, como critérios de parada; número de grupos  $k = 2$  para a base de dados sintética (a) e  $k = 3$  para as bases de dados sintéticas (b), (c), (d) e (e). A escolha de  $k$  foi realizada com a prerrogativa de que o algoritmo *k-means* deveria encontrar os grupos considerando todos os atributos descritivos, e desta forma, agrupar os dados de acordo com a similaridade total. Assim,  $k$

foi escolhido a partir da quantidade de agrupamento de linhas presente na base de dados. A escolha de  $k$  maiores levaria o algoritmo a, necessariamente, dividir em grupos diferentes os dados que deveriam pertencer a um mesmo agrupamento de linhas. Foram executadas 10 rodadas do algoritmo, com inicialização de centróides aleatória, sendo que o melhor modelo resultante nessas rodadas foi escolhido para ilustrar a avaliação da qualidade da reconstrução.

As reconstruções obtidas a partir dos centróides encontrados pelo algoritmo *k-means* são ilustradas na figura 9. As matrizes originais são repetidas na figura, nas cinco primeiras posições, de forma a facilitar a análise visual dos resultados. Para um melhor entendimento da representação gráfica, note que cada linha recebe cores de acordo com sua pertinência a um agrupamento de linhas (ou grupo, na visão de agrupamento) representado por um centróide. Ou seja, se a linha 10 pertencer ao centróide 1, então os valores da linha 10 serão substituídos pelos valores do centróide 1. Note que o centróide é um vetor com o mesmo número de coordenadas de um dado da base de dados, sendo assim, a substituição é direta. O centróide que representa o agrupamento de linhas referente ao cogrupo de interesse (em azul na figura), possui, claramente, valores similares aos dados que pertencem a esse cogrupo, e por isso o procedimento proposto pode ser visto como uma representação de reconstrução.

Figura 7 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo *k-means*.



Fonte: Lucas Fernandes Brunialti, 2016

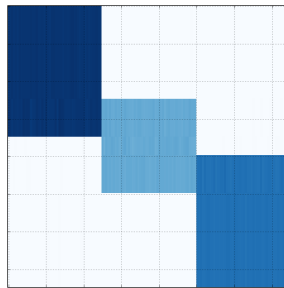
O modelo resultante da execução do *k-means* permite representar perfeitamente a reconstrução para as bases de dados (a) à (d). Porém, é preciso lembrar que não

há informação sobre agrupamento de colunas no resultado do algoritmo, sendo que a visualização gráfica de cada coagrupamento é possível apenas por conta do algoritmo levar em conta todos os atributos para realização do agrupamento.

O modelo não permite a descoberta de grupos com interseção de linhas (um dado não pode pertencer a mais de um grupo), e portanto não permite uma boa representação para a reconstrução no caso da base de dados (e). Esse resultado já era esperado devido à natureza da solução apresentada pelo *k-means*.

No entanto, se mudar o valor de  $k$  para descobrir 5 grupos ao invés dos 3 grupos desejados, é possível realizar a reconstrução da base de dados (e), como mostra a figura 8. Porém, não é uma reconstrução natural, ou seja, a informação que seria desejado é que há 3 grupos de linhas com interseção parcial entre o primeiro e o segundo grupos, e interseção parcial entre o segundo e terceiro grupos, diferentemente do obtido, que existe 5 grupos diferentes.

Figura 8 – Resultado da reconstrução da base de dados (e) utilizando o algoritmo *k-means* com  $k = 5$ .



Fonte: Lucas Fernandes Brunialti, 2016

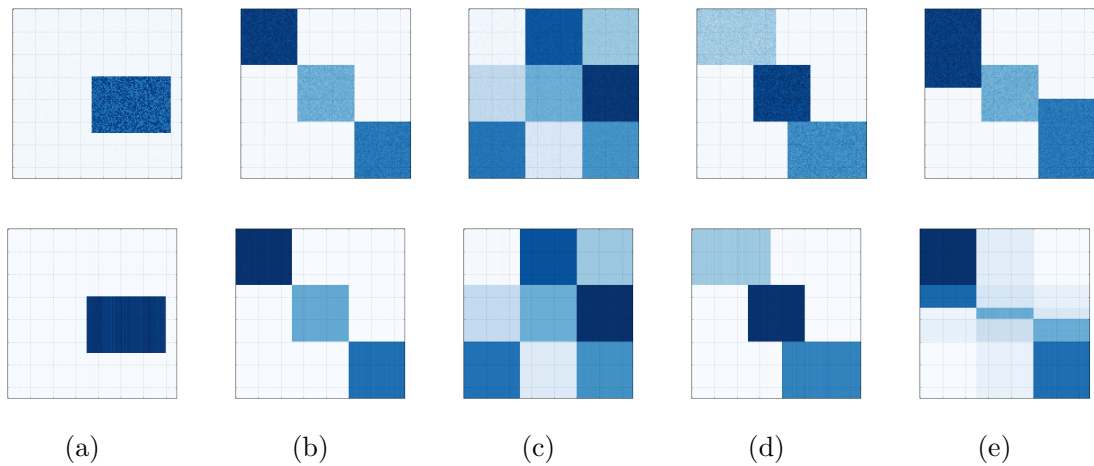
### 5.1.3 Reconstrução a partir dos resultados do algoritmo *fuzzy k-means*

Os mesmos critérios de parada e número de grupos usados para o *k-means* foram também usados nos experimentos com o *fuzzy k-means*. O valor para o parâmetro de fuzzificação  $m$  foi mantido em 2, como indicado em Rocha et al. (). Também, 10 execuções foram realizadas, com iniciação aleatória de pesos e o melhor modelo obtido foi usado para avaliação da qualidade de reconstrução obtida.

As reconstruções apresentadas na figura 9 foram obtidas por meio da multiplicação da matriz de partição pela matriz dos protótipos  $UC$ , resultantes da execução do algoritmo, e as matrizes originais foram repetidas para facilitar a análise visual. Note que a matriz  $U$

do *fuzzy k-means* tem justamente o mesmo papel que a matriz de pertinências de linhas  $U$  das fatorações *BVD*, *ONMTF* e *OvNMTF*.

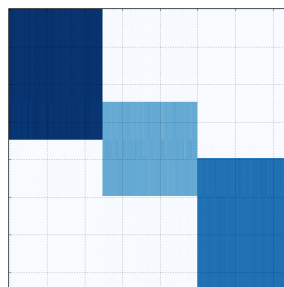
Figura 9 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo *fuzzy k-means*.



Fonte: Lucas Fernandes Brunialti, 2016

De forma semelhante à análise de reconstrução proporcionada pelo *k-means*, o *fuzzy k-means* também permite representar perfeitamente a reconstrução para as bases de dados (a) à (d). Porém, a reconstrução para o caso (e) não foi obtida com sucesso. No entanto, o algoritmo preservou parte da interseção, isso pode ser visto pelas regiões sombreadas onde há interseção de grupos.

Figura 10 – Resultado da reconstrução da base de dados (e) utilizando o algoritmo *fuzzy k-means* com  $k = 5$ .



Fonte: Lucas Fernandes Brunialti, 2016

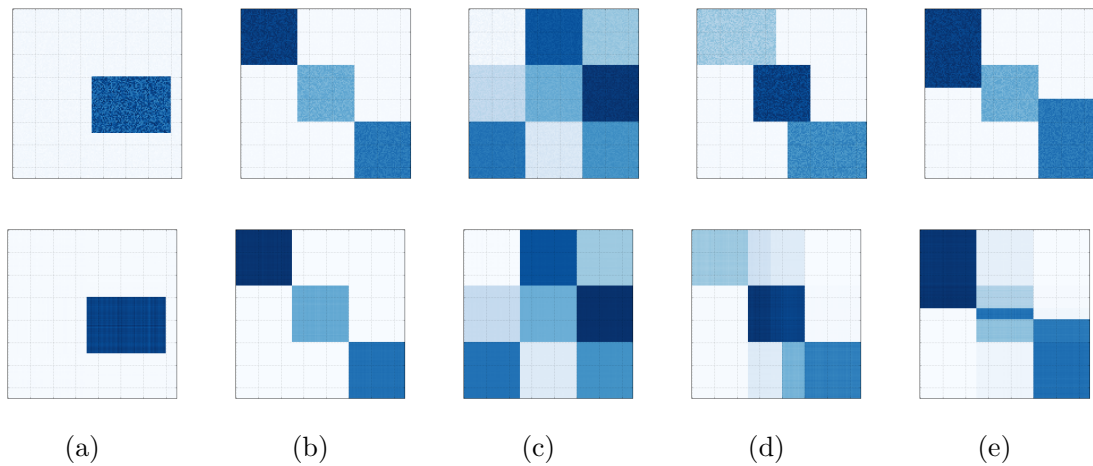
Da mesma forma que o *k-means*, é possível mudar o valor de  $k$  para 5, para assim, obter uma reconstrução perfeita (figura 10). Fazendo as mesmas ressalvas, que a informação obtida a partir da interpretação desse modelo, não é a desejada.

### 5.1.4 Reconstrução a partir dos resultados do algoritmo *ONMTF*

Para execução dos experimentos com o *ONMTF* a seguinte parametrização foi estabelecida: número máximo de iterações (1000) ou a diferença do erro de quantização em duas iterações consecutivas ( $\leq 1 \times 10^{-4}$ ), o que ocorrer primeiro, como critérios de parada; o número de cogrupos de linhas ( $k$ ) e colunas ( $l$ ) configurados da seguinte maneira:  $k = l = 2$  para (a) e  $k = l = 3$  (b), (c), (d) e (e). Novamente, as escolhas são baseadas no conhecimento *a priori* que se tem sobre a quantidade de cogrupos, e quais cogrupos, deseja-se obter. Para a visualização da reconstrução do algoritmo *ONMTF* mantendo os valores (e portanto as cores) das matrizes de dados, não foi realizada a normalização baseada em probabilidade proposta por Yoo e Choi (2010) (apresentada no algoritmo 3), pois esta muda a escala dos valores da reconstrução, dada a interpretação probabilística que pode ser feita.

A partir da realização de 10 execuções do algoritmo, foi escolhido o modelo que alcançou o menor erro de quantização e a partir do resultado obtido a reconstrução foi avaliada. A qualidade das resconstruções pode ser visualmente observada na figura 11. A resconstrução foi obtida a partir da multiplicação das matrizes fatoradas, ou seja,  $USV^T$ , conforme explicado no capítulo 3.

Figura 11 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas resconstruções, realizadas a partir dos resultados obtidos com o algoritmo *ONMTF*.



Fonte: Lucas Fernandes Brunialti, 2016

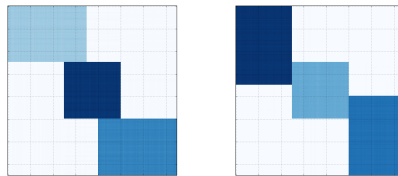
Analisando os resultados, é possível perceber que a reconstrução é realizada com êxito nos casos (a), (b) e (c). O algoritmo falhou em reconstruir a matriz no caso (d)

pois não foi capaz de associar corretamente as colunas que pertencem a mais de um agrupamento de colunas (interseção parcial nas colunas). O mesmo efeito ocorre com o caso (e), no qual há interseção parcial de linhas. A falha da reconstrução é percebida na coloração diferenciada, formando regiões com sombras, nas colunas, ou linhas, envolvidas nas interseções.

A coloração diferenciada daquela esperada (seguindo a matriz original) é resultante do estado da matriz  $V$ , que contém as relações de associação de linhas e/ou colunas aos cogrupos. Para o caso (d), por exemplo, é possível deduzir, a partir da análise visual, que há valores de magnitude diferentes estabelecendo a associação das colunas 50 à 79 e 80 à 99, que pertencem aos primeiro e terceiro cogrupos, respectivamente, com o cogrupos de colunas do meio. Ou seja, a reconstrução preservou parte da interseção. Esse comportamento é observado pois o algoritmo força a ortogonalidade entre cogrupos de linhas e cogrupos de colunas.

Ainda, é possível fazer a reconstrução perfeita, se mudar o valor de  $l$  para 5 na base de dados (d), e o valor de  $k$  para 5 na base de dados (e) (figura 12). Porém, como nos outros casos, a interpretação do modelo gerado não é seria a desejada.

Figura 12 – Resultado da reconstrução da base de dados (d) com  $k = 5$  e (e) com  $l = 5$ , respectivamente, utilizando o algoritmo *ONMTF*.



Fonte: Lucas Fernandes Brunialti, 2016

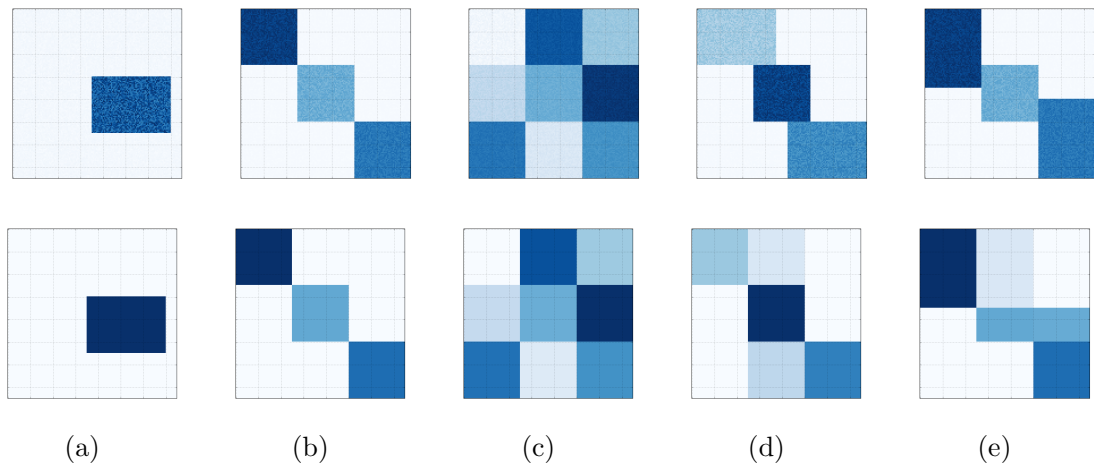
### 5.1.5 Reconstrução a partir dos resultados do algoritmo *FNMTF*

Para execução dos experimentos com o *FNMTF* a seguinte parametrização foi estabelecida: número máximo de iterações (300) ou a diferença do erro de minimização em duas iterações consecutivas ( $\leq 1 * 10^{-4}$ ), o que ocorrer primeiro, como critérios de parada; o número de cogrupos de linhas ( $k$ ) e colunas ( $l$ ) configurados da seguinte maneira:  $k = l = 2$  para (a),  $k = l = 3$  para (b), (c), (d) e (e).

A partir da realização de 10 execuções do algoritmo, foi escolhido o modelo que alcançou o menor erro de quantização e a partir do resultado obtido a reconstrução foi

avaliada. A qualidade das reconstruções pode ser visualmente observada na figura 13. A reconstrução foi obtida a partir da multiplicação das matrizes fatoradas, ou seja,  $USV^T$ , conforme explicado no capítulo 3.

Figura 13 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo *FNMTF*.

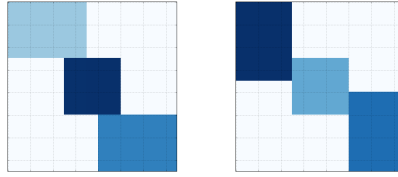


Fonte: Lucas Fernandes Brunialti, 2016

Este caso é semelhante ao anterior (algoritmo *ONMTF*). O algoritmo permitiu a reconstrução com êxito dos casos (a), (b) e (c), falhando na reconstrução dos casos (d) e (e), nos quais há intersecção de colunas ou linhas nos cogrupos de interesse. No entanto, nesse caso o algoritmo restringe a associação de linhas (ou colunas), a apenas um agrupamento de linhas (ou de colunas), e por isso a informação referente a intersecção em cogrupos é totalmente descaracterizada. Observe, por exemplo, que a reconstrução no caso (d) se assemelha à reconstrução do caso (c), embora as matrizes originais, em cada caso, representem informação sobre associação de dados e atributos em cogrupos também diferenciada.

Ainda, é possível fazer a reconstrução perfeita, se mudar o valor de  $l$  para 5 na base de dados (d), e o valor de  $k$  para 5 na base de dados (e) (figura 12). Porém, como nos outros casos, a interpretação do modelo gerado não é seria a desejada.

Figura 14 – Resultado da reconstrução da base de dados (d) com  $k = 5$  e (e) com  $l = 5$ , respectivamente, utilizando o algoritmo *FNMTF*.



Fonte: Lucas Fernandes Brunialti, 2016

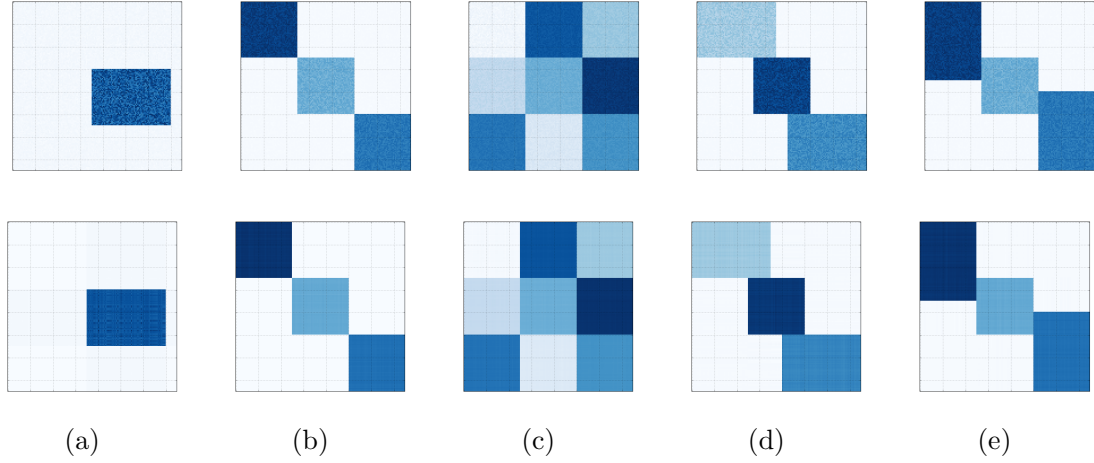
### 5.1.6 Reconstrução a partir dos resultados do algoritmo *OvNMTF*

Para execução dos experimentos com o *OvNMTF* a seguinte parametrização foi estabelecida: número máximo de iterações (1000) ou a diferença do erro de quantização ( $\leq 1 * 10^{-4}$ ), o que ocorrer primeiro, como critérios de parada; o número de cogrupo de linhas ( $k$ ) e colunas ( $l$ ) configurados da seguinte maneira:  $k = l = 2$  para (a),  $k = 3$  e  $l = 2$  para (b), (d) e (d), e  $k = l = 3$  para (c). Note que o parâmetro  $l$  mudou em relação às outras reconstruções, pois no caso dos algoritmos propostos, para escolher o parâmetro  $l$ , é necessário responder a pergunta: quantos cogrupo de coluna existem para cada cogrupo de linhas?

A partir da realização de 10 execuções do algoritmo, foi escolhido o modelo que alcançou o menor erro de quantização e a partir do resultado obtido a reconstrução foi avaliada. A qualidade das reconstruções pode ser visualmente observada na figura 15. A reconstrução foi obtida a partir da multiplicação das matrizes fatoradas, ou seja,  $U \sum_{p=1}^k I_{(p)} S V_{(p)}^T$ , conforme explicado no capítulo 4.



Figura 15 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo *OvNMTF*.



Fonte: Lucas Fernandes Brunialti, 2016

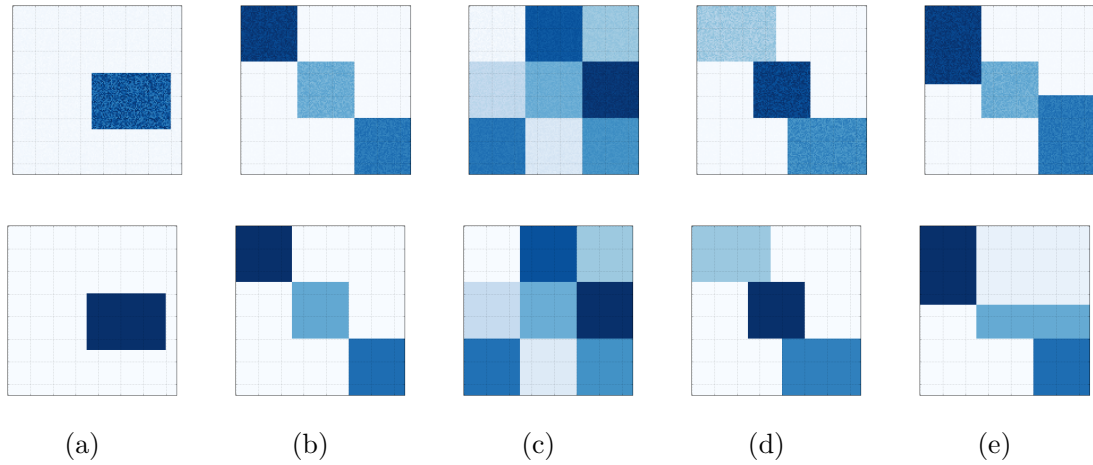
O algoritmo conseguiu realizar a reconstrução com êxito em todos os casos. Note que o algoritmo *OvNMTF* é capaz de lidar com cogrupos com intersecção de colunas, e portanto é capaz de resolver o caso (d). Ainda, como não há restrições de ortogonalidade e existe fatores para controlar a pertinência de uma mesma linha ou coluna à múltiplos cogrupos, o algoritmo é capaz de reconstruir perfeitamente a base de dados (e).

### 5.1.7 Reconstrução a partir dos resultados do algoritmo *BinOvNMTF*

Para execução dos experimentos com o *BinOvNMTF* a seguinte parametrização foi estabelecida: número máximo de iterações (300) ou a diferença do erro de quantização ( $\leq 1 * 10^{-4}$ ), o que ocorrer primeiro, como critérios de parada; o número de cogrupos de linhas ( $k$ ) e colunas ( $l$ ) configurados da mesma maneira que a reconstrução com o algoritmo *OvNMTF*:  $k = l = 2$  para (a),  $k = 3$  e  $l = 2$  para (b), (d) e (d), e  $k = l = 3$  para (c).

A partir da realização de 10 execuções do algoritmo, foi escolhido o modelo que alcançou o menor erro de quantização e a partir do resultado obtido a reconstrução foi avaliada. A qualidade das reconstruções pode ser visualmente observada na figura 16. A reconstrução foi obtida a partir da multiplicação das matrizes fatoradas, ou seja,  $U \sum_{p=1}^k I_{(p)} S V_{(p)}^T$ , conforme explicado no capítulo 4.

Figura 16 – As primeiras cinco matrizes são as matrizes originais, as demais são suas respectivas reconstruções, realizadas a partir dos resultados obtidos com o algoritmo *BinOvNMTF*.

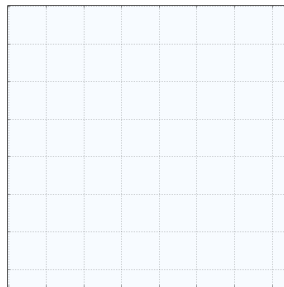


Fonte: Lucas Fernandes Brunialti, 2016

O algoritmo *BinOvNMTF* conseguiu realizar a reconstrução com êxito dos casos (a), (b), (c) e (d). Como o algoritmo é capaz de lidar com cogrupos com intersecção de colunas, ele é capaz de resolver o caso (d). No entanto, o algoritmo não é capaz de lidar com cogrupos com intersecção nas linhas, então, não é capaz de realizar a reconstrução do caso (e). E nesse caso, também não preserva qualquer informação sobre a intersecção, já que possui a restrição de associação binária (cada linha/coluna associadas a um único agrupamento de linhas/colunas).

Fazendo as mesmas ressalvas das reconstruções nos casos com o algoritmos *kmeans*, *fuzzy kmeans*, *ONMTF* e *FNMTF*, é possível fazer a reconstrução perfeita utilizando  $k = 5$  para a base de dados (e) (figura 17).

Figura 17 – Resultado da reconstrução da base de dados (e) utilizando o algoritmo *BinOvNMTF* com  $k = 5$ .



Fonte: Lucas Fernandes Brunialti, 2016

### 5.1.8 Análise da Capacidade de Quantização

Nos experimentos com as bases de dados sintéticas, fazendo o uso da mesma parametrização descrita para a análise da reconstrução, foi feita a medida do erro de quantização

A tabela 2 sumariza o melhor erro de quantização resultante das 10 execuções de cada algoritmo para cada base de dados.

Tabela 2 – Avaliação da capacidade de quantização segundo o erro de quantização com os melhores destacados em negrito.

	base (a)	base (b)	base (c)	base (d)	base (e)
<i>k-means</i>	30.336, 0	62.776, 5	184.992, 8	<b>79.238, 5</b>	2.886.245, 5
<i>fuzzy k-means</i>	29.402, 8	63.768, 5	183.991, 4	<b>78.340, 0</b>	2.307.168, 6
<i>BVD</i>	29.260, 7	61.536, 5	179.804, 0	<b>75.575, 0</b>	<b>77.135, 5</b>
<i>ONMTF</i>	30.555, 4	60.794, 8	184.255, 9	579.136, 7	781.131, 6
<i>FNMTF</i>	30.924, 7	64.636, 1	186.224, 4	1.634.328, 5	2.881.172, 0
<i>OvNMTF</i>	30.439, 8	61.863, 2	178.886, 8	<b>75.533, 6</b>	<b>75.931, 2</b>
<i>BinOvNMTF</i>	31.239, 0	63.660, 4	187.579, 8	<b>79.968, 0</b>	3.160.391, 0

Fonte: Lucas Fernandes Brunialti, 2016

Através dos resultados obtidos é possível perceber que a capacidade de quantização para todos os algoritmos para as bases de dados (a), (b) e (c), são semelhante, assim como a foi percebido na análise da reconstrução.

Para a base de dados (d), que contém cogrupos de colunas organizados com intersecção parcial, os algoritmos que obtiveram erro de quantização entre 75.000 e 80.000 foram destacados na tabela 2, isso indica que eles foram capazes de quantizar a informação de intersecção parcial. É importante perceber que capacidade de quantização é diferente de possibilidade de interpretação da quantização realizada, isto é, os algoritmos *kmeans* e *fuzzy k-means* foram capazes de quantizar a intersecção parcial de cogrupos de colunas, mas isso não quer dizer que é possível interpretar essa informação. Como foi mostrado na análise de reconstrução, o algoritmo *ONMTF* não quantizou toda a informação de intersecção parcial da base de dados (d), mas parte dessa informação foi quantizada, diferentemente do algoritmo *FNMTF*, que perdeu totalmente essa informação. Isso também está evidenciado na diferença do erro de quantização entre os algoritmos *ONMTF* e *FNMTF*.

O único algoritmo que foi capaz de quantizar a intersecção parcial entre cogrupos de linhas da base de dados (e) foi o *OvNMTF*, isso ocorre pois não há restrição de

ortogonalidade entre cogrupos. Da mesma forma que na base de dados (d), o algoritmo *ONMTF* não quantizou toda a informação de intersecção, evidenciado pela diferença do erro de quantização obtido pelo algoritmo *OvNMTF*.

## 5.2 Experimentos com Bases de Dados Reais

Três bases de dados reais foram usadas a fim de ilustrar a aplicação dos algoritmos de coagrupamento no contexto da resolução de uma tarefa de mineração de texto. Dessas bases, uma (*NIPS14-17*) é base de referência, usada pela comunidade científica para experimentação de algoritmos de mineração de textos. Duas das bases de dados, *IG* e sua versão reduzida *IG toy*, foram elaboradas no contexto deste trabalho, e constituem-se como uma das contribuições desta pesquisa, já que podem constituir-se como mais duas bases de referência baseadas em dados reais, em língua portuguesa. Essas bases de dados estão descritas nessa seção, juntamente com o pré-processamento realizado sobre elas, e as análises quantitativas e qualitativas obtidas a partir da aplicação dos algoritmos de coagrupamento sobre elas.

### 5.2.1 Descrição das bases de dados

As três bases de dados reais são compostas por textos referentes à textos acadêmicos e notícias de portais postagem. O contexto referente a cada uma das bases é brevemente apresentados nesta seção, sendo que na tabela 3 são listadas algumas estatísticas para cada um deles. A esparsidade foi calculada da seguinte forma:  $1 - \frac{\# \text{ Total de palavras}}{\# \text{ Documentos} \times \# \text{ Palavras únicas}}$ .

1. ***NIPS14-17***<sup>4</sup> Esta base de dados contém uma coleção de trabalhos acadêmicos publicados no congresso *NIPS* (*Neural Information Processing Systems*) no período de 2001 a 2003, dos volumes 14 a 17. A construção da base de dados *NIPS14-17* foi realizada por *Sam Roweis*<sup>5</sup>, a partir de um processamento aplicado aos dados adquiridos por *Yann LeCun* usando um dispositivo de reconhecimento ótico de caracteres (OCR) ([GLOBERSON et al., 2007](#)). A fonte de dados original, usada na construção desta base, possui os trabalhos científicos publicados em 18 volumes (0 a 17), porém apenas os trabalhos dos volumes 14 a 17 estão rotulados. Tais

<sup>4</sup> <http://robotics.stanford.edu/~gal/data.html>

<sup>5</sup> <http://www.cs.nyu.edu/~roweis/data.html>

documentos estão organizados sob tópicos que compreendem áreas técnicas (teoria de aprendizado, neurociência, algoritmos e arquiteturas e etc), e estão distribuídos de forma desbalanceada entre os grupos caracterizados por cada um desses tópicos, por isso, foram usados documentos das 9 áreas técnicas com mais documentos, das 13 áreas técnicas originalmente disponibilizadas.

2. **IG** Esta base de dados foi criada como uma contribuição deste trabalho, e consiste em uma coleção de notícias extraídas do portal iG<sup>6</sup>. Cada documento nesta base contém o endereço eletrônico no qual a notícia está publicada, título, subtítulo, corpo e canal da notícia, sendo que o canal representa uma classificação para a notícia, atribuída pelos construtores do portal. As notícias que compõem essa base foram publicadas no período de 2 de janeiro de 2012 à 11 de outubro de 2014 e estão distribuídas em 13 canais, de maneira desbalanceada.
3. **IG toy (ou reduzido)** Esta base de dados é um subconjunto do conjunto de dados IG, composto por 300 notícias dispostas de forma balanceada em três canais (*esporte*, *arena* e *jovem*).

Tabela 3 – Estatísticas das bases de dados usadas nos experimentos.

Base de dados	# Palavras únicas	# Total de palavras	# Documentos	# Grupos	Esparsidade
<i>NIPS14-17</i>	6.881	746.826	555	9	0,804
<i>IG</i>	19.563	1.187.334	4.593	13	0,987
<i>IG toy</i>	6.764	70.169	300	3	0,965

Fonte: Lucas Fernandes Brunialti, 2016

A extração das notícias do portal iG foi realizada através da implementação de um *web crawler* utilizando a linguagem Python<sup>7</sup>. As notícias foram capturadas a partir de uma página de início, fornecida para o *web crawler*, que era selecionada a fim de equalizar a distribuição de notícias por ano e por categoria (canal), mostradas nas tabelas 4 e 5. Todas as notícias coletadas para compor a base tem no mínimo 250 caracteres no corpo.

Analisando a distribuição de notícias por ano foi possível verificar que os links escolhidos como partida para o *web crawler* realizar a extração de notícias foram efetivos para deixar a distribuição perto de uniforme, com média e desvio padrão de aproximadamente  $1.531 \pm 337$  notícias. Contrariamente, a distribuição de notícias por canal não ficou perto

<sup>6</sup> <http://ig.com.br/>

<sup>7</sup> <https://www.python.org/>

Tabela 4 – Distribuição de notícias por ano (base de dados *IG*).

Ano	# Notícias
2012	1.551
2013	1.933
2014	1.109

Fonte: Lucas Fernandes Brunialti, 2016

Tabela 5 – Distribuição de notícias por canal (base de dados *IG*).

Canal	# Notícias
<i>economia</i>	907
<i>ultimosegundo</i>	555
<i>igirl</i>	527
<i>jovem</i>	524
<i>arena</i>	421
<i>tecnologia</i>	359
<i>esporte</i>	342
<i>delas</i>	252
<i>igay</i>	210
<i>gente</i>	196
<i>deles</i>	141
<i>saude</i>	88
<i>luxo</i>	71

Fonte: Lucas Fernandes Brunialti, 2016

do uniforme, com média e desvio padrão de aproximadamente  $353 \pm 225$  notícias. Uma hipótese é que isso se deve à idade e popularidade do canal, por exemplo, os canais *luxo* e *deles* são muito mais recentes e menos populares que o *ultimosegundo*.

Para a base de dados *NIPS*, da mesma forma, é mostrada, nas tabelas 6 e 7, as distribuições de trabalhos acadêmicos por ano e trabalhos acadêmicos por áreas técnicas, respectivamente. Note que os documentos rotulados como as seguintes áreas técnicas não foram usados nos experimentos: *Control & Reinforcement Learning*, *Vision (Biological)*, *Brain Imaging* e *Speech and Signal Processing*.

Tabela 6 – Distribuição de trabalhos acadêmicos por ano (base de dados *NIPS*).

Ano	# Notícias
2001	192
2002	204
2003	197

Fonte: Lucas Fernandes Brunialti, 2016

Tabela 7 – Distribuição de trabalhos acadêmicos por áreas técnicas (base de dados *NIPS*).

Áreas técnicas	# Trabalhos acadêmicos
<i>Algorithms &amp; Architectures</i>	209
<i>Cognitive Science &amp; AI</i>	68
<i>Implementations</i>	66
<i>Vision</i>	50
<i>Neuroscience</i>	47
<i>Vision (Machine)</i>	40
<i>Emerging Technologies</i>	33
<i>Applications</i>	26
<i>Learning Theory</i>	16
<i>Control &amp; Reinforcement Learning</i>	15
<i>Vision (Biological)</i>	9
<i>Brain Imaging</i>	7
<i>Speech and Signal Processing</i>	7

Fonte: Lucas Fernandes Brunialti, 2016

### 5.2.2 Pré-processamento

A fim de estruturar a informação nas bases de dados para construção dos modelos, foi necessária uma fase de pré-processamento. As tarefas de pré-processamento comumente executadas em dados textuais são necessárias para melhorar a qualidade dos dados que serão submetidos à análise e também adequar a representação dos textos às necessidades dos algoritmos de análise (no caso deste trabalho, é necessário criar uma representação numérica e vetorial para os textos).

As ações executadas neste trabalho para realizar o pré-processamento dos textos foram:

- tokenização: o objetivo nesta ação é criar um “dicionário de termos” para a coleção de documentos. Para isso, um procedimento de quebra do texto em termos (*tokens* ou palavras) é realizado por meio da determinação de caracteres delimitadores (espaço em branco) e eliminação de caracteres que não se constituem como termos significativos para representação do texto (pontuação, caracteres especiais, etc). A decisão sobre como decidir os caracteres delimitadores e quais símbolos serão considerados insignificantes depende do contexto dos textos. Para a base de dados *IG* e *IG toy* foi usada uma expressão regular que separa caracteres não contíguos: (w+). Essa fase de pré-processamento não foi necessária para a base de dados *NIPS*, já que a base de dados é fornecida com essa etapa.

- filtragem: na filtragem são retirados os termos (*stopwords*) que não contribuem para a descrição ou identificação de um texto. Tradicionalmente, palavras das seguintes classes gramaticais são retiradas por esse filtro: conjunções, artigos, preposições e advérbios. Podem ser adicionadas a essa lista estão outras classes de palavras como numerais, nomes próprios, elementos da web, tokens monetários, e também palavras que aparecem em todos os textos por uma questão de padrão/formato dos mesmos. Algumas palavras foram acrescentadas à lista tradicional de *stopwords* no tratamento da base de dados IG: *leia, lendo, twitter, facebook, mais, divulgação, agnews, ler, ig, reprodução, siga, curta, images, imagens, veja, tambem, também, saiba, informações, thinkstock, getty, photos, foto, fotos, ap, gettyimages, infográfico, aí*. Para a base de dados *NIPS*, não foi necessário a etapa de filtragem de *stopwords*. Ainda, como parte da filtragem, palavras cuja frequência de ocorrência nos documentos é muito pequena ou muito grande, podem ser acrescentadas à lista. Em todas as bases usadas neste trabalho, todas as palavras com ocorrência menor ou igual à dois em todos os documentos, foram retiradas.

A representação vetorial objetivada é composta por uma relação de documentos e termos, organizada em modelo conhecimento como modelo do espaço vetorial, ou *Vector Space Model* (SALTON; WONG; YANG, 1975; SEBASTIANI, 2002; LOPS; GEMMIS; SEMERARO, 2011). Nesta representação, cada documento é representado por um vetor composto por tantas dimensões quanto forem os termos presentes no “dicionário de termos”. Formalmente, há um conjunto de  $n$  documentos  $\{d_1, \dots, d_n\}$ , e um conjunto de  $m$  termos  $\{t_1, \dots, t_m\}$ , e para cada par  $(d_i, t_j)$ , sendo os índices  $i = \{1, \dots, n\}$  e  $j = \{1, \dots, m\}$ , estabelece-se uma relação que expressa a maneira como um termo será usado na representação de um documento.

A relação  $tf(t_j, d_i)$  expressa a frequência de ocorrência do termo  $t_j$  no documento  $d_i$ , e o vetor de representação de um documento, que será uma linha da matriz de dados  $X$ , e portanto, a entrada para os algoritmos, é representado por  $\mathbf{x}_i = [tf(d_i, t_1), \dots, tf(d_i, t_m)]$ ,  $\forall i$ . Contudo, Salton, Wong e Yang (1975) mostram a partir de experimentação em diversos conjuntos de dados textuais, que o uso da relação conhecida como Frequência de Termos-Frequência de Documentos Inversa (*Term Frequency-Inversed Document Frequency - tfidf*) é capaz de melhorar a separação de documentos. Essa relação é calculado como descrito na equação 16, e tem o efeito de fazer com que a frequência dos termos que aparecem em



muitos documentos seja enfraquecida, e a frequência dos termos que aparecem em alguns raros documentos seja fortalecida, gerando um efeito de normalização.

$$\begin{aligned} tfidf(d_i, t_j) &= tf(d_i, t_j) \times idf(t_j) \\ &= tf(d_i, t_j) \times \left( \log_2 \frac{n}{df(t_j)+1} \right) \end{aligned} \quad (16)$$

em que  $idf(t_j)$  representa a frequência de documentos inversa do termo  $t_j$ , e  $df(t_j)$  a frequência de documentos que contém  $t_j$ .

Nos experimentos presentes neste trabalho, também é usada normalização norma- $L_2$  para que todos os vetores de termos  $\mathbf{x}_i$  tenham comprimento iguais, ou seja,  $\|\mathbf{x}_i\| = 1$ . A partir dessa normalização aplicada aos vetores gerados com as relações  $tf(d_i, t_j)$  e  $tfidf(d_i, t_j)$ , obtém-se respectivamente, as relações  $tf_{norm}(d_i, t_j)$  e  $tfidf_{norm}(d_i, t_j)$ .

Sendo assim, para os experimentos presentes neste trabalho, foram utilizadas quatro formas de representação de documentos:  $tf$ ,  $tfidf$ ,  $tf_{norm}$  e  $tfidf_{norm}$ .

### 5.2.3 Experimentos quantitativos

Para avaliar os algoritmos propostos de forma quantitativa foram feitos experimentos considerando a tarefa de agrupamento de documentos.

Através do uso das classes presentes nas bases de dados *NIPS*, *IG* e *IG toy* foi possível avaliar e comparar os algoritmos *k-means*, *fuzzy k-means*, *ONMTF* (algoritmo 3 de Yoo e Choi (2010)), *FNMTF*, *OvNMTF* e *BinOvNMTF* perante métricas clássicas de agrupamento (descritas no capítulo 2): Índice de Rand e Informação Mútua Normalizada.

Primeiramente são apresentadas as configurações usadas para esses experimentos para então apresentar os resultados obtidos para cada base e uma discussão acerca dos resultados obtidos.

#### 5.2.3.1 Configuração dos experimentos

Os algoritmos usados nesse experimento foram os mesmos na experimentação com base de dados reais (seção 5.1), porém, utilizando implementações diferentes com a justificativa de viabilizar a geração de modelos para esses experimentos com bases de dados reais. Para os experimentos com os algoritmos de agrupamento tradicional (*kmeans* e *fuzzy k-means*), não foi necessário o uso de outra implementação, visto que as disponíveis

já são bem estabelecidas. Além disso, esses algoritmos têm menor complexidade de tempo comparando com os algoritmos para coagrupamento. Então, foram utilizadas as mesmas bibliotecas e, portanto, a mesma linguagem de programação: Python com as bibliotecas *scikit-learn* para o algoritmo *k-means*, e *scikit-fuzzy* para o algoritmo *fuzzy k-means*.

Para os experimentos com os algoritmos de FM para coagrupamento (capítulos 3 e 4), foram feitas implementações utilizando a linguagem de programação *C++*<sup>8</sup>. Os algoritmos *FNMTF* e *BinOvNMTF* foram reimplementados usando a biblioteca de álgebra linear *Armadillo*<sup>9</sup> (SANDERSON, 2010), para manipulação de vetores e matrizes (álgebra linear), e como motor de multiplicação de matrizes dessa biblioteca, foram usadas as rotinas do software *OpenBLAS*<sup>10</sup> (*Open Basic Linear Algebra Subprograms* (WANG et al., 2013)), que proveram aos algoritmos, computação paralela de alta performance, de forma transparente.

Já os algoritmos *ONMTF* e *OvNMTF*, que envolvem intensas multiplicações de matrizes no processo de otimização, além da biblioteca *Armadillo* para manipulação de vetores e matrizes, também foi usada a plataforma *CUDA*<sup>11</sup> com a biblioteca *cuBLAS*<sup>12</sup> para multiplicação de matrizes. A plataforma *CUDA* é uma interface para interação com Unidades de Processamento Gráfico (*Graphics Processing Unit* - GPU) da *Nvidia*, então, como a sua arquitetura é massivamente paralela, problemas paralelizáveis, como multiplicação de matrizes, se beneficiam em performance com esse hardware (FATAHALIAN; SUGERMAN; HANRAHAN, 2004). A biblioteca *cuBLAS* fornece a mesma interface de rotinas para multiplicação de matrizes que as rotinas do software *OpenBLAS*. Então, o uso dessas estratégias, foi uma maneira efetiva para melhorar a performance computacional dos algoritmos *ONMTF* e *OvNMTF*.

Ainda para os algoritmos *ONMTF* e *OvNMTF*, foi utilizado um algoritmo (CORMEN et al., 2001) para otimização da ordem das multiplicações de matrizes. Para verificar a utilidade dessa estratégia, considere o seguinte exemplo:

$$UU^T X$$

em que  $U \in \mathbb{R}^{n \times k}$ ,  $S \in \mathbb{R}^{k \times l}$  e  $V \in \mathbb{R}^{m \times l}$ , sendo  $k \approx l$ ,  $n \approx m$ ,  $n \gg k$  e  $m \gg l$ . Se a ordem das multiplicações for  $(U(U^T X))$ , serão realizadas  $nmk + nk^2 \approx n^2k + nk^2$  operações

<sup>8</sup> <https://isocpp.org/>

<sup>9</sup> <http://arma.sourceforge.net/>

<sup>10</sup> <http://www.openblas.net/>

<sup>11</sup> <https://developer.nvidia.com/cuda-zone>

<sup>12</sup> <http://docs.nvidia.com/cuda/cublas/>

de multiplicação de ponto flutuante; contrariamente, se a ordem das multiplicações for  $((UU^T)X)$ , serão realizadas  $n^2k + n^2m \approx n^3 + n^2k$  operações, que são, claramente, mais operações que as  $n^2k + nk^2$  operações.

Ainda na configuração dos experimentos para parametrização dos algoritmos, foi usado o número de classes presentes em cada base de dados, para determinar o número de grupos de documentos ( $k$ ). Já o número de grupos de termos ( $l$ ), presentes nos algoritmos *ONMTF*, *FNMTF*, *OvNMTF* e *BinOvNMTF*, como não existe nenhuma evidência de quantos grupos de termos existem nas bases de dados, foram gerados modelos de agrupamento considerando os seguintes valores de  $l$  para cada base de dados:

- *NIPS*:  $l \in \{6, 9, 12, 15, 18\}$ ;
- *IG*:  $l \in \{7, 10, 13, 16, 19\}$ ;
- *IG toy*:  $l \in \{2, 3, 4, 5, 6\}$ .

Os valores de  $l$  foram escolhidos desta forma a fim de verificar como os algoritmos iriam se comportar para  $l = k$ ,  $l < k$  e  $l > k$ .

Para todas as bases de dados foram criadas quatro representações que serviram de entrada para os algoritmos:  $tf$ ,  $tfidf$ ,  $tf_{norm}$  e  $tfidf_{norm}$ , visando comparar o efeito que os algoritmos têm sob cada representação.

Como nenhum dos algoritmos são capazes de garantir um mínimo global para minimização do erro de quantização, foram gerados 10 modelos para cada configuração, lembrando que serão geradas diferentes soluções, devido à inicialização aleatória das matrizes  $U$ ,  $C$ ,  $S$  e  $V$  (como descrito nos algoritmos apresentados nos capítulos 2, 3 e 4).

Como condições de para dos algoritmos, foi configurado o número máximo de iterações como 1.000 para os experimentos com a base de dados *IG toy* e 10.000 para as bases de dados *NIPS* e *IG*, foi observado que bases de dados com maior número de grupos de documentos necessitam de mais iterações para convergência. Também, o valor  $1 \times 10^{-4}$  foi configurado para a diferença do erro de quantização em duas iterações consecutivas.

### 5.2.3.2 Resultados

Com a parametrização dos experimentos quantitativos realizada, foram criados 10 modelos para cada combinação, por exemplo, para a representação  $tf$ , com  $k = 3$ ,  $l = 3$ , para a base de dados *IG toy* e o algoritmo *ONMTF*, foram gerados 10 modelos.

Os resultados obtidos para a base de dados *IG toy* são sumarizados nas tabelas 8 e 9, e nas figuras 18 e 19. Analisando as tabelas 8 e 9, a primeira conclusão que se pode fazer, é que o índice de Rand é consistente com a informação mútua normalizada, na maioria dos casos, os melhores resultados representados em uma também são os melhores resultados representados em outra, considerando as configurações de representação e número de grupos de termos ( $l$ ), a única exceção é para o resultado do algoritmo *k-means*, que para os resultados com o índice de Rand a melhor representação foi  $tf_{norm}$ , enquanto nos resultados com a informação mútua normalizada, o melhor resultado foi com a representação  $tf$ .

Ainda nas tabelas que mostram os resultados para a base de dados *IG toy*, é possível ver que as representações que tem maior capacidade de separar os grupos de notícias são as  $tf$  e  $tf_{norm}$ , que apresentam os melhores resultados para todos os algoritmos, especialmente a representação  $tf_{norm}$ . Os algoritmos que mais sofreram com ...

Analisando o parâmetro de número de grupos de termos, para os algoritmos de FM para coagrupamento, é possível perceber que a melhor escolha é quando  $l \geq k$ , pois não há nenhum resultado com  $l = 2$ , que teve média do índice de Rand ou informação mútua normalizada para uma representação. Isso impacta diretamente na qualidade dos grupos de documentos, pois os grupos de termos são usados para formar os grupos de documentos.

O algoritmo *OvNMTF*, proposto nesse trabalho, teve o melhor resultado nos experimentos com a base de dados *IG toy*, além de ter maior estabilidade, isso pode ser visto nas tabelas e especialmente nas figuras 18 e 19, que mostram as distribuições das medidas índice de Rand e informação mútua normalizada para a melhor configuração de representação e número de grupos de termos para cada algoritmo.

Os resultados para a base de dados *IG toy* também mostram a instabilidade e capacidade de agrupamento inferior dos algoritmos que tem as restrições binárias nas matrizes indicadoras de grupos de documentos e termos. É interessante ver que esse comportamento não acontece com os algoritmos clássicos, isto é, o algoritmo *k-means* tem capacidade de agrupamento superior ao *fuzzy k-means*. Uma possível razão para isso, pode ser as restrições binárias nas matrizes que representam os termos, enquanto o *k-means* não faz nenhuma compactação no espaço de termos, então, é capaz de representar esse espaço exatamente como ele é, os algoritmos *FNMTF* e *BinOvNMTF*, fazem essa compactação e além de fazer essa compactação, restringem essa representação para ser binária, podendo levá-los à uma representação de termos de má qualidade, e portanto, grupos de termos de

má qualidade, que diretamente influencia nos grupos de documentos. Esse efeito parece ser amenizado para o algoritmo *BinOvNMTF*, já que a capacidade de agrupamento é maior que que o *FNMTF*. Isso provavelmente ocorre pois o *BinOvNMTF* tem maior capacidade de representação dos termos devido ao uso das  $k$  matrizes  $V$ .

Tabela 8 – Índice de Rand médio por representação do conjunto de dados *IG toy* com  $k = 3$ , destacando os melhores resultados por representação.

Algoritmo	$tf$	$tf_{norm}$	$tfidf_{norm}$	$tfidf$
<i>k-means</i>	0,7017	0,7086	0,4701	0,3869
<i>fuzzy k-means</i>	0,4966	0,4970	0,4673	0,4390
<i>ONMTF</i>	0,3372 ( $l = 5$ )	0,6479 ( $l = 5$ )	0,5717 ( $l = 3$ )	0,1758 ( $l = 3$ )
<i>FNMTF</i>	0,2615 ( $l = 3$ )	0,2590 ( $l = 3$ )	0,1535 ( $l = 6$ )	0,1543 ( $l = 3$ )
<i>OvNMTF</i>	<b>0,7466</b> ( $l = 4$ )	<b>0,7487</b> ( $l = 3$ )	<b>0,6755</b> ( $l = 6$ )	<b>0,6674</b> ( $l = 6$ )
<i>BinOvNMTF</i>	0,4360 ( $l = 5$ )	0,4818 ( $l = 5$ )	0,2943 ( $l = 5$ )	0,4079 ( $l = 3$ )

Fonte: Lucas Fernandes Brunialti, 2016

Tabela 9 – Informação Mútua Normalizada média por representação do conjunto de dados *IG toy* com  $k = 3$ , destacando os melhores resultados por representação.

Algoritmo	$tf$	$tf_{norm}$	$tfidf_{norm}$	$tfidf$
<i>k-means</i>	0,7169	0,7134	0,5467	0,4668
<i>fuzzy k-means</i>	0,0694	0,5421	0,4701	0,0836
<i>ONMTF</i>	0,3704 ( $l = 5$ )	0,6720 ( $l = 5$ )	0,6411 ( $l = 3$ )	0,2143 ( $l = 3$ )
<i>FNMTF</i>	0,2770 ( $l = 3$ )	0,2734 ( $l = 3$ )	0,2039 ( $l = 6$ )	0,1690 ( $l = 3$ )
<i>OvNMTF</i>	<b>0,7257</b> ( $l = 4$ )	<b>0,7288</b> ( $l = 3$ )	<b>0,7033</b> ( $l = 6$ )	<b>0,6964</b> ( $l = 6$ )
<i>BinOvNMTF</i>	0,4975 ( $l = 5$ )	0,5500 ( $l = 5$ )	0,3559 ( $l = 5$ )	0,4343 ( $l = 3$ )

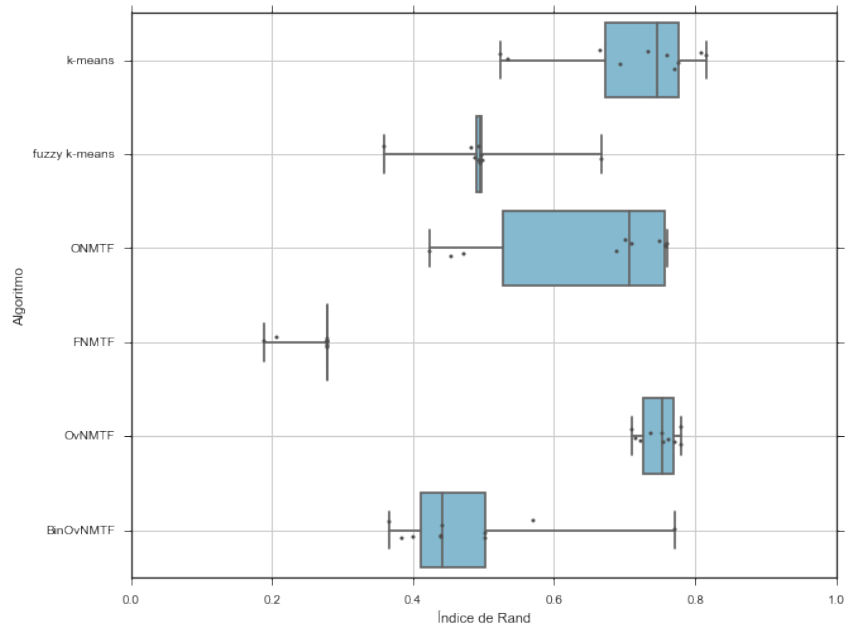
Fonte: Lucas Fernandes Brunialti, 2016

Tabela 10 – Melhores resultados obtidos para o conjunto de dados *IG toy* com  $k = 3$ , destacando os melhores resultados por algoritmo.

Algoritmo	Índice de Rand	Informação Mútua Normalizada
<i>k-means</i>	0,8152* ( $tf_{norm}$ )	<b>0,8110</b> ( $tf_{norm}$ )
<i>fuzzy k-means</i>	0,6669 ( $tf_{norm}$ )	0,6785 ( $tf$ )
<i>ONMTF</i>	0,7885 ( $tfidf_{norm}$ , $l = 5$ )	0,7719 ( $tfidf_{norm}$ , $l = 5$ )
<i>FNMTF</i>	0,4502 ( $tfidf$ , $l = 4$ )	0,5041 ( $tfidf$ , $l = 4$ )
<i>OvNMTF</i>	0,8208* ( $tf_{norm}$ , $l = 2$ )	0,7855* ( $tf_{norm}$ , $l = 2$ )
<i>BinOvNMTF</i>	<b>0,8261</b> ( $tfidf$ , $l = 3$ )	0,8024* ( $tfidf$ , $l = 3$ )

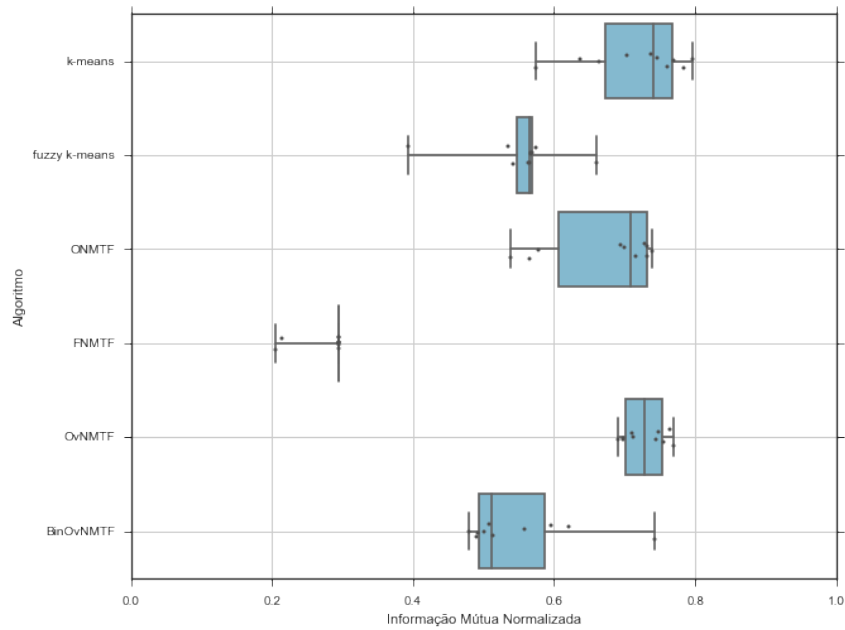
Fonte: Lucas Fernandes Brunialti, 2016

Figura 18 – Distribuições em formato de diagrama de caixa do Índice de Rand para as melhores configurações de  $l$  e representação para cada algoritmo na base de dados *IG toy*.



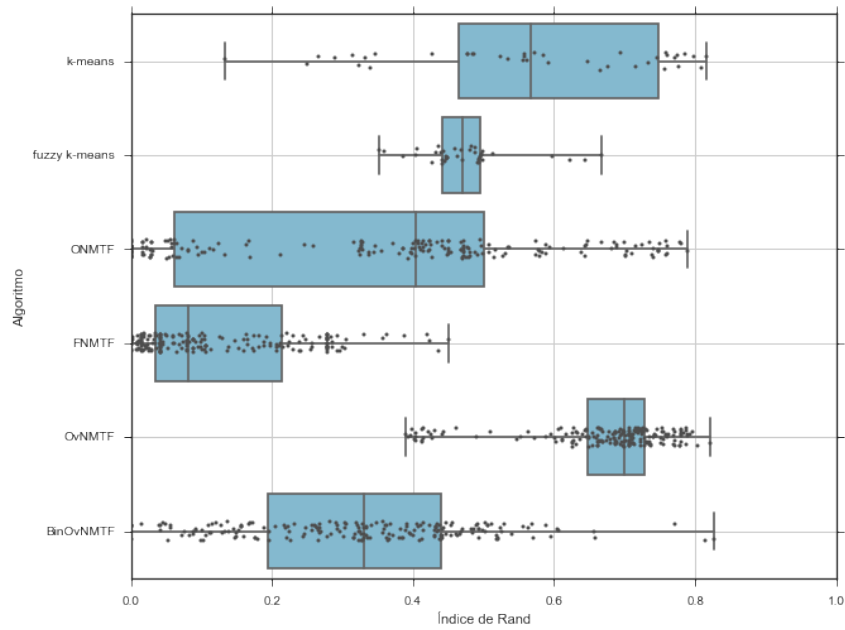
Fonte: Lucas Fernandes Brunialti, 2016

Figura 19 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada para as melhores configurações de  $l$  e representação para cada algoritmo na base de dados *IG toy*.



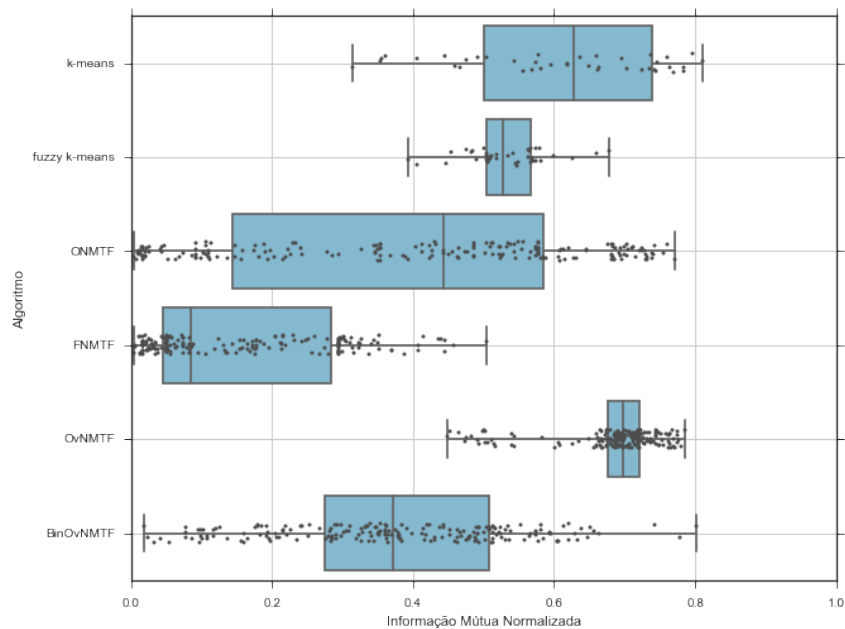
Fonte: Lucas Fernandes Brunialti, 2016

Figura 20 – Distribuições em formato de diagrama de caixa do Índice de Rand de todas as configurações para cada algoritmo na base de dados *IG toy*.



Fonte: Lucas Fernandes Brunialti, 2016

Figura 21 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada de todas as configurações para cada algoritmo na base de dados *IG toy*.



Fonte: Lucas Fernandes Brunialti, 2016

Tabela 11 – Melhores resultados do Índice de Rand para experimentos com a base de dados IG. \* - testes não finalizados

Algoritmo	Representação	$l$	Rand Index
<i>k-means</i>	TF	-	0.3082
<i>k-means</i>	TF norm	-	0.3350
<i>k-means</i>	TF-IDF	-	0.2750
<i>k-means</i>	TF-IDF norm	-	0.2262
<i>fuzzy k-means</i>	TF	-	0.2681
<i>fuzzy k-means</i>	TF norm	-	0.1757
<i>fuzzy k-means</i>	TF-IDF	-	0.2347
<i>fuzzy k-means</i>	TF-IDF norm	-	0.3832
<i>ONMTF</i>	TF	19	0.1883
<i>ONMTF</i>	TF norm	13	0.2023
<i>ONMTF</i>	TF-IDF	13	0.1814
<i>ONMTF</i>	TF-IDF norm	7	0.2037
<i>FNMTF</i>	TF	13	0.2367
<i>FNMTF</i>	TF norm	19	0.2678
<i>FNMTF</i>	TF-IDF	7	0.2410
<i>FNMTF</i>	TF-IDF norm	16	0.2480
<i>OvNMTF</i>	TF	19	0.3520*
<i>OvNMTF</i>	TF norm	16	0.3360*
<i>OvNMTF</i>	TF-IDF	10	0.3682*
<i>OvNMTF</i>	TF-IDF norm	10	0.3724*
<i>BinOvNMTF</i>	TF	16	0.4415
<i>BinOvNMTF</i>	TF norm	7	0.4081
<i>BinOvNMTF</i>	TF-IDF	7	0.3681
<i>BinOvNMTF</i>	TF-IDF norm	16	0.3883

Tabela 12 – Índice de Rand médio por representação do conjunto de dados *NIPS* com  $k = 9$ , destacando os melhores resultados por representação.

Algoritmo	$tf$	$tf_{norm}$	$tfidf_{norm}$	$tfidf$
<i>k-means</i>	0.1573	0.1527	0.1519	0.1368
<i>fuzzy k-means</i>	0.1223	0.1240	0.1736	0.1882
<i>ONMTF</i>	0, 1579 ( $l = 6$ )	0, 1352 ( $l = 15$ )	0, 1318 ( $l = 9$ )	0, 1442 ( $l = 18$ )
<i>FNMTF</i>	0, 1293 ( $l = 18$ )	0, 1325 ( $l = 18$ )	0, 2128 ( $l = 18$ )	0, 2199 ( $l = 18$ )
<i>OvNMTF</i>	0, 1672 ( $l = 6$ )	0, 1641 ( $l = 9$ )	0, 1742 ( $l = 12$ )	0, 1711 ( $l = 9$ )
<i>BinOvNMTF</i>	<b>0, 2247</b> ( $l = 9$ )	<b>0, 2118</b> ( $l = 15$ )	<b>0, 2811</b> ( $l = 6$ )	<b>0, 2813</b> ( $l = 15$ )

Fonte: Lucas Fernandes Brunialti, 2016

Tabela 13 – Informação Mútua Normalizada média por representação do conjunto de dados *NIPS* com  $k = 9$ , destacando os melhores resultados por representação.

Algoritmo	$tf$	$tf_{norm}$	$tfidf_{norm}$	$tfidf$
<i>k-means</i>	0.3226	0.3106	0.3506	0.3476
<i>fuzzy k-means</i>	0.1876	0.1929	0.2575	0.2496
<i>ONMTF</i>	0, 2930 ( $l = 15$ )	0, 2832 ( $l = 18$ )	0, 3361 ( $l = 18$ )	0, 3441 ( $l = 18$ )
<i>FNMTF</i>	0, 2272 ( $l = 18$ )	0, 2312 ( $l = 18$ )	0, 3109 ( $l = 18$ )	0, 3017 ( $l = 18$ )
<i>OvNMTF</i>	0, 3090 ( $l = 6$ )	0, 3092 ( $l = 9$ )	0, 3541 ( $l = 12$ )	0, 3493 ( $l = 15$ )
<i>BinOvNMTF</i>	<b>0, 3255</b> ( $l = 15$ )	<b>0, 3139</b> ( $l = 9$ )	<b>0, 4013</b> ( $l = 12$ )	<b>0, 4009</b> ( $l = 15$ )

Fonte: Lucas Fernandes Brunialti, 2016

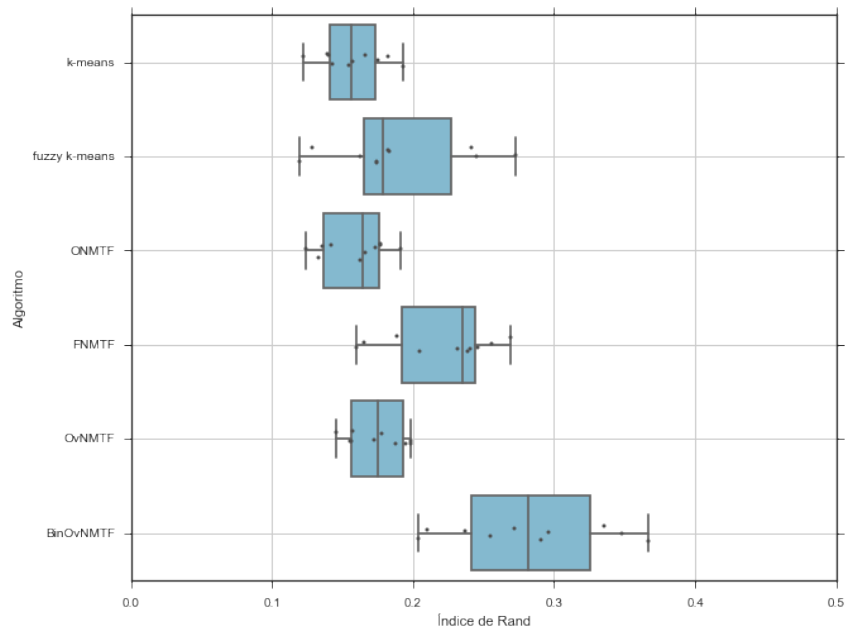


Tabela 14 – Melhores resultados obtidos para o conjunto de dados *NIPS* com  $k = 9$ , destacando os melhores resultados por algoritmo.

Algoritmo	Índice de Rand	Informação Mútua Normalizada
<i>k-means</i>	0,2006 ( $tfidf_{norm}$ )	0,3952 ( $tfidf_{norm}$ )
<i>fuzzy k-means</i>	0,2728 ( $tfidf$ )	0,3115 ( $tfidf$ )
<i>ONMTF</i>	0,2704 ( $tfidf$ , $l = 18$ )	0,3997* ( $tfidf$ , $l = 18$ )
<i>FNMTF</i>	0,3009* ( $tfidf$ , $l = 15$ )	0,3744 ( $tfidf_{norm}$ , $l = 18$ )
<i>OvNMTF</i>	0,1992 ( $tfidf$ , $l = 9$ )	0,3870 ( $tfidf$ , $l = 9$ )
<i>BinOvNMTF</i>	<b>0,3670</b> ( $tfidf$ , $l = 15$ )	<b>0,4589</b> ( $tfidf_{norm}$ , $l = 9$ )

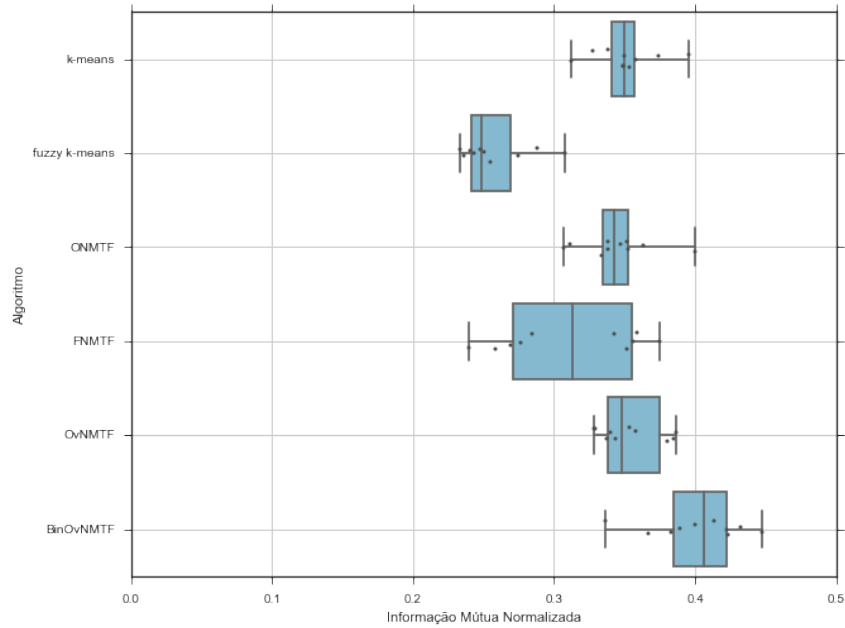
Fonte: Lucas Fernandes Brunialti, 2016

Figura 22 – Distribuições em formato de diagrama de caixa do Índice de Rand para as melhores configurações de  $l$  e representação para cada algoritmo na base de dados *NIPS*.



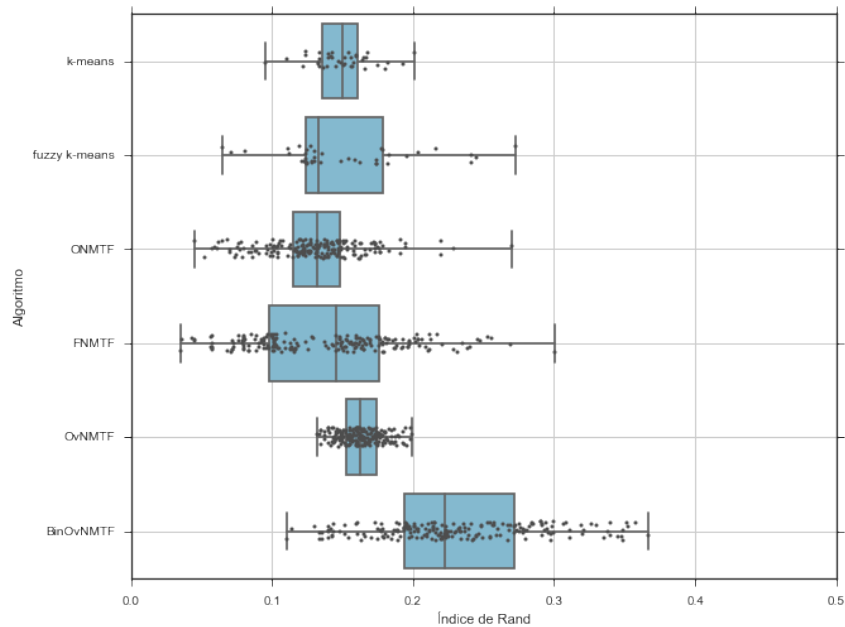
Fonte: Lucas Fernandes Brunialti, 2016

Figura 23 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada para as melhores configurações de  $l$  e representação para cada algoritmo na base de dados *NIPS*.



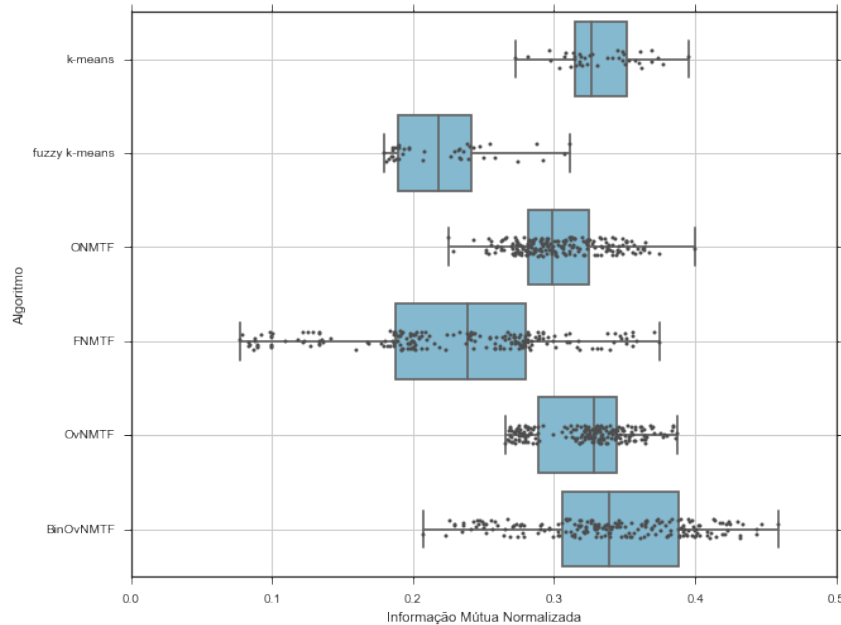
Fonte: Lucas Fernandes Brunialti, 2016

Figura 24 – Distribuições em formato de diagrama de caixa do Índice de Rand de todas as configurações para cada algoritmo na base de dados *NIPS*.



Fonte: Lucas Fernandes Brunialti, 2016

Figura 25 – Distribuições em formato de diagrama de caixa da Informação Mútua Normalizada de todas as configurações para cada algoritmo na base de dados *NIPS*.



Fonte: Lucas Fernandes Brunialti, 2016

#### 5.2.4 Análise qualitativa

Estes experimentos mostram a aplicabilidade dos algoritmos no domínio de mineração de textos, com ênfase nas informações que são possíveis de extrair dos modelos gerados por algoritmos de coagrupamento baseados em fatoração de matrizes: *ONMTF* e *OvNMTF*.

Essa análise qualitativa se faz importante, para mostrar que os algoritmos são capazes de encontrar tópicos (grupos de palavras), e possivelmente, explicar os grupos de notícias formados, ou até mesmo, realizar rotulação dos grupos de notícias formados.

Foi necessária a construção de um estudo de caso para avaliar os algoritmos de Coagrupamento baseados em Fatoração de Matrizes. O estudo de caso é composto por uma análise dos cogrupos de palavras correlacionando-os com os cogrupos de notícias.

Para isso, foi usado o conjunto de dados *IG toy*, que foi construído, justamente, para fazer esse tipo de análise em um ambiente com maior controle.

O portal IG<sup>13</sup> é um dos maiores portais de notícias brasileiro (SITES..., 2016) que é composto por sites importantes como o noticiário Último Segundo, o iG Gente, o iG

<sup>13</sup> <http://www.ig.com.br/>

Esportes, a TV iG, o iG Economia, o Delas e etc. Cada um desses sites é caracterizado como um canal que contém notícias de um determinado assunto macro.

O conjunto de dados *IG toy* é composto por três desses canais: esportes, que contém notícias, principalmente, dos esportes mais populares no Brasil, como o futebol e o UFC; jovem, que contém notícias mais interessantes para o público jovem em geral, como notícias sobre filmes, esportes e músicas voltados para o público jovem; arena, que compõem notícias de games, novidades de todos os tipos de games, consoles e coberturas de eventos de games. Sendo 100 notícias para cada canal para compor o conjunto de dados, totalizando 300 notícias.

Esses canais foram escolhidos para formar o conjunto de dados *IG toy* por serem de assuntos similares, notícias do canal jovem podem ter semelhanças com notícias do canal esportes, como notícias sobre surfe, por exemplo, notícias do canal arena podem ser similares com notícias do canal de esportes, por existirem games que simulam esportes, ou até mesmo, possuírem similaridades com notícias do canal jovem, pois games é um assunto ligado ao público jovem na sua maioria. Essa escolha torna possível verificar como os algoritmos tratam essa intersecção entre palavras nas notícias.

As subseções a seguir irão mostrar como os algoritmos *ONMTF* e *OvNMTF* são úteis para análise de tópicos e palavras no conjunto de dados *IG toy*.

#### 5.2.4.1 Análise de dados utilizando *ONMTF*

Para as análises construídas com o *ONMTF* foi utilizado o modelo que obteve o melhor resultado considerando o *Índice de Rand* na base de dados *IG toy*, descrito na seção 5.2.3.2, que foi o modelo com a configuração:  $k = 3$ ,  $l = 5$  e representação  $tfidf_{norm}$ .

O algoritmo *ONMTF* é capaz de encontrar cogrupos de notícias e cogrupos de palavras, além disso, cada cogrupo de notícias é relacionado com os cogrupos de palavras por um fator, assim como cada cogrupo de palavras está relacionado com os cogrupos de notícias pelo mesmo fator.

Os fatores que conectam os cogrupos de notícias com os cogrupos de palavras podem ser observados na matriz  $S$ . A matriz  $S$  que foi usada nesta análise pode ser observada na Tabela ?? . Note que os valores estão normalizados para que a soma dos elementos de cada linha seja igual à 1. A normalização foi feita para tornar claro que cada cogrupo de palavras foi usado pelo algoritmo para caracterizar um grupo de notícias, os valores

mostram que o cogruppo de notícias “arena” está associado com um maior fator ao cogruppo de palavras #4, o cogruppo de notícias “esporte” está associado com fatores semelhantes aos cogruppo de palavras #2 e #5, e o cogruppo de notícias “jovem” está associado com fatores semelhantes aos cogruppos de palavras #1 e #3. No entanto, essa normalização muda a solução obtida anteriormente, então, foi realizada apenas como caráter ilustrativo.

Tabela 15 – Matriz  $S$  para o algoritmo  $ONMTF$  com  $k = 3$  e  $l = 5$  na base de dados  $IG\ toy$  normalizada para que a soma da linha seja igual à 1

	CP #1	CP #2	CP #3	CP #4	CP #5
CN “esportes”	0,0	<b>0,5</b>	0,1	0,0	<b>0,4</b>
CN “arena”	0,0	0,05	0,05	<b>0,9</b>	0,0
CN “jovem”	<b>0,4</b>	0,1	<b>0,5</b>	0,0	0,0

Fonte: Lucas Fernandes Brunialti, 2016

A figura 26 exemplifica as informações que um modelo gerado pelo algoritmo  $ONMTF$  provê. A notícia “Avaliação do FIFA 15 por um jogador fanático” foi usada como exemplo. O resultado do coagrupamento de notícias foi que dos três cogruppos de notícias, a notícia pertence ao cogruppo (rotulado manualmente através da análise de palavras) “esportes” com um fator equivalente à 40%, também pertence ao cogruppo rotulado como “arena” com um fator equivalente à 60%, e não pertence ao cogruppo rotulado como “jovem”. Cada cogruppo de notícias é formado pela combinação dos cogruppos de palavras com os fatores que indicam a relação entre cogruppos de notícias e cogruppos de colunas, denotado pelas linhas que os conecta. Cada cogruppo de palavras é ilustrado pelas palavras mais relevantes que o compõem, ou seja, as palavras que contém os maiores fatores (representado em porcentagem) para aquele determinado cogruppo.

Ainda sobre a Figura 26, note que as palavras do corpo da notícia foram coloridas de acordo com as cores dos cogruppos de palavras as quais pertencem. É possível perceber que existem mais palavras que caracterizam a notícia como sendo sobre o assunto “games”, o que vai de acordo com o coagrupamento realizado.

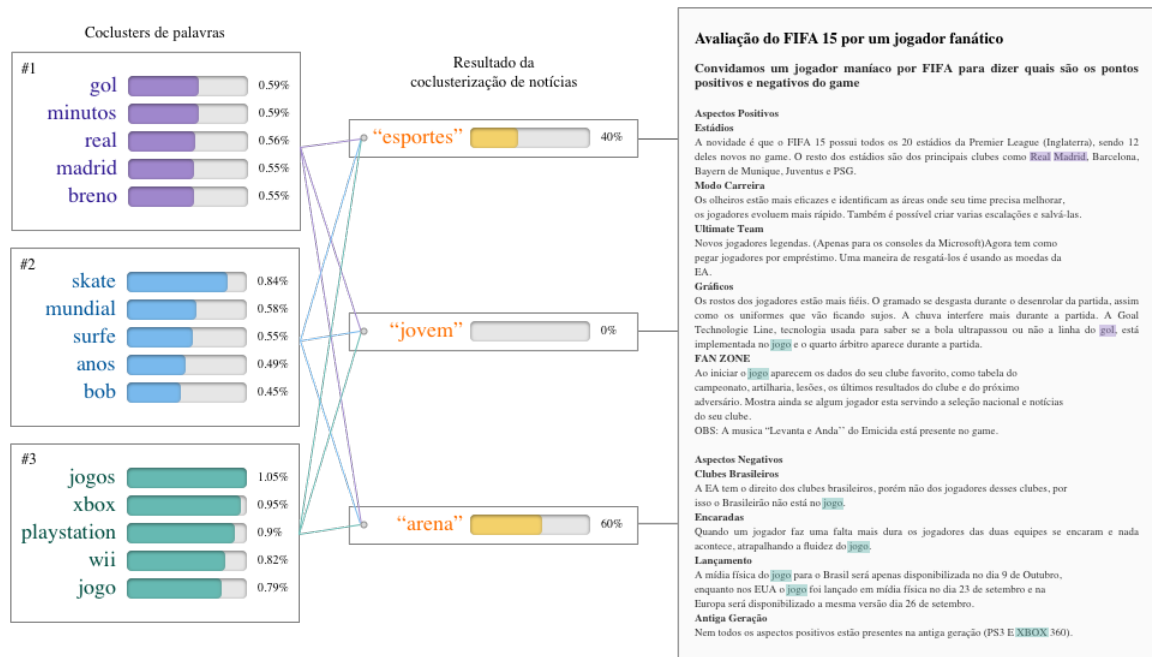


Figura 26 – Exemplo de uma notícia do canal arena e sua disposição quanto aos cogrupos de notícias e palavras.

Fonte: Lucas Fernandes Brunialti, 2016

A rotulação dos grupos de notícias foi possível através da análise de palavras em cada cogrupos de palavras e da matriz  $S$ , utilizando o particionamento (descrito no capítulo 3), que atribui cada palavra ao cogrupos que possui o maior fator na matriz  $V$ . O resultado das palavras mais relevantes para cada cogrupos de palavras, gerado pelo algoritmo  $ONMTF$ , é mostrado na tabela 16. A ordenação foi feita através dos fatores de cada palavra para cada cogrupos de palavras, presentes na matriz  $V$ .

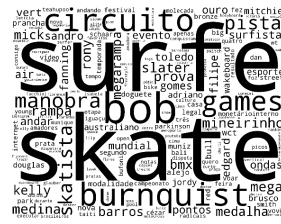
Tabela 16 – top-20 palavras para cada cogruppo de palavras, após a realização do particionamento

CP #1 “esportes radicais”	CP #2 “futebol”	CP #3 “esportes em geral”	CP #4 “games”	CP #5 “futebol”
skate	real	anos	jogos	gol
surfe	breno	mundial	xbox	madrid
bob	time	brasil	playstation	gols
burnquist	barcelona	etapa	wii	messi
circuito	partida	brasileiro	jogo	euro
games	equipe	jovem	of	guardiola
mineirinho	minutos	rio	console	ronaldo
slater	jogador	dias	sony	itália
rampa	campeonato	música	legends	cristiano
medina	liga	pedro	nintendo	bola
manobras	futebol	atleta	game	bayern
mega	casa	americano	league	pontos
megarampa	temporada	gente	one	espanhol
kelly	grupo	esporte	novo	zagueiro
prova	feira	campeão	ps	atacante
bmX	dois	competição	microsoft	defesa
pista	vitória	mundo	the	técnico
barros	tempo	ano	usmonetáriointerno	placar
skatistas	área	janeiro	controle	atlético
rony	copa	ufc	versão	polônia

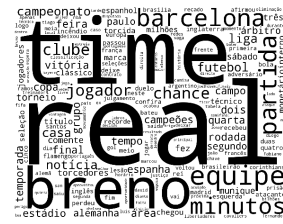
Fonte: Lucas Fernandes Brunialti, 2016

A tabela 16 mostra claramente que os cogrupos de palavras # 2 e # 5 (rotulados como “futebol”) foram responsáveis por caracterizar o cogruppo de notícias do canal “esportes”, o cogruppo de palavras # 4 (rotulado como “games”) foi responsável por caracterizar o cogruppo de notícias do canal “arena”, e por fim, os cogrupos de palavras # 1 e # 3 (rotulados como “esportes radicais” e “esportes em geral”, respectivamente) foram responsáveis por caracterizar as notícias do cogruppo do canal “jovem”.

Foi criada outra forma de visualização para as palavras do resultado do coagrupamento nesse experimento, essa visualização apresenta-se como *word cloud*, em que o tamanho das palavras é definido pelo seu fator de pertinência ao cogruppo correspondente. A Figura 27 mostra essa visualização contendo as 100 palavras com maior fator para cada coluna da matriz  $V$ .



(a) CP #1



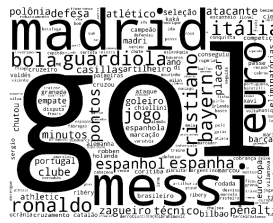
(b) CP #2



(c) CP #3



(d) CP #4



(e) CP #5

Figura 27 – Visualização *word cloud* das top-100 palavras, dispostas em uma área quadrada, para cada cogruppo gerados pelo algoritmo *ONMTF*

Fonte: Lucas Fernandes Brunialti, 2016

Note que neste caso uma palavra se apresenta em diferentes cogrupos de palavras, como exemplo a palavra *jogo* presente nos cogrupos # 4 e # 5 (a maior palavra do cogruppo # 4 e no centro do ‘o’ da palavra *gol* para o cogruppo # 5), que significa que existem notícias em que a palavra ocorre para os cogrupos de notícias do canal “arena” e “esportes”, respectivamente, com maior ocorrência no primeiro do que no segundo. Isso ocorre pois foram usados os 100 maiores de cada uma das  $l$  colunas na matriz  $V$ , gerada pela fatoração *ONMTF*. Note que é possível dizer que a palavra *jogo*, por exemplo, possui diferentes significados para cada cogruppo de palavras. No entanto, como há interdependência entre cogrupos de palavras, essa afirmação não é totalmente verdadeira, pois a palavra *jogo* carregando o significado de “futebol” pode afetar negativamente o cogruppo de palavras # 4, assim como, a palavra *jogo* com o significado de “game” pode afetar negativamente o cogruppo de palavras # 5.



### 5.2.4.2 Análise de dados utilizando *OvNMTF*

O algoritmo *OvNMTF* é capaz de encontrar cogrupos de notícias e cogrupos de palavras, semelhante aos outros.

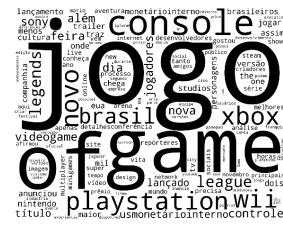
Tabela 17 – Matriz  $S$  para execução do algoritmo *OvNMTF* com  $k = 3$  e  $l = 2$  na base de dados *IG toy*.

	CP #1, #3, #4	CP #2, #4, #6
CN “arena”	0.38	0.62
CN “jovem”	0.46	0.54
CN “esportes”	0.94	0.06

Fonte: Lucas Fernandes Brunialti, 2016



(a) #1



(b) #2

Figura 28 – Visualização *word cloud* de palavras para cada cogrupos de palavras do cogrupos de notícias “arena”, gerados pelo algoritmo *OvNMTF*.

Fonte: Lucas Fernandes Brunialti, 2016



(a) #3



(b) #4

Figura 29 – Visualização *word cloud* de palavras para cada cogrupos de palavras do cogrupos de notícias “jovem”, gerados pelo algoritmo *OvNMTF*.

Fonte: Lucas Fernandes Brunialti, 2016

Tabela 18 – top-20 palavras para cada cogrupo de palavras, após a realização do particionamento

CP #1 “ <i>esportes radicais</i> ”	CP #2 “ <i>futebol</i> ”	CP #3 “ <i>esportes em geral</i> ”	CP #4 “ <i>games</i> ”	CP #5 “ <i>futebol</i> ”	CP #6 “ <i>x</i> ”
jogos	jogo	games	anos	time	breno
sony	of	jovem	skate	real	casa
ps	playstation	brasileiro	mundial	feira	gol
the	game	paulo	brasil	final	bayern
pessoas	novo	dia	surfe	gols	minutos
wii	console	mundo	música	madrid	clube
microsoft	xbox	ano	rio	jogador	partida
nintendo	games	esporte	conta	tempo	técnico
estúdio	league	vai	dias	pontos	título
one	legends	janeiro	primeira	grupo	livre
arena	brasil	bem	final	liga	jogo
melhor	além	burnquist	atleta	fez	meia
apenas	nova	além	pessoas	brasileiro	volta
lançamento	jogadores	gente	casa	jogadores	segunda
versão	dia	bob	paulista	campo	técnica
opiniões	lançado	série	ficou	rodada	casillas
caio	usmonetáriointerno	etapa	fim	três	espanha
site	personagens	história	monetáriointerno	cristiano	equipe
and	feira	melhor	melhores	copa	argentino
forma	dois	circuito	amigos	deixe	semana

Fonte: Lucas Fernandes Brunialti, 2016



(a) #5



(b) #6

Figura 30 – Visualização *word cloud* de palavras para cada cogruppo de palavras do cogruppo de notícias “esportes”, gerados pelo algoritmo *OvNMTF*.

Fonte: Lucas Fernandes Brunialti, 2016

## Referências<sup>14</sup>

- CABANES, G.; BENNANI, Y.; FRESNEAU, D. Enriched topological learning for cluster detection and visualization. *Neural Networks*, v. 32, p. 186–195, 2012. Citado na página 29.
- CHENG, Y.; CHURCH, G. M. Biclustering of expression data. In: *Procedures of the 8th ISMB*. [S.l.]: AAAI Press, 2000. p. 93–103. Citado 2 vezes nas páginas 28 e 29.
- CORMEN, T. H. et al. *Introduction to Algorithms*. 2nd. ed. [S.l.]: McGraw-Hill Higher Education, 2001. ISBN 0070131511. Citado 2 vezes nas páginas 36 e 81.
- DING, C.; HE, X.; SIMON, H. D. On the equivalence of nonnegative matrix factorization and spectral clustering. In: *in SIAM International Conference on Data Mining*. [S.l.: s.n.], 2005. Citado na página 60.
- DING, C. et al. Orthogonal nonnegative matrix tri-factorizations for clustering. In: *In SIGKDD*. [S.l.]: Press, 2006. p. 126–135. Citado 5 vezes nas páginas 31, 32, 37, 38 e 40.
- FATAHALIAN, K.; SUGERMAN, J.; HANRAHAN, P. Understanding the efficiency of gpu algorithms for matrix-matrix multiplication. In: *Proceedings of the ACM SIGGRAPH/EUROGRAPHICS Conference on Graphics Hardware*. New York, NY, USA: ACM, 2004. (HWWS '04), p. 133–137. ISBN 3-905673-15-0. Disponível em: <<http://doi.acm.org/10.1145/1058129.1058148>>. Citado na página 81.
- FELDMAN, R.; SANGER, J. *The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data*. Cambridge, MA, USA: Cambridge University Press, 2006. Hardcover. Citado na página 103.
- FRANCA, F. de. *Biclusterização na Análise de Dados Incertos*. Tese (Doutorado) — Universidade Estadual de Campinas, Campinas, SP, BR, 11 2010. Citado 2 vezes nas páginas 20 e 26.
- FRANÇA, F. de; ZUBEN, F. V. Finding a high coverage set of 5-biclusters with swarm intelligence. In: *Evolutionary Computation (CEC), 2010 IEEE Congress on*. [S.l.: s.n.], 2010. p. 1–8. Citado na página 29.
- GETZ, G.; LEVINE, E.; DOMANY, E. Coupled two-way clustering analysis of gene microarray data. *Proc. Natl. Acad. Sci. USA*, v. 97, p. 12079–12084, 2000. Citado 2 vezes nas páginas 27 e 28.
- GLOBERSON, A. et al. Euclidean embedding of co-occurrence data. *The Journal of Machine Learning Research*, MIT Press Cambridge, MA, USA, v. 8, p. 2265–2295, 2007. Citado na página 75.
- HAN, J.; KAMBER, M. *Data mining: Concepts and Techniques*. 2. ed. [S.l.]: Morgan Kaufmann San Francisco, Calif, USA, 2006. Citado na página 19.
- HAN, J.; KAMBER, M.; PEI, J. *Data Mining: Concepts and Techniques*. 3rd. ed. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2011. Citado na página 105.

<sup>14</sup> De acordo com a Associação Brasileira de Normas Técnicas. NBR 6023.

HAYKIN, S. *Neural Networks and Learning Machines (3rd Edition)*. 3. ed. [S.l.]: Prentice Hall, 2008. Hardcover. Citado na página 105.

HO, N.-D. *Nonnegative Matrix Factorization Algorithms and Applications*. Tese (Doutorado) — Université Catholique de Louvain, Louvain-la-Neuve, Belgique, 6 2008. Citado 2 vezes nas páginas 20 e 31.

HOCHREITER, S. et al. Fabia: factor analysis for bicluster acquisition. *Bioinformatics*, v. 26, n. 12, p. 1520–1527, 2010. Citado 3 vezes nas páginas 25, 29 e 30.

HOTH, A.; NÜRNBERGER, A.; PAAß, G. A brief survey of text mining. *LDV Forum - GLDV Journal for Computational Linguistics and Language Technology*, v. 20, n. 1, p. 19–62, maio 2005. Citado na página 103.

JAIN, A. K.; MURTY, M. N.; FLYNN, P. J. Data clustering: A review. *ACM Computing Surveys*, ACM, v. 31, p. 264 – 323, September 1999. Citado na página 19.

KLUGER, Y. et al. Spectral biclustering of microarray data: Coclustering genes and conditions. *Genome Research*, v. 13, n. 4, p. 703–716, 2003. Citado na página 31.

KOREN, Y. *1 The BellKor Solution to the Netflix Grand Prize*. 2009. Citado na página 31.

KUANG, D. *Nonnegative Matrix Factorization For Clustering*. Tese (Doutorado) — University of Tampere, Tampere, Finlândia, 2014. Citado na página 31.

LEE, D. D.; SEUNG, H. S. Learning the parts of objects by nonnegative matrix factorization. *Nature*, v. 401, p. 788–791, 1999. Citado 2 vezes nas páginas 21 e 31.

LEE, D. D.; SEUNG, H. S. Algorithms for non-negative matrix factorization. In: *NIPS*. [s.n.], 2000. p. 556–562. Disponível em: <[citeseer.ist.psu.edu/lee01algorithms.html](http://citeseer.ist.psu.edu/lee01algorithms.html)>. Citado na página 21.

LI, T.; DING, C. The relationships among various nonnegative matrix factorization methods for clustering. In: *Proceedings of the Sixth International Conference on Data Mining*. Washington, DC, USA: IEEE Computer Society, 2006. (ICDM '06), p. 362–371. ISBN 0-7695-2701-9. Disponível em: <<http://dx.doi.org/10.1109/ICDM.2006.160>>. Citado na página 32.

LONG, B.; ZHANG, Z. M.; YU, P. S. *Co-clustering by block value decomposition*. [S.l.]: ACM Press, 2005. 635–640 p. Citado 6 vezes nas páginas 31, 32, 33, 34, 35 e 50.

LOPS, P.; GEMMIS, M. de; SEMERARO, G. Content-based recommender systems: State of the art and trends. In: RICCI, F. et al. (Ed.). *Recommender Systems Handbook*. [S.l.]: Springer, 2011. p. 73–105. Citado 3 vezes nas páginas 79, 103 e 104.

MADEIRA, S. C.; OLIVEIRA, A. L. Biclustering algorithms for biological data analysis: A survey. *IEEE Transactions on Computational Biology and Bioinformatics*, IEEE Computer Society Press, Los Alamitos, CA, USA, v. 1, p. 24–45, January 2004. Citado 3 vezes nas páginas 26, 27 e 61.

MINER, G. et al. *Practical Text Mining and Statistical Analysis for Non-structured Text Data Applications*. 1st. ed. [S.l.]: Academic Press, 2012. Citado na página 105.

MURPHY, K. P. *Machine Learning: A Probabilistic Perspective*. [S.l.]: The MIT Press, 2012. Citado na página 105.

PEDREGOSA, F. et al. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, v. 12, p. 2825–2830, 2011. Citado na página 63.

PRELIĆ, A. et al. A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics*, Oxford University Press, Oxford, UK, v. 22, n. 9, p. 1122–1129, maio 2006. Citado na página 29.

ROCHA, T. et al. Tutorial sobre fuzzy-c-means e fuzzy learning vector quantization: Abordagens híbridadas para tarefas de agrupamento e classificação. *RITA - Revista de Informática Teórica e Aplicada*, v. 19, n. 1, p. 120–163, March. Citado na página 66.

SALAKHUTDINOV, R.; MNIH, A. Probabilistic matrix factorization. In: *Advances in Neural Information Processing Systems*. [S.l.: s.n.], 2008. v. 20. Citado na página 31.

SALTON, G.; WONG, A.; YANG, C. S. A vector space model for automatic indexing. *Communications of the ACM*, ACM, New York, NY, USA, v. 18, n. 11, p. 613–620, 1975. Citado 3 vezes nas páginas 79, 103 e 104.

SANDERSON, C. *Armadillo: An Open Source C++ Linear Algebra Library for Fast Prototyping and Computationally Intensive Experiments*. [S.l.], 2010. Citado na página 81.

SANTAMARÍA, R.; MIGUEL, L.; THERÓN, R. Methods to bicluster validation and comparison in microarray data. *Lecture Notes in Computer Science: Proceedings of IDEAL'07*, v. 4881, p. 780–789, 2007. Citado 2 vezes nas páginas 25 e 29.

SEBASTIANI, F. Machine learning in automated text categorization. *ACM Comput. Surv.*, ACM, New York, NY, USA, v. 34, n. 1, p. 1–47, 2002. Citado 3 vezes nas páginas 79, 103 e 104.

SHAHNAZ, F. et al. Document clustering using nonnegative matrix factorization. *Information Processing & Management*, v. 42, n. 2, p. 373 – 386, 2006. Citado na página 21.

SITES mais acessados do Brasil. 2016. Disponível em: <<http://www.alexa.com/topsites/countries;1/BR>>. Citado na página 90.

TANAY, A.; SHARAN, R.; SHAMIR, R. Biclustering algorithms: A survey. In: *In Handbook of Computational Molecular Biology Edited by: Aluru S. Chapman & Hall/CRC Computer and Information Science Series*. [S.l.: s.n.], 2005. Citado na página 27.

TJHI, W.-C.; CHEN, L. Dual fuzzy-possibilistic coclustering for categorization of documents. *IEEE Transactions on Fuzzy Systems*, IEEE, v. 17, p. 533 – 543, June 2009. Citado na página 21.

WANG, H. et al. Fast nonnegative matrix tri-factorization for large-scale data co-clustering. In: *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Two*. AAAI Press, 2011. (IJCAI'11), p. 1553–1558. ISBN 978-1-57735-514-4. Disponível em: <<http://dx.doi.org/10.5591/978-1-57735-516-8-IJCAI11-261>>. Citado 3 vezes nas páginas 32, 41 e 55.

WANG, Q. et al. Augem: Automatically generate high performance dense linear algebra kernels on x86 cpus. In: *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. New York, NY, USA: ACM, 2013. (SC '13), p. 25:1–25:12. ISBN 978-1-4503-2378-9. Disponível em: <http://doi.acm.org/10.1145/2503210.2503219>. Citado na página 81.

WEISS, S. M.; INDURKHYA, N.; ZHANG, T. *Fundamentals of predictive text mining*. London; New York: Springer-Verlag, 2010. Citado na página 104.

XU, R.; WUNSCH, D. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, IEEE, v. 16, p. 645 – 678, May 2005. Citado na página 19.

XU, W.; LIU, X.; GONG, Y. Document clustering based on non-negative matrix factorization. In: *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Informaion Retrieval*. New York, NY, USA: ACM, 2003. (SIGIR '03), p. 267–273. ISBN 1-58113-646-3. Disponível em: <http://doi.acm.org/10.1145/860435.860485>. Citado na página 21.

YANG, J.; LESKOVEC, J. Overlapping community detection at scale: A nonnegative matrix factorization approach. In: *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*. New York, NY, USA: ACM, 2013. (WSDM '13), p. 587–596. Citado na página 29.

YOO, J.; CHOI, S. Orthogonal nonnegative matrix tri-factorizations for co-clustering: multiplicative updates on stiefel manifolds. *Information Proessing and Management*, v. 46, p. 559–570, 2010. Citado 10 vezes nas páginas 21, 31, 33, 37, 39, 40, 51, 60, 68 e 80.

## Apêndice A – Mineração de Texto

Técnicas de Mineração de Texto são muito usadas para SRs baseados em conteúdo textual (LOPS; GEMMIS; SEMERARO, 2011), principalmente quando o contexto do SR trata de informações não-estruturadas. Mineração de Texto lida com análise de texto, suportando a sua natureza não-estruturada, imprecisa, incerta e difusa, para extração de informação e conhecimento (HOTH0; NÜRNBERGER; PAAß, 2005). Além disso, a área de Mineração de Texto utiliza de técnicas das áreas de Recuperação de Informação e Processamento de Linguagem Natural (PLN), conectando essas técnicas com algoritmos e métodos de Descoberta de Conhecimento em Banco de Dados, Mineração de Dados, Aprendizado de Máquina e Estatística (HOTH0; NÜRNBERGER; PAAß, 2005).

Feldman e Sanger (2006) apresentam uma arquitetura geral para aplicações de Mineração de Textos composta por quatro etapas: *tarefas de pré-processamento*, que preparam os dados para a central de operações de mineração; *central de operações de mineração*, que incluem algoritmos para a descoberta de padrões, tendências e conhecimentos por meio de técnicas e algoritmos; *componentes de apresentação*, que incluem interfaces para o usuário, apresentando visualizações dos conhecimentos gerados na etapa anterior; e *técnicas de refinamento*, também descritas como uma fase de pós-processamento, que incluem métodos para filtrar informações redundantes.

### A.1 Tarefas de pré-processamento

As tarefas de pré-processamento incluem rotinas, processos e métodos para a estruturação dos textos presentes nos documentos. A estruturação se faz necessária para a extração de informações e descoberta de conhecimento por meio de técnicas e algoritmos (HOTH0; NÜRNBERGER; PAAß, 2005).

#### A.1.1 Representação textual

Para a estruturação dos textos é necessário a definição da representação textual dos documentos. O vetor de termos, ou *Vector Space Model* (SALTON; WONG; YANG, 1975), é a representação clássica usada para representar documentos textuais (SEBASTIANI, 2002; LOPS; GEMMIS; SEMERARO, 2011). Cada dimensão desse vetor está associada a um termo, sendo que todas as dimensões representam todos os termos do conjunto de documentos. Formalmente, há um conjunto de documentos  $D = \{d_1, d_2, \dots, d_n\}$ , em que  $d_i$  representa um documento e  $n$  o número total de documentos, e um conjunto de termos  $\mathcal{T} = \{t_1, t_2, \dots, t_m\}$ , em que  $t_j$  representa um termo e  $m$  o número de termos presentes em todos os documentos. Representando a frequência de um termo pelo número de vezes



que  $t_j$  aparece em um documento  $d_i$ , denotado por  $ft(t_j, d_i)$ , o vetor de termos pode ser construído e representado da seguinte forma:  $\vec{v}_{t_{d_i}} = (TF(t_1, d_i), TF(t_2, d_i), \dots, TF(t_m, d_i))$ . Salton, Wong e Yang (1975) argumentam que a representação textual de documentos em vetor de termos é suficiente para separar documentos. Ao invés de frequência de termos, também é usado, a representação binária (SEBASTIANI, 2002), ou seja,  $t_j$  aparecendo em  $d_i$  corresponde à entrada 1 na dimensão  $j$  em  $\vec{v}_{t_{d_i}}$ . Há também outros métodos para representação textual, como *n-gramas* e *ontologias* (LOPS; GEMMIS; SEMERARO, 2011).

Ainda sobre o vetor de termos, Salton, Wong e Yang (1975) mostram com experimentos em diversos conjuntos de dados, que o uso da normalização nos vetores usando a técnica de Frequência de Termos-Frequência de Documentos Inversa (*Term Frequency-Inversed Document Frequency* – TF-IDF) é capaz de melhorar a separação de documentos:

$$\begin{aligned} TF-IDF(t_j, d_i) &= TF(t_j, d_i) \cdot IDF(t_j) \\ TF-IDF(t_j, d_i) &= TF(t_j, d_i) \cdot \left( \log_2 \frac{n}{DF(t_j) + 1} \right) \end{aligned} \quad (17)$$

em que  $IDF(t_j)$  representa a frequência de documentos inversa do termo  $t_j$ , e  $DF(t_j)$  a frequência de documentos que contém  $t_j$ . Essa normalização faz com que a frequência dos termos que aparecem em muitos documentos seja reduzida, e a frequência dos termos que aparecem em alguns raros documentos seja aumentada, com um fator de  $\log_2$ .

### A.1.2 Tokenização

Para realizar a estruturação de textos e representar os textos dos documentos em vetores de termos, o primeiro processo a ser realizado é a *tokenização*, que cria um dicionário de termos para cada documento através da quebra dos textos desses documentos. A quebra do texto pode ser feita através de caracteres delimitadores de termos, como espaços em branco, pontuações e etc. No entanto, existem casos que esses caracteres podem não ser delimitadores de termos, como por exemplo os termos *Prof.* e *Sr.*. Este problema é chamado de determinação de fim de sentença, e pode ser resolvido por métodos estáticos (*hard-coded*), baseados em regras e métodos de Aprendizado de Máquina (WEISS; INDURKHYA; ZHANG, 2010).

### A.1.3 Filtragem

Métodos de filtragem têm a função de retirar termos do conjunto  $\mathcal{T}$  que não contribuem para distinguir ou identificar documentos, como exemplo, conjunções (*e*, *pois*, *que*), artigos (*um*, *o*, *a*), preposições (*de*, *para*) e etc. A técnica de retirar determinados

termos de  $\mathcal{T}$  a partir de uma lista, é chamada de *stopwords*. Também são usadas outras técnicas, como a eliminação de termos com a frequência muito alta ou muito baixa.

#### A.1.4 Stemming

A fim de reduzir a ambiguidade de termos, o método de *stemming* é capaz de juntar, em uma única forma, termos relacionados (MINER et al., 2012). Por exemplo, o verbo *fazer* pode se apresentar em diversas formas, como *fazendo*, *fez*, etc. Esse processo é capaz de aumentar a capacidade da representação em distinguir ou identificar documentos, além de reduzir a dimensionalidade, reduzindo também a esparsidade.

#### A.1.5 Redução de Dimensionalidade

A representação em vetor de termos pode resultar em vetores esparsos num espaço de alta dimensão, que pode fazer com que algoritmos sofram do problema de *Maldição de Dimensionalidade*, que diz respeito à perda de densidade em espaços de alta dimensão, isto significa que medidas de distância se tornam incapazes de detectar padrões em um conjunto de dados (HAYKIN, 2008). Para amenização desse problema, são usados métodos de *redução de dimensionalidade*. A técnica mais comum de *redução de dimensionalidade* é chamada *Análise dos Componentes Principais* (*Principal Component Analysis - PCA*) (MURPHY, 2012). Esta técnica tem o objetivo de encontrar uma representação compacta através da descoberta de  $k$  vetores  $n$ -dimensionais ortogonais aos dados ( $\vec{v}$ ), em que  $k \leq m$ . Os vetores são encontrados a partir da minimização da projeção dos dados em  $\vec{v}$ . Depois de encontrados os vetores  $\vec{v}$ , é feita a projeção dos dados nesses vetores, resultando em uma representação num espaço mais compacto (HAN; KAMBER; PEI, 2011). É possível aplicar o algoritmo *PCA*, no vetor de termos, diminuindo a dimensionalidade e esparsidade, superando o problema de *Maldição de Dimensionalidade*.