

ANÁLISE FUZZY CLUSTER EM AMBIENTE R: UMA APLICAÇÃO DOS ALGORITMOS FANNY, CMEANSnE C-MEDOIDS NA CLASSIFICAÇÃO DOS MUNICÍPIOS BRASILEIROS SEGUNDO INDICADORES DE BEM-ESTAR-SOCIAL

Rodrigo Nunes Ferreira*

Resumo

O artigo demonstra o uso da análise fuzzy cluster através do uso de três algoritmos fanny, cmean e c-medoids na classificação dos municípios brasileiros segundo indicadores de bem estar- social. São utilizados índices de validação de agrupamentos para comparar os três algoritmos utilizados quanto a capacidade de gerar agrupamentos fuzzy com coesão interna e isolamento dos grupos formados. Como estratégia complementar de validação foi gerada uma análise discriminante que busca reproduzir o modelo de agrupamento gerado pela análise de cluster e assim medir o grau de confiabilidade, ou qualidade, do esquema classificatório. Todos os cálculos foram operacionalizados no ambiente R, fazendo uso dos recursos disponíveis na biblioteca do ambiente. Como indicador de bem estar social dos municípios brasileiros foram selecionados os sete indicadores que compõe o IDH-M acrescidos de três indicadores: percentual de pessoas acima da linha de pobreza, percentual de domicílios em condições adequadas de moradia e percentual de pessoas que demoram até trinta minutos no deslocamento casatrabalho. Os resultados mostraram melhores resultados para o algoritmo c-means, que não está implementado nos pacotes estatísticos mais utilizados em estudos socioeconômicos. O que mostra a importância do uso dos recursos do ambiente R em estudos desta natureza.

Palavras-chave: análise de agrupamento, fuzzy cluster, validação de agrupamentos, fanny, c-means, cmedoids

*Doutorando do Programa de Pós-Graduação em Geografia/IGC/UFMG.

1-INTRODUÇÃO

Em 1930 o matemático polonês Jan Lukasiewicz desenvolveu os princípios modernos da lógica nebulosa. Na década de 1950 A. Kaplan e H. F. Schott formalizam os princípios da álgebra dos conjuntos nebulosos (fuzzy), partindo do princípio de que os processos de decisão são marcados, muitas vezes, pela “incerteza” presente nos fatores a serem analisados, podendo afetar a exatidão das respostas e a coerência das ações a eles relacionados. A lógica fuzzy, ao utilizar a linguagem matemática - grau de pertinência - associada à representações linguísticas de agrupamentos, fundamenta modelagens como instrumentos de resolução de processos decisórios com fatores de incerteza. Nos conjuntos fuzzy as fronteiras não são nitidamente definidas, um elemento x_i poderá pertencer, com certo grau de pertinência, a conjuntos diferentes. O grau de pertencimento é dado por uma função “grupo-conceito”, que imputa, em relação a esse “grupo-conceito”, graus de pertencimento no intervalo $[0,1]$ a cada elemento x_i do universo de discurso. As funções de pertencimento (μ) podem ter forte conteúdo subjetivo em sua definição. Assim, um determinado “grupo-conceito” poderá ter diferentes funções de pertencimento; ou, uma função padrão (Jané,2004; Oliveira, 2003; Cruz, 2004)

A teoria dos conjuntos fuzzy foi proposta por Zadeh (1965), que desenvolveu o conceito de incerteza de pertença descrito por uma função de pertinência. As primeiras aplicações da teoria dos conjuntos fuzzy em análise de cluster surgiram nos trabalhos de Bellman, Kalaba e Zadeh (1966) e Ruspini (1969), que desenvolveram o algoritmo FANNY. A concepção de Dunn (1974), aprimorada por Bezdek (1981), permitiu o desenvolvimento do algoritmo Fuzzy C-Means, como uma versão fuzzy do algoritmo K-means de particionamento rígido apresentado por J. MacQueen (1967) (Yang, 1993). Atualmente existe uma grande número de algoritmos fuzzy cluster disponíveis na literatura, em sua maioria variações e aprimoramento dos algoritmos

inicialmente desenvolvidos na décadas de 1960 e a1970.

Sobre o uso da análise Fuzzy Cluster em estudos geográficos, Camara et. al. (2003, p. 94) avalia que a partir da década de 1990 o desenvolvimento da análise de sistemas geográficos deu grande ênfase no uso de técnicas de estatística espacial e lógica nebulosa (fuzzy). Burrough (1986)

destaca da seguinte forma a importância desta técnica para a modelagem dos sistemas geográficos:

Os limites desenhados em mapas temáticos (como solo, vegetação, o geologia) raramente são precisos e desenhá-los como linhas finas muitas vezes não representa adequadamente seu caráter. Assim, talvez não nos devamos preocupar tanto com localizações exatas e representações gráficas elegantes. (Burrough, 1986, apud Camara et. al., 2003).

O uso de algoritmos fuzzy para classificação é comum em estudos ambientais, trabalhos como os de Guerra e Caldas (2011) e Ortega (2001) ilustram possíveis usos de algoritmo fuzzy na classificação do uso e ocupação do solo a partir de imagens digitais. Já em estudos socioeconômicos o uso da classificação fuzzy é mais restrito, e poucos são os trabalhos publicados no Brasil que fazem uso desta técnica. Simões (2003) utiliza a técnica de conjuntos nebulosos para identificação de complexos industriais espaciais a partir de uma matriz de acessibilidade espacial. Miranda-Ribeiro e Garcia (2005; 2008) utilizaram a técnica de FuzzyCluster para investigar a relação entre segregação social e segregação espacial em Belo Horizonte, enquanto Nogueira, Garcia e Horta (2011) classificaram os países da América Latina quanto ao grau de urbanização a partir de uma matriz de particionamento fuzzy. Na mesma linha dos trabalhos anteriores, Tavares e Porto Junior (2010) utilizaram as técnicas de FuzzyCluster para classificar os municípios da região sul do Brasil de acordo com o nível de desenvolvimento.

Em comum a todos os trabalhos na área socioeconômica acima citados tem o fato de fazerem uso do algoritmo Fanny, implementado no S-PLUS. O fato de o pacote estatístico que está entre os mais utilizados nas ciências humanas, o SPSS, não disponibilizar um algoritmo de classificação fuzzy, pode ser um dos fatores a colaborar para o baixo uso desta técnica em estudos socioeconômicos. Assim como a preferência pelo algoritmo Fanny pode ser, em parte, explicada pelo fato de ser este o único algoritmo fuzzy cluster implementado em um pacote estatístico comercial e de interface gráfica, o S-PLUS. Esta constatação reforça a importância deste trabalho, ao demonstrar o uso de outros algoritmos Fuzzy Cluster mediante o uso dos recursos disponíveis na biblioteca do pacote R.

O R é classificado pela comunidade de usuários como uma linguagem e um ambiente de desenvolvimento integrado, para cálculos estatísticos e gráficos. É distribuído mediante licença GPL (Licença Pública Geral), e possui uma série de pacotes adicionais com funções específicas.

O uso do R na análise de cluster é facilitado pela disponibilidade na biblioteca do ambiente (CRAN) de uma série de pacotes com diversos algoritmos e técnicas de análise,

permitindo a exploração de diversas metodologias de agrupamento disponíveis na literatura, o que é inviável nos pacotes estatísticos comerciais.

2 – METODOLOGIA

A análise de agrupamentos, também conhecida como análise de conglomerados, classificação ou cluster, tem como objetivo dividir os elementos da amostra, ou população, em grupos de forma que os elementos pertencentes a um grupo sejam similares entre si com respeito às variáveis (características) que neles foram medidas, e os elementos em grupos diferentes sejam

heterogêneos em relação a estas mesmas características (MINGOTI, 2005, p. 155). Todo processo de classificação requer a escolha das unidades observacionais e das variáveis subscritas, que são, na visão de Faissol (1972, p. 78), as duas decisões arbitrárias fundamentais. Estas escolhas refletem o julgamento do investigador sobre os aspectos da realidade que são relevantes para o propósito da classificação desejada, sendo, portanto, uma caracterização inicial do dado, sem direcionamento matemático ou estatístico (Ferreira, Lima, 1979, p. 114). A complexidade do problema de agrupamentos advém da sua natureza não supervisionada, pois não se dispõe de um resultado final desejado (meta concreta a ser alcançada). Assim, a definição do que se entende por grupo (cluster) é carregada de elevado grau de subjetividade, e existe uma variedade de categorizações possíveis para um mesmo conjunto de dados, pois os objetos podem ser agrupados de diferentes maneiras dependendo da perspectiva (Vendramin, 2012).

Fuzzy Cluster é um processo de particionamento não-hierárquico, no qual busca-se a divisão do conjunto de entidades em um número de grupos homogêneos, com relação a uma medida de similaridade apropriada. Há uma diferença fundamental entre um agrupamento tradicional *hard* e um agrupamento *fuzzy*: enquanto no agrupamento *hard* o objeto pertence a um único cluster, no agrupamento *fuzzy* um objeto pode pertencer a mais de um grupo, mas com diferentes graus de pertinência (u). Portanto, o que define o agrupamento *fuzzy* é que sua aplicação a uma base de dados com n objetos $X = \{X_1, \dots, X_n\}$, tem como resultado final uma matriz de partição *fuzzy* desses objetos em um dado número de k de grupos, tal que $U = U_{ij}$ $k \times n$, no qual U é uma matriz de partição *fuzzy* $k \times n$ cujos elementos U_{ij} representam o grau de pertinência (ou simplesmente pertinência) do j -ésimo objeto ao i -ésimo grupo *fuzzy*.

Nos algoritmos tradicionais tipo *hard* assume-se que os objetos pertencem exclusivamente a um único grupo, mas em algumas situações esta hipótese pode não ser realista, pois as distribuições estatísticas possuem sobreposições, ou seja, as categorias se sobrepõem umas às outras (Vendramin, 2012). Os algoritmos de agrupamento difusos foram desenvolvidos para lidar exatamente com este tipo de situação. Atualmente existe uma variedade de algoritmos baseados em modificações e/ou extensões da formulação original do algoritmo mais conhecido, o Fuzzy *c*-Means. Esta diversidade de opções à disposição do pesquisador é ilustrada no trabalho de Vendramin (2012) que apresenta as principais características e potencialidades do ponto de vista processamento computacional de um total de 23 diferentes algoritmos para agrupamento fuzzy de dados.

2.1 - Algoritmos selecionados

Para os propósitos deste trabalho foram selecionados apenas três algoritmos para a classificação dos municípios brasileiros segundo um conjunto de indicadores de bem estar social: os dois mais tradicionais, *fanny* e *c-means*, e uma modificação do segundo, o *c-medoids*. A escolha deste último deve-se a tentativa de explorar a peculiaridade deste algoritmo para a classificação de municípios, pois, como discutido adiante, os grupos são formados tendo como centro médio não os valores médios do grupo segundo as *p*-variáveis selecionadas, mas um dos elementos do próprio grupo.

Na operacionalização no R foram utilizados os pacotes *cluster* para a classificação *fanny*, *fclust* para a *c-means* e *vegclust* para a *c-medoids*. O algoritmo *c-means* também está disponível nos pacotes *vegclust* e *e1071*, que produziram os mesmos resultados do pacote *fclust*. Também é possível encontrar o algoritmo *c-medoids* no pacote *fclust*, mas neste caso a função disponível produziu resultados inferiores aos encontrados via *vegclust*, uma análise da função de cálculo do *c-medoids* no *fclust* mostrou divergências em relação à notação do artigo de referência.

2.1.1 - Algoritmo FANNY

O método FANNY, desenvolvido nos trabalhos de Bellman, Kalaba e Zadeh (1966) e Ruspini (1969), define que para cada objeto *i* e cada cluster *v*, existe uma associação u_{iv} que indica o grau de pertencimento do objeto ao cluster. As associações são definidas através de processos iterativos, buscando a minimização da função objetivo *J*, que segundo notação de Kaufman e Rousseeuw (1990) é dada por:

$$J = \sum_{v=1}^k \frac{\sum_{i,j=1}^n u_{iv}^2 u_{jv}^2 d(i,j)}{2 \sum_{j=1}^n u_{jv}^2}, \quad (1)$$

na qual para cada elemento i e cada cluster v há uma pertinência, inicialmente desconhecida, u_{iv} , dada porque indica quão fortemente i pertence a v , se satisfeitas as condições de $0 \leq u_{iv} \leq 1$ para todo $i=(1,\dots,n)$, e $\sum_{i=1}^n u_{iv} = 1$ para todo $v=(1,\dots,k)$. Sendo $d(i,j)$ uma medida de dissimilaridade entre os objetos i e j .

2.1.2 - Algoritmo Fuzzy C-Means (FCM)

O algoritmo FCM, desenvolvido por Dunn (1974) e Bezdek (1981, 1984), é, segundo Jain et al. (1999), um dos mais facilmente compreendidos, bem documentados e intensamente utilizados algoritmo fuzzy cluster, e consiste numa versão fuzzy do algoritmo de agrupamentos rígido k-Means. No processo de partição do FCM o conjunto $X = \{x_1, x_2, \dots, x_n\}$ é dividido em p grupos, dado o grau de pertinência u_{ij} da amostra x_i ao j -ésimo grupo na matriz u . Via processos iterativos busca-se minimizar a função objetivo dada por:

$$J(U; c) = \sum_{i=1}^n \sum_{j=1}^p u_{ij}^m d(x_i, c_j)^2 \quad (2)$$

Sendo $\sum_{i=1}^n u_{ij} > 0$ para todo j , $\sum_{i=1}^n u_{ij} = 1$ para todo j , $m \in [1, \infty)$ é o parâmetro de fuzziness², c_j é o centro de um agrupamento fuzzy ($j=1,\dots,p$) e $d(x_i, c_j)$ a distância entre x_i e c_j . O centro do agrupamento j é definido como:

$$c_j = \frac{\sum_{i=1}^n u_{ij}^m x_i}{\sum_{i=1}^n u_{ij}^m}, \quad (3)$$

e a função de pertinência fuzzy u é dada por:

¹Na notação do pacote *cluster* do R (Maechler, 2013), o valor 2 do exponencial da função de pertinência u_{iv} , denominado *exponente de pertinência (membership exponent)* ou parâmetro de imprecisão (*parameter of fuzziness*), é substituído por uma constante r , que tem valor maior que 1 e por padrão valor 2. Aproximando o valor de r a 1, gera-se agrupamentos mais rígidos (*crisper clusterings*), com o aumento o valor de r em direção ao ∞ , aumenta-se o grau de fuzzyness (nebulosidade) do agrupamento, até atingir a completa fuzzyness, no qual todos os elementos pertencem igualmente a cada um dos agrupamentos gerados (quando todo $u_{iv}=1/k$). Nos casos em que pelo menos um elemento tem o valor máximo do grau de pertencimento (u_{iv}) igual em dois ou mais grupos, não é possível a conclusão do processo de classificação, via a desfuzzificação. Nestes casos é necessário reduzir o valor de r para aumentar o grau de separação entre os grupos, por isso a notação proposta por Maechler (2013) é mais adequada. Na sistematização feita Kaufman e Rousseeuw (1990, p. 191), os autores reconhecem a possibilidade de variação do valor do expoente r , semelhante a opção adotada pelos desenvolvedores do algoritmo *c-means*.

²Pal, James e Bezdek (1995, 370) sugerem que a melhor escolha para m é, provavelmente, no intervalo $[1,25; 2,5]$, sendo o ponto médio $m = 2$ a escolha preferida dos usuários do FCM.

$$u_{ij} = \frac{\left(\frac{1}{r(x_i, c_j)}\right)^{2/(m-1)}}{\sum_{k=1}^c \left(\frac{1}{r(x_i, c_k)}\right)^{2/(m-1)}} \quad (4)$$

A principal diferença do FCM em relação ao Fanny está na forma de se medir na função objetivo a dissimilaridade intra-grupos. Enquanto o segundo utiliza uma matriz D_{x_i, x_j} , ou seja, a distância entre os elementos que compõe o grupo, o segundo simplificar o cálculo e faz uso de uma matriz D_{x_i, c_j} , calculando apenas a distância entre os elementos do grupo e o centroide do agrupamento.

2.1.3 - Algoritmo Fuzzy C- Medoids - FCMdd

Krishnapuram et al. (1999) e Krishnapuram et al. (2001) criaram a versão Fuzzy do algoritmos de particionamento rígido PAM e CLARA (Kaufman e Rousseeuw, 1990). O algoritmo Fuzzy cmedoids (FCMdd) tem por base uma função objetivo para agrupamento relacional difuso (fuzzy) a partir da identificação de k objetos representativos em cada agrupamento (medoids) que minimizem a dissimilaridade dentro de cada grupo. A principal diferença em relação ao FCM é a seleção como centro do agrupamento não o ponto médio, mas de um elemento dentro do próprio conjunto de dados que minimize a dissimilaridade. O algoritmo foi proposto para aplicação em grandes bancos de dados voltados ao agrupamento de documentos eletrônicos, mas o princípio de seleção de um dos elementos como representativos do Grupo pode ter interessantes aplicações nos estudos socioeconômicos, pois facilita o entendimento e a ilustração das características do Grupo ao selecionar um dos elementos n como síntese das características do Grupo.

Seja $X = \{X_i | i=1, \dots, n\}$ um conjunto de n objetos, $r(X_i, X_j)$ expressa a distância entre o objeto X_i e o objeto X_j . Enquanto $V = \{v_1, v_2, \dots, v_c\}$, $v_i \in X$, representa um subconjunto de X com cardinalidade c , isto é, V é um c -subgrupo de X . Enquanto X^c representa o grupo de todos os subgrupos V de X . O algoritmo FCMdd é minimizado pela seguinte função:

$$J_m(V; X) = \sum_{j=1}^n \sum_{i=1}^c u_{ij}^m r(X_j, V_i), \quad (5)$$

onde a minimização é realizada para todo V em X^c . Na função (5) u_{ij}^m representa a matriz de pertinência fuzzy de X_j no cluster i , tendo m como expoente de pertinência (parâmetro

defuzziness) com valor no intervalo $[1, \infty)^3$. A função de pertinência u_{ij} é semelhante a utilizada noFCM, dada por:

$$u_{ij} = \frac{\left(\frac{1}{r(x_j, v_i)}\right)^{1/(m-1)}}{\sum_{k=1}^c \left(\frac{1}{r(x_j, v_k)}\right)^{1/(m-1)}}$$

onde $m [1, \infty]$ é o parâmetro de fuzziness.

2.1.4 - Medidas de Dissimilaridade

Para agrupar um conjunto de dados constituídos de n elementos com p -variáveis aleatórias em k grupos, é necessário definir uma medida síntese de dissimilaridade entre os objetos segundo as p variáveis, de tal forma que quanto menor o valor mais similares/próximos os elementos no plano p -dimensional. Nesta situação para cada elemento j tem-se o vetor de medidas X_j , definido por $X_j = [X_{1j}, X_{2j} \dots X_{pj}]'$, $j = 1, 2, \dots, n$, onde X_{1j} representa o valor observado da variável 1 medida no elemento j . Existem várias medidas para cálculo da dissimilaridade, e a escolha da medida tem impacto direto no resultado do agrupamento (Mingoti, 2005). Para operacionalização dos algoritmos selecionados será utilizada a Distância Euclidiana, a “mais popular medida para dados contínuos” (Jain et al., 1999, p. 271, tradução livre), que define a distância entre dois elementos X_i e X_j , $i \neq j$ e $K = (1, \dots, k)$, como:

$$d_2(x_i, x_j) = \left(\sum_{k=1}^d (x_{i,k} - x_{j,k})^2\right)^{1/2}. \quad (7)$$

2.2 - Base de dados utilizada

Todo processo de classificação requer “duas decisões arbitrárias fundamentais”: a definição das unidades observacionais e das variáveis subscritas (Faissol, 1972, p. 78). Estas escolhas refletem o julgamento do investigador sobre os aspectos da realidade que são relevantes para o propósito da classificação desejada, sendo, portanto, uma caracterização inicial do dado, sem direcionamento matemático ou estatístico (Ferreira, Lima, 1979, p. 114).

Neste trabalho propõe-se uma aplicação da técnica fuzzy cluster na classificação dos municípios brasileiros segundo indicadores de bem-estar social. Este enunciado já delimita a primeira das escolhas: as unidades observacionais serão formadas pelo conjunto dos municípios

³Os autores recomendam como padrão o valor de $m=1,5$.

brasileiros. Quanto aos indicadores de bem-estar-social a serem utilizados, faz-se necessário uma brevediscussão sobre o conceito em tela.

2.2.1 – Indicadores de bem-estar-social ou de qualidade de vida?

Segundo Ferrán Casas (1999) o debate sobre a questão da qualidade de vida, enquanto conceitooperacionalizável do ponto de vista da pesquisa acadêmica e das intervenções de políticaspúblicas, teve origem nos países desenvolvidos nas décadas de 1960 e 1970, coincidindo com acrise do Estado de bem-estar-social nestes países. Ainda segundo Casas (1999), o avanço dodebate neste período leva ao rompimento com a fundamentação conceitual centrada nascondições materiais de vida (bem-estar-social), sedimentando o entendimento da qualidade devida como impregnada de componentes subjetivos (psicossociais segundo o autor).

Joseph Stiglitz, Amartya Sen e Jean-Paul Fitoussi (2009), em relatório com as conclusões deuma comissão internacional de notáveis sobre como medir o progresso das sociedades,sintetizaram os desafios da mensuração do progresso social em 12 recomendações. Dentre elas asque reforçam a importância dos indicadores de qualidade de vida, que devem incluir indicadores subjetivos de satisfação individual. Como resumido por Veiga (2010, p. 47), para a Comissão “aquidade de vida só pode ser medida por um índice composto bem sofisticado, que incorporeaté mesmo as recentes descobertas desse novo ramo que é a economia da felicidade”. Trata-se de um reconhecimento explícito das limitações da abordagem objetivista baseada apenas nas escolhas feitas pelos indivíduos e que podem ser observadas, as denominadas preferencias observadas (Sen, 1986; IPEA, 2012).

A escolha do termo “indicadores de bem-estar-social” é um cuidado quanto ao uso do termomais popular, qualidade de vida, sem que algumas premissas sejam observadas. Portanto, concorda-se com a definição proposta por Casas (1999, p. 5), para quem o bem-estar-socialrepresenta as condições materiais, objetivamente observadas da qualidade de vida. Enquanto aquidade de vida é função do entorno material (bem estar social) e do entorno psicossocial (bemestar psicológico).

2.2.2 – Indicadores de bem-estar-social selecionados

Mesmo que não seja o interesse deste artigo explorar as possibilidades e limitações inerentes aoprocesso de escolha de um indicador de bem estar social adequado aos municípios brasileiros,convém uma breve exposição das orientações que embasaram a escolha dos indicadores para oteste aqui proposto.

Parte-se das pistas deixadas por dois importantes estudiosos sobre o desenvolvimento e urbanização brasileiros. Primeiramente o economista Celso Furtado, um dos maiores estudiosos da história social e econômica do Brasil, que em depoimento a uma Comissão do Congresso Nacional em 1999 expôs da seguinte maneira sua visão sobre a questão da pobreza no Brasil:

Podemos abordar o problema da pobreza de ângulos diferentes. Três são as dimensões que têm preocupado os estudiosos da matéria: 1) a questão da fome endêmica, que está presente, em graus diversos, em todo mundo; 2) a questão da habitação popular, que alguns países já encontraram solução; e 3) a questão da insuficiência de escolaridade, que contribuiu para perpetuar a pobreza (Furtado, 2002, p. 12).

Recentemente, após a onda de protestos que varreram o Brasil em junho de 2013, a urbanista Ermínia Maricato, em artigo com um sugestivo título “É a questão urbana, estúpido!”, manifestou sua visão sobre o fundamental da questão urbana brasileira:

As cidades são o principal local onde se dá a reprodução da força de trabalho. Nem toda melhoria das condições de vida é acessível com melhores salários ou com melhor distribuição de renda. Boas condições de vida dependem, frequentemente, de políticas públicas urbanas – transporte, moradia, saneamento, educação, saúde, lazer, iluminação pública, coleta de lixo, segurança. Ou seja, a cidade não fornece apenas o lugar, o suporte ou o chão para essa reprodução social. Suas características e até mesmo a forma como se realizam fazem a diferença. (...) A cidade constitui um grande patrimônio construído historicamente, mas sua apropriação é desigual e o nome do negócio é renda imobiliária ou localização, pois ela tem um preço devido a seus atributos. (Maricato, 2013, p. 6)

Como mensurar o bem estar social brasileiro considerando a perspectiva apontada pelos autores? Sem dúvida o mais conhecido e utilizado índice para comparação dos municípios brasileiros, o IDH, não é suficiente. É consenso, inclusive entre os próprios formuladores do índice, que o IDH tem seus limites enquanto indicador de progresso social (Sen, 1999). No contexto aqui proposto, o IDH não contempla, pelo menos diretamente, a questão urbana, seja do ponto de vista das condições de moradia ou do acesso aos serviços urbanos básicos⁴. A variável para o índice de renda, a renda per capita, também não contempla a questão da desigualdade na distribuição intragrupo da renda auferida (Gadrey e Jany-Catrice, 2006; Guimarães e Jannuzzi, 2011).

⁴Os indicadores que compõem o IDH-M 2013 são os sete primeiros do Quadro 1.

Entretanto, como não há espaço no âmbito deste trabalho para discussões sobre os indicadores alternativos, propõe-se partir dos indicadores do IDH e acrescentar três novos indicadores: a taxa de pessoas acima da linha de pobreza (inverso da taxa de pobreza); o percentual de domicílios em condições adequadas de moradia; e o percentual de pessoas que demoram até trinta minutos no deslocamento casa-trabalho (QUADRO 1). O primeiro visa corrigir a deficiência do indicador de renda per capita, apontando os locais onde, embora sejam elevadas as rendas médias, encontra-se elevado índice de pobreza⁵. O segundo visa sintetizar as condições de moradia no Brasil, trata-se de um indicador sugerido pela ONU para monitoramento da meta 11 dos Objetivos de Desenvolvimento do Milênio (ODMs), com metodologia definida para o Brasil no Relatório Nacional de Monitoramento dos ODMs (IPEA, 2007)⁶. Por fim, o terceiro indicador busca captar os custos de aglomeração dos grandes centros urbanos, fruto da expansão horizontal da área ocupada e da concentração das atividades geradoras de emprego nas áreas centrais, que implicam em tempos cada vez maiores de deslocamento no trajeto casa trabalho⁷.

⁵O indicador de pobreza exige a definição de uma linha de corte. A linha aqui utilizada considerou o critério do principal programa federal de transferência de renda, o Bolsa Família, que em 2010 considerava como pobres, e público alvo do Programa, todos aqueles que possuíam renda per capita de até R\$140,00. Para uma discussão sobre os resultados para o contexto brasileiro da mensuração da pobreza segundo diferentes linhas de corte ver IPEA (2010). Para um indicador alternativo considerando o patamar médio da renda local via correção da linha por uma

proxy do custo de vida nos municípios brasileiros ver Garcia e Matos (2007).

⁶ Na versão do indicador aqui adotado foram feitas pequenas alterações no indicador nacional, conforme metodologia adaptada para o monitoramento local dos ODMs no Município de Belo Horizonte (PBH, 2008). Considera-se como adequados o percentual de domicílios particulares permanentes que atendem simultaneamente

todas as seguintes condições: acesso a rede de esgoto geral ou fossa séptica, acesso a rede geral de água com canalização em pelo menos um cômodo, coleta de lixo por serviços de limpeza, iluminação via rede geral, existência de banheiro e média de até 3 moradores por cômodo servindo como dormitório.

⁷ Este mesmo indicador, mas um pouco menos restritivo considerando o tempo de deslocamento de até uma hora, foi utilizado no cálculo do índice de bem-estar urbano produzido por pesquisadores do Observatório da Metrópole, os autores do índice consideraram este indicador como “uma boa proxy das condições de deslocamento, apesar de não se referir a outros elementos importantes da mobilidade urbana, a exemplo da sua qualidade, da segurança dos serviços prestados e da infraestrutura disponível” (Ribeiro e Ribeiro, 2012, p. 19). Para exemplificar a importância dos diferenciais entre os municípios brasileiros registrados por este indicador, dados organizados por Matos e Ferreira (2013) mostram que se o tempo gasto no deslocamento casa-trabalho fosse remunerado resultaria em um acréscimo médio da remuneração de 22,8% no entorno metropolitano de SP e de 11,5% nas cidades médias, áreas nas quais gasta-se em média, respectivamente, 101 e 51 minutos por dia no deslocamento casa-trabalho-casa.

QUADRO 1 – Indicadores de bem-estar social selecionados

Sigla	Indicador	Fonte
ESPVIDA	Esperança de vida ao nascer	FJP, IPEA, PNUD. Atlas IDH 2012
T_FREQ5a6	Percentual da população de 5 a 6 anos de idade frequentando a escola	
T_FUND11a13	Percentual da população de 11 a 13 anos de idade frequentando os anos finais do fundamental ou que já concluiu o fundamental	
T_FUND15A17	Percentual da população de 15 a 17 anos com fundamental completo	
T_FUND18M	Percentual da população de 18 anos ou mais com fundamental completo	
T_MED18a20	Percentual da população de 18 a 20 anos de idade com o ensino médio completo	
RDPC	Renda per capita média	
PNPOB	Percentual de não pobres (renda per capita maior que R\$140,00)	
DOM_ADEQ	Percentual de Domicílios com Condições Adequadas de Moradia	IBGE. Censo 2010 Microdados da Amostra (elaboração própria)
DESL_ATE30	Percentual de pessoas que gastam até 30 minutos no deslocamento casa trabalho	

2.3 - Quantos grupos?

Um dos principais problemas na análise de cluster é a definição do melhor algoritmo e do número total de cluster k a utilizar. Embora nos métodos não-hierárquicos o número de clusters seja pré-definido, é necessário avaliar em que medida o total de agrupamentos selecionados se adequa ao agrupamento natural das variáveis selecionadas no espaço p -dimensional. Não existe uma resposta exata para esta questão, e na ausência de critérios estatísticos internos usado para inferência, encontra-se na literatura uma série de critérios ad hoc que podem auxiliar na decisão final (Mingoti, 2005; Hair Jr. Et al., 2005).

Uma abordagem utilizada para contornar este problema consiste em executar o algoritmo de agrupamento diversas vezes com matrizes de distância e de partições/protótipos iniciais diferentes e número de grupos variados e então escolher a partição mais adequada de acordo com um critério⁸. É possível utilizar diversas medidas como regra de parada (Stopping Rule) que

⁸G. Milligan e M. Cooper testaram um total de 30 medidas de validação de agrupamento visando determinar o número ideal de cluster, denominadas *stopping rule* (regra de parada), no intervalo [2, 5] para cada uma das medidas utilizadas em um processo de agrupamento hierárquico. Os resultados apontaram para os métodos de validação externa Calinski and Harabasz Index (PseudoF) e Duda and Hart index (Pseudo T2) como bons indicadores do

possibilitam a identificação do número ótimo de grupos a serem definidos para o conjunto de dados. A Figura 1 apresenta a aplicação no banco de dados utilizado, via a função `clValid` do pacote `clValid` disponível no R, de duas medidas para avaliar o grau de separação das partições: o Silhouette Width e o Dunn Index⁹. As medidas foram computadas para os resultados dos agrupamentos gerados por 10 diferentes métodos no intervalo de 2 a 8 partições geradas. Como ambos os índices devem ser maximizados, os resultados mostram melhores resultados do Silhouette Width para partições com menor número de grupos, enquanto para o Dunn Index, no geral, partições com maior número de grupos geraram melhores resultados, embora com maior variação entre os diversos algoritmos utilizados.

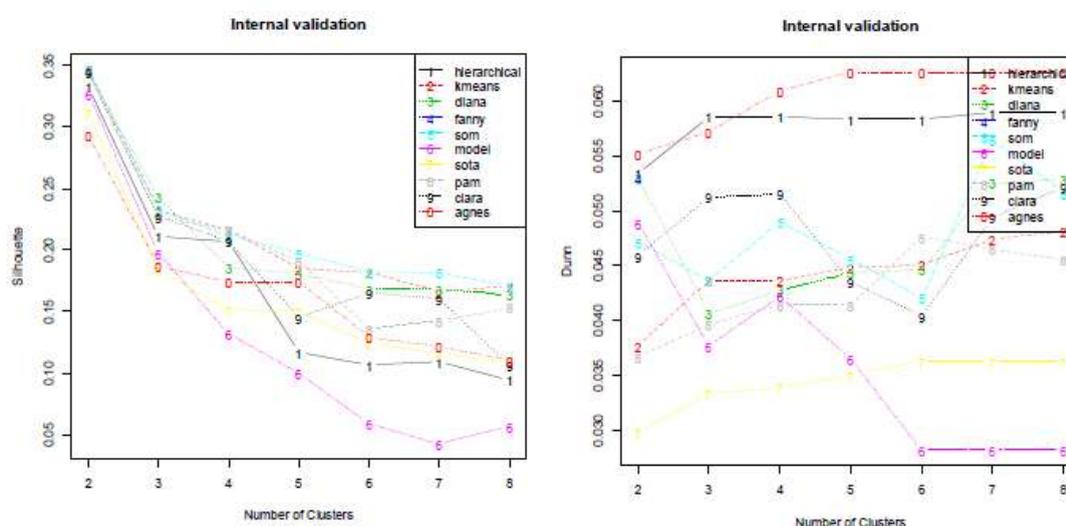


Figura 1: estimativa dos coeficientes de silhueta e de Dunn segundo valores de k no intervalo [2-8] para dez diferentes algoritmos.

Optou-se por utilizar uma classificação em três grupos. Considerando o tipo de observação utilizada no agrupamento, a geração de três grupos de municípios ilustra melhor as diferenças de nível de desenvolvimento dos municípios brasileiros, permitindo separar dois grupos com os melhores e os piores indicadores em dois polos, separados por um grupo intermediário, de transição.

3 - RESULTADOS

número de Grupos (Milligan e Cooper, 1985, p. 169; Mingotti, 2005).

⁹A *Silhouette Width* é a largura média da silhueta das observações (ver detalhamento no item 3, função (8)). O *Dunn Index* é a razão entre a menor distância entre as observações que não estão no mesmo cluster e a maior distância

$$D = \min_{i \in I_k, j \in I_{k'}} \|M_i^{(k)} - M_j^{(k')}\| / \max_{i, j \in I_k, i \neq j} \|M_i^{(k)} - M_j^{(k)}\|$$

entre as observações dentro do mesmo clusters, dada por

Tem valor entre 0 e o infinito, e deve ser maximizada (Brock et al., 2013).

Os resultados gerados pelo agrupamento podem ser visualizados nos gráficos da figura 2, que plotam o agrupamento num plano bidimensional tendo com eixos dois componentes que explicam 68% da variabilidade original da matriz de variáveis. Observa-se que o desenho do agrupamento intermediário foi o quem mais sofreu alterações entre os três algoritmos¹⁰.

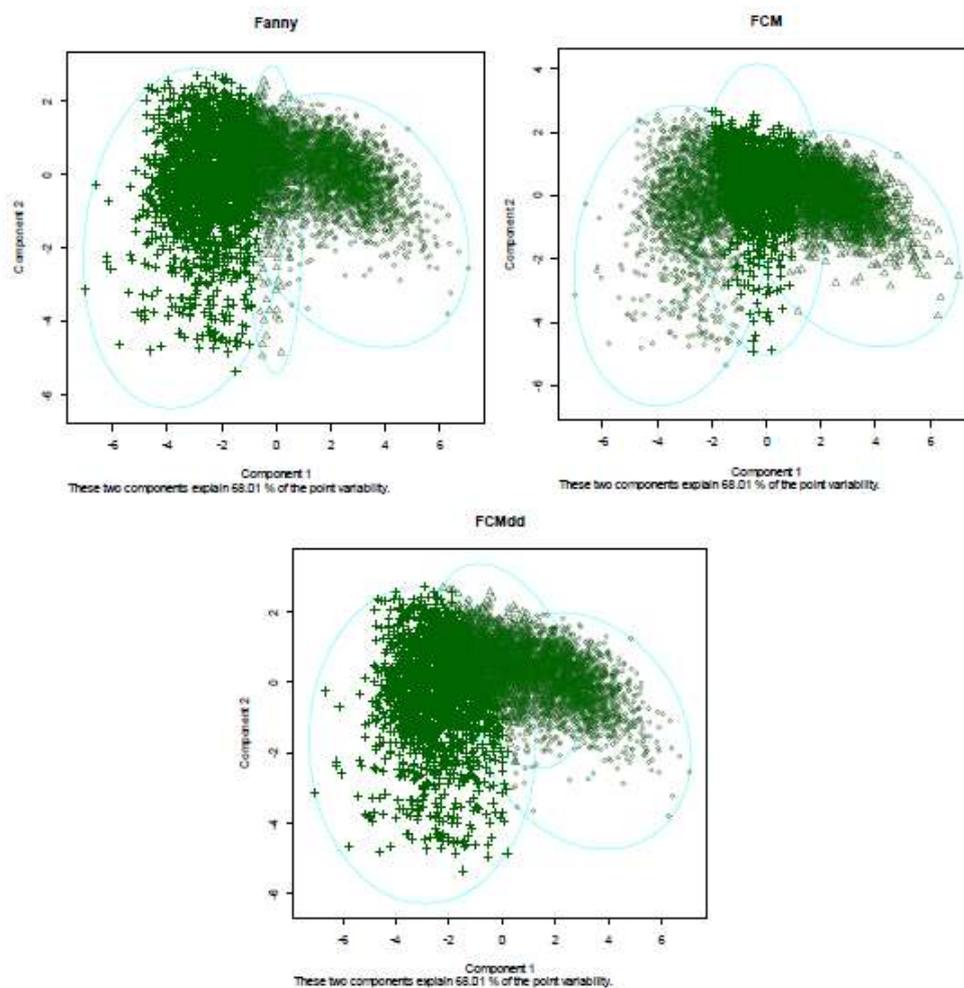


Figura 2: Gráficos Clusplot para resultados da classificação gerada pelo algoritmos selecionados.

O comportamento espacial das mudanças no grupo intermediário (cluster 2 no Fanny e FCMdd e 3 no FCM) pode ser observado na figura 3, que apresenta os resultados da classificação dos municípios brasileiros segundo as três técnicas utilizadas.

¹⁰Foi utilizado o valor do parâmetro de fuzziness $m=1,9$ em todos os algoritmos. O valor padrão 2 não foi utilizado pois no algoritmo fanny gerava casos de fuzzyness.

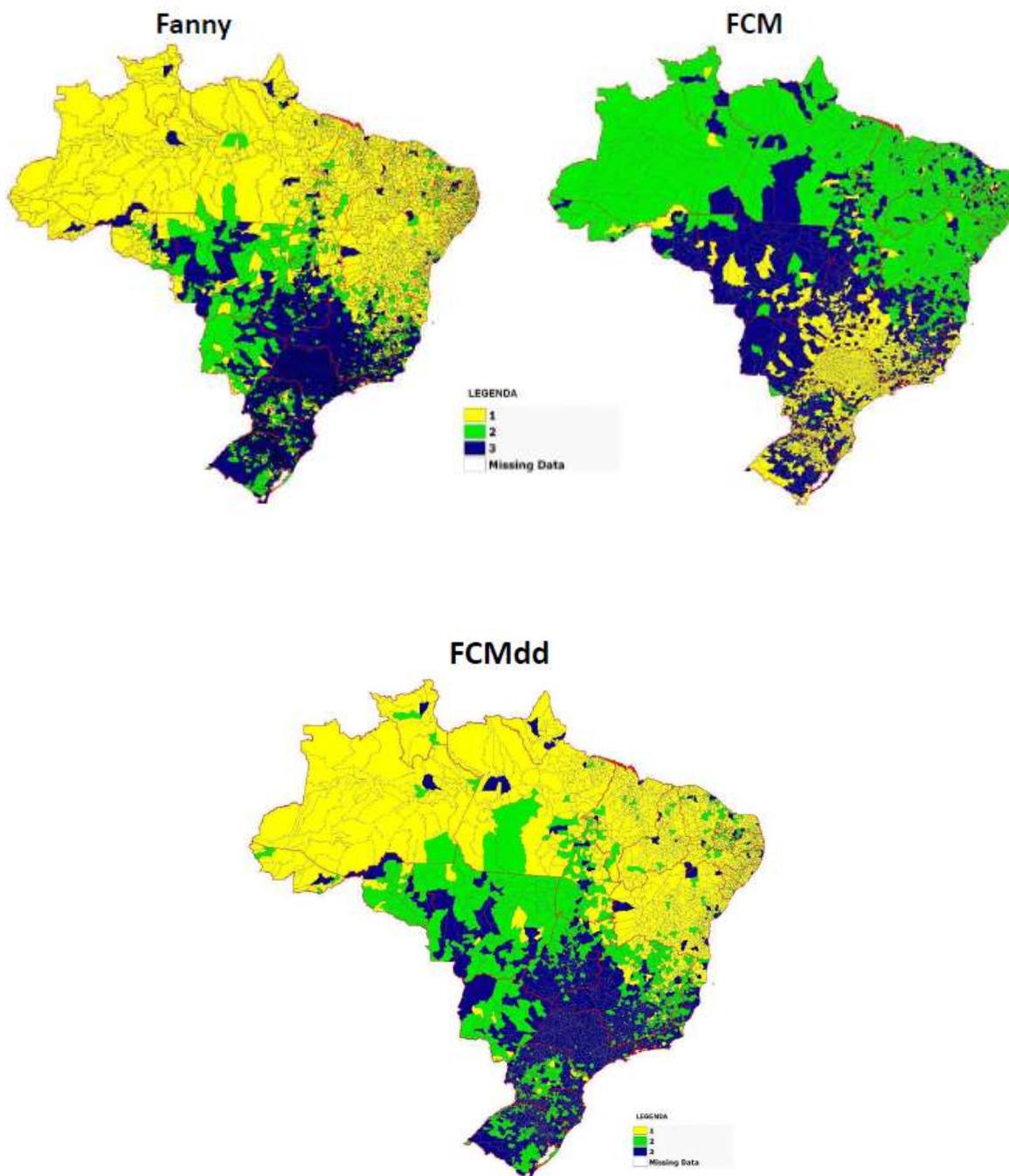


Figura 3: classificação dos municípios brasileiros quanto ao cluster de pertencimento segundo os algoritmos utilizados.

A classificação apresentada na Figura 1 pode ser interpretada a luz dos valores médios (centroides) de cada agrupamento para as variáveis selecionadas (Tabela 1). O FCMdd selecionou como centroides (medoids) os municípios de Cravolândia-BA (cluster 1), Presidente Kennedy-ES (cluster 2) e Pirai do Sul-RS (cluster 3).

Tabela 1: valores médios das variáveis (centroides) dos agrupamentos gerados pelos algoritmos selecionados.

Variável	Fanny			FCM			FCMdd		
	C1	C2	C3	C1	C2	C3	C1	C2	C3
ESPVIDA	70,7	73,8	75,3	75,3	70,4	73,7	70,9	73,5	73,9
T_FREQ5A6	90,8	87,8	93,8	94,0	92,0	90,1	93,1	89,8	91,0
T_FUND11A13	78,8	87,2	90,1	89,9	78,0	87,2	82,4	87,4	89,5
T_FUND15A17	41,6	56,9	67,2	67,2	40,1	57,9	42,0	55,8	62,5
T_FUND18M	31,8	38,4	47,9	48,9	31,2	39,6	33,8	36,9	46,1
T_MED18A20	23,9	35,8	48,4	48,5	23,0	37,7	27,9	38,0	43,6
RDPC	281,4	483,0	709,5	715,5	274,7	509,6	259,9	420,0	631,8
PNPOB	59,1	82,2	92,8	92,3	57,6	81,5	57,5	79,5	88,1
DOM_ADEQ	13,2	24,3	52,2	57,8	12,8	27,1	18,8	27,9	47,0
DESL_ATE30	80,6	82,0	80,4	80,5	80,1	82,3	78,2	85,1	78,4

Fonte: citadas no Quadro 1

Outra forma de visualizar graficamente o resultado do particionamento é através do gráfico desilhueta (Rousseeuw, 1987) com a medição da largura média da silhueta das observações dadas por $s(j)$. No processo de cálculo de $s(j)$ primeiro é obtido o valor individual de silhueta para cada observação $s(i)$, dado pela dissimilaridade média do objeto i para os demais objetos do Grupo A , indicado como $a(i)$, e a dissimilaridade média do objeto i para os objetos do grupo B , $d(i,B)$. Após computar $d(i,B)$ para todos os grupos $B \neq A$, o menor valor é selecionado, isto é, $b(i) = \min\{d(i,B), B \neq A\}$, que representa a dissimilaridade média do de i ao grupo mais próximo. Assim o valor $s(i)$ é dado por:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

sendo que a média dos valores de $s(i)$, dada $s(j) = \sum_{i=1}^n s(i)/n$, confiança do agrupamento, com valores no intervalo $[-1,1]$, quanto mais próximo de 1 melhor a qualidade do agrupamento e quanto mais próximo de -1 pior.

Os resultados de $s(i)$ e $s(j)$ apresentados na Figura 4 mostram melhores resultados de $s(j)$ para o método FCM ($s(j) = 0,228$), no qual apenas no cluster 1 é possível encontrar objetos com valores negativos de $s(i)$, que indica objetos que, na média, estão mais próximos de objetos externos ao grupo de classificação. Os algoritmos Fanny e FCMdd obtiveram, respectivamente, valores de $s(j) = 0,174$ e $s(j) = 0,208$. Outra observação interessante é a capacidade do FCM gerar grupos de tamanhos semelhantes, o que pode ser uma característica interessante do algoritmo

para a classificação em questão se esta vier acompanhada de bons resultados na análise de validação.

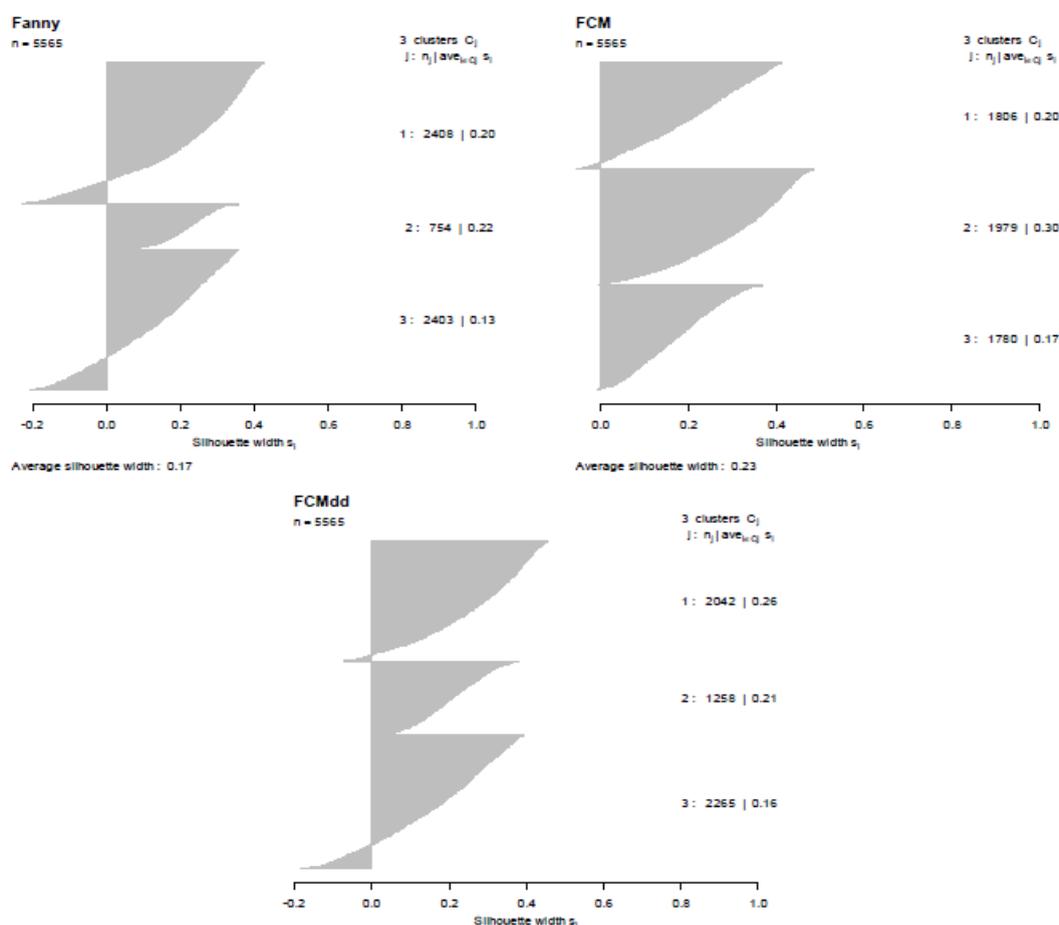


Figura 4: Gráfico de Silhueta para resultados da classificação gerada pelo algoritmosselecionados.

3.1 - Validação dos Agrupamentos

O objetivo de uma partição não-hierarquica é sempre produzir uma partição dos n elementos em k grupos de modo que a partição satisfaça dois requisitos básicos: coesão interna e isolamento dos grupos formados (Mingoti, 2005, p. 192). Segundo Mingoti (2005), quando os dados apresentam uma partição natural no espaço em relação as p -variáveis selecionadas, os resultados do agrupamento não serão muito diferentes entre os diversos métodos. Mas quando a partição natural apresenta grande interseção entre grupos, os resultados serão diferentes de acordo com o método e parâmetros escolhidos. Neste contexto o termo validação é utilizado para caracterizar, de forma ampla, os diferentes procedimentos para avaliar de maneira objetiva e quantitativa os resultados da análise de agrupamento. Os índices e critérios podem ser de três tipos (Vendraminet al., 2010; Campello, 2010; Desgraupes, 2013):

a. Externos: Avalia o grau de correspondência entre a estrutura de grupos (partição ou hierarquia) sob avaliação e informação a priori na forma de uma solução de agrupamento

esperada ou conhecida ou mesmo uma outra partição para comparação.

b. Internos: Avalia o grau de compatibilidade entre a estrutura de grupos sob avaliação e os dados, usando apenas os próprios dados

c. Relativos: Avaliam qual dentre duas ou mais estruturas de grupos é melhor sob algum

aspecto. É uma classe particular de critérios com habilidade para indicar qual a melhor dentre duas ou mais partições.

O Objetivo é sempre avaliar individualmente uma única partição e quantificar esta avaliação através de um valor que possa ser comparado relativamente. No clássico trabalho de Milligan e Cooper (1985) os autores compararam 30 diferentes medidas de validação de agrupamentos através de um extenso conjunto de experimentos envolvendo vários conjuntos de dados. Vendramin et al. (2010) testou 40 critérios de validação relativos com foco na complexidade computacional requerida no processamento. Desgraupes (2013) sistematizou e disponibilizou o pacote clusterCrit para R um total de 27 índices de validação interna ou relativos e outros 11 de validação externa. A existência na literatura de dezenas de tais critérios torna uma tarefa difícil para o usuário a escolha de um critério específico, e a alternativa sugerida tem sido o uso combinado de diversos critérios como forma de avaliar o agrupamento ideal.

A maioria dos índices referidos acima não fazem distinção entre o tipo de agrupamento hard ou fuzzy, pois não consideram em seus cálculos a matriz de pertinência U. Embora em menor número, também existem diversos critérios para avaliação de partições fuzzy (Campello, 2010). Vendramin (2012) apresenta uma proposta de agregação dos principais critérios de validação relativa de agrupamentos fuzzy em dois grandes grupos¹¹:

a. Índices baseados apenas na matriz de partição: avaliam a qualidade das soluções utilizando as informações apenas da matriz de partição U. Não utilizam informações sobre os atributos dos objetos e das dissimilaridade entre eles. Ex: Partition Coefficient (PC), Partition Entropy (PE) e Modified Partition Coefficient (MPC).

¹¹Embora Vendramin (2012) detalhe um total de 13 índices de validação de agrupamento fuzzy, optou por exemplificar em cada categoria apenas os índices a serem utilizados para análise subsequente. A escolha destes índices baseou-se em dois critérios: primeiro a disponibilidade de funções nos pacotes R utilizados que permitissem o cálculo, e segundo a viabilidade de uso das funções disponíveis para cálculo do índice na comparação de diferentes algoritmos. Algumas funções, da forma como foram implementadas nos pacotes R, permitem apenas a avaliação de resultados FCM.

b. Índices Baseados na Matriz de Dados: usam tanto a matriz de partição quanto a matriz de dados. EX: Xie-Beni (XB), Fuzzy Simplified Silhouette (FSS).

Para os testes de validação aqui propostos foram selecionados cinco índices disponíveis via função `fclust.index` do pacote `fclust` no R. O Quadro 2 detalha as fórmulas de cálculo e as referências de cada um dos índices utilizados. As seguintes notações são utilizadas no Quadro 2:

N é número de elementos, U_{ij} os valores da matriz de pertencimento, v_j o centro do cluster e k o número de clusters.

QUADRO 2 - Fórmulas de cálculo e referências dos índices de validação de agrupamentos fuzzy utilizados

Nome	Fórmula	Referência Bibliográfica
PC (Partition Coefficient)	$V_{PC} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^N (u_{ij})^2$	Bezdek, 1981; Dunn, 1974
MPC (Modified Partition Coefficient)	$V_{MPC} = 1 - \frac{k}{k-1} (1 - V_{PC})$	Dunn, 1974; Dave, 1996
PE (Partition Entropy)	$V_{PE} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^N u_{ij} \log_a(u_{ij})$	Bezdek, 1981
Xie Beni Index	$V_{XB} = \frac{\sum_{i=1}^k \sum_{j=1}^N (u_{ij})^m \ x_j - v_i\ ^2}{N \min_{t \neq s} \ v_t - v_s\ ^2}$	Xie & Beni, 1991;
FSS (Fuzzy Simplified Silhouette index)	$V_{FSS} = \frac{\sum_{j=1}^N (u_{pj} - u_{qj})^\alpha S_j}{\sum_{j=1}^N (u_{pj} - u_{qj})^\alpha}$ Sendo S_j a silhueta individual do objeto X_j conforme definição da função (8).	Campello & Hruschka, 2006

Fonte: adaptado de Vendramin (2012), Giordani e Ferraro (2013) e Meyer (2013)

O Quadro 3 apresenta os resultados dos índices selecionados. Observa-se que nos cinco critérios utilizados o FCM apresentou melhores resultados, seguido pelo FCMdd que não superou o Fanny apenas no Xie Beni Index.

QUADRO 3 – Resultados dos índices de validação de agrupamento fuzzy selecionados

Nome	Objetivo	Fanny	FCM	FCMdd
PC (Partition Coefficient)	max	0,353	0,519	0,372
MPC (Modified Partition Coefficient)	max	0,030	0,279	0,058
PE (Partition Entropy)*	min	1,069	0,818	1,041
Xie Beni Index**	min	0,806	0,680	1,188
FSS (Fuzzy Simplified Silhouette Index)***	max	0,271	0,370	0,333

Notas sobre as adaptações realizadas:

* utilizou-se exponencial de 1 como base do logaritmo. Para o FCMdd não é possível calcular o termo $\sum_{j=1}^k \frac{1}{N_j} \sum_{i=1}^n u_{ij}^2$ quando $N_j = 0$, ou seja, para os elementos tomados como centroides do agrupamento, que possuem $u_{ij}=1$ para o cluster de centro e $u_{ij}=0$ para os demais. Os elementos centroides de grupos foram excluídos do cálculo e utilizou-se $N-k$.

** a função objetivo fanny não utiliza a distância em relação ao centro médio do grupo, por isso foi necessário calcular os centros médios dos agrupamento gerados por este método.

*** utilizou-se $\alpha=1$

3.1.1 - Análise discriminante para validação de agrupamento

Como critério adicional para validação do agrupamento é possível implementar uma Análise Discriminante (AD) a partir do agrupamento gerado pela análise de cluster. Segundo Mingoti(2005) a AD é uma técnica que pode ser utilizada para classificação de elementos, mas requer a definição dos grupos, e de suas características, para os quais cada elemento pode ser classificado. O que permite a elaboração de uma função matemática, chamada regra de classificação ou de discriminação, que pode ser utilizada para classificação dos novos elementos.

Aqui propõe-se o uso da técnica de análise discriminante para mensuração do grau de ajuste da classificação, mas sem entrar em maiores detalhes sobre esta técnica de análise que o objetivo deste artigo. A técnica de validação adotada consiste em dividir o grupo previamente classificado pelos algoritmos selecionados em dois subgrupos, um para treinamento e elaboração da função de discriminação, e outro para aplicação da função elaborada para o primeiro grupo.

De forma semelhante à análise de cluster, na análise discriminante pretende-se explorar os dados de forma a maximizar a diferença entre os K grupos existentes. Para isso são criadas funções preditivas, combinações lineares das variáveis preditivas, que discriminam os K grupos maximizando a diferença entre eles. Assim, dado n elementos, K grupos e p variáveis preditivas $X=(X_1, \dots, X_p)$, o objetivo é estabelecer uma combinação linear dada por $Y_j = v_0 + \sum_{i=1}^p v_{ji} X_i$, onde Y é o vetor das funções lineares Y_j , v a matriz dos coeficientes lineares v_{ji} e X o vetor das variáveis preditivas (Mingoti, 2005).

Os resultados oferecem mais uma medida para avaliar o grau de confiabilidade, ou qualidade, do esquema classificatório da análise de cluster. A análise discriminante, quando comparada aos índices de validação apresentados no tópico anterior, oferece como vantagem o uso de uma medida de validação que vai além do uso apenas das distâncias matemáticas, inserindo uma regra de classificação fundamentada na teoria das probabilidades, com o objetivo de básico de derivar combinações lineares das p variáveis iniciais que maximizem a

diferenciação entre os grupos. A técnica tem como principais pressupostos a normalidade das variáveis, a homogeneidade das matrizes de variância-covariância e a existência de diferenças significativas entre os grupos (Mingoti, 2005; Ferreira e Lima, 1978).

Os resultados da AD aplicada aos três agrupamentos gerados são apresentados no Quadro 4, enquanto a Figura 5 ilustra a separação dos grupos pela primeira função discriminante. Os resultados da avaliação da reprodutibilidade do modelo de agrupamento para três clusters em que a previsão para o grupo é a classificação prevista pelas funções discriminantes, mostram que, embora os três algoritmos tenham permitido a geração de modelos com alto percentual de variância explicada pela primeira função, observa-se um maior grau de acerto para as funções geradas a partir do FCM (94,4%), confirmando os resultados das técnicas de validação apresentadas no item anterior.

Quadro 4: Síntese dos resultados da análise discriminante para os agrupamentos gerados pelos algoritmos selecionados

Medida	FANNY	FCM	FCMdd
% da variância explicada pela função 1	92.3%	91.9%	90.8%
% de classificação correta	85.6%	94.4%	91.9%

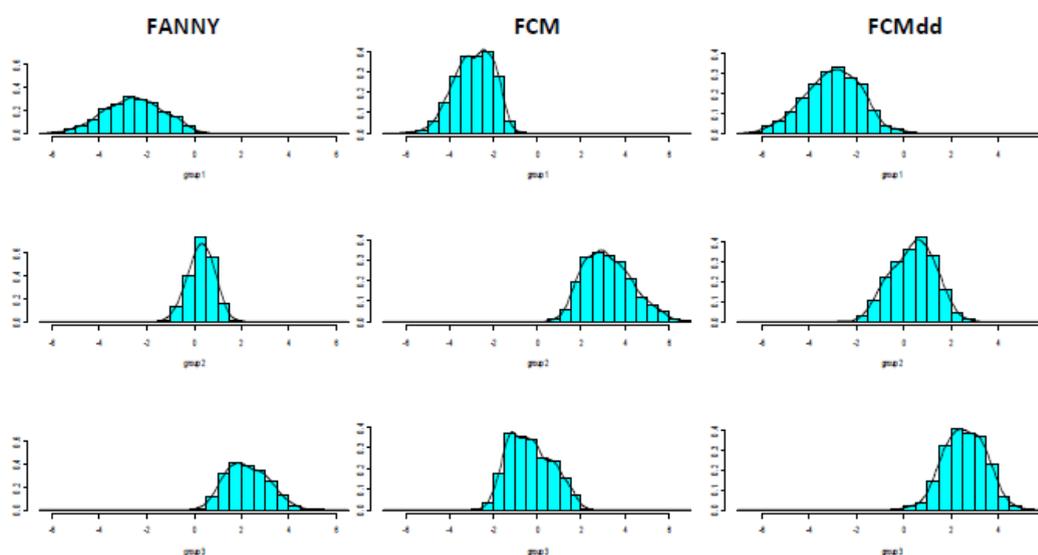


Figura 5: Histogramas da distribuição dos elementos dos grupos segundo scores da função discriminante 1 por algoritmo utilizado.

A AD também permite identificar as variáveis que mais contribuem para a discriminação dos grupos, através do índice de correlação linear das variáveis originais com os scores da função 1, conforme apresentado no Quadro 5. Observa-se que as variáveis _FREQ5A6 e DESL_ATE30 foram as que apresentaram os menores índices de correlação com a função discriminante de maior poder explicativo do agrupamento, um indicativo de que a inclusão do indicador de tempo de deslocamento pouco contribuiu com os resultados do agrupamento gerado.

Quadro 5: Coeficiente de correlação linear de Pearson entre as variáveis e os scores da função discriminante 1 por algoritmo utilizado

Variável	FANNY	FCM	FCMdd
ESPVIDA	0.877	-0.875	0.884
T_FREQ5A6	0.169	-0.155	0.095
T_FUND11A13	0.664	-0.672	0.666
T_FUND15A17	0.844	-0.843	0.833
T_FUND18M	0.765	-0.792	0.767
T_MED18A20	0.853	-0.845	0.825
RDPC	0.892	-0.885	0.876
PNPOB	0.965	-0.952	0.969
DOM_ADEQ	0.675	-0.712	0.667
DESL_ATE30	-0.002	0.028	-0.027

4 - CONCLUSÕES

Os resultados da análise fuzzy cluster, através do uso de três algoritmos fanny, c-means e cmedoid para classificação dos municípios brasileiros segundo indicadores de bem-estar social, mostraram melhores resultados do algoritmo FCM. O FCM obteve melhor desempenho em todas as medidas utilizadas para validação de agrupamento difuso quanto a coesão interna e ao isolamento dos grupos formados. O FCM também obteve bom desempenho nos resultados da avaliação ad hoc via análise discriminante, que foi capaz de reproduzir a classificação gerada pelo FCM com 94,5% de classificação correta. Como este algoritmo não está implementado nos

pacotes estatísticos mais utilizados estudos na área social, o SPSS e o S-PLUS, o trabalho aqui apresentado destaca a importância do uso do ambiente R. Que além de ampliar o leque de algoritmos disponíveis para a classificação fuzzy cluster, permite ao usuário acessar uma série de médias de validação de agrupamento e assim explorar toda a potencialidade desta técnica de classificação.

REFERÊNCIAS BIBLIOGRÁFICAS

Bellman , R., Kalaba, R., and Zadeh , L.A., (1966) Abstraction and pattern classification, I. *Math. Anal. and Appl.* 2, 581-586.

Bezdek, J. (1981). *Pattern Recognition with Fuzzy Objective Function Algorithms*. Kluwer Academic Publishers, Norwell, MA, USA.

Bezdek, J. C. (1981). *Pattern Recognition with Fuzzy Objective Function Algorithms*. Plenum Press.

Bezdek, James C.; Ehrlich, R.; Full, W. FCM: The fuzzy c-means clustering algorithm, *Computers & Geosciences*, Volume 10, Issues 2–3, 1984, Pages 191-203

Brock, G., Pihur, V., Datta , S., Datta, S, Package ‘clValid’. August 29, 2013, <http://cran.rproject.org/web/packages/clValid/index.html>

Brock, G., Pihur, V., Datta, S. and Datta, S. (2008) clValid: An R Package for Cluster Validation *Journal of Statistical Software* 25(4) <http://www.jstatsoft.org/v25/i04>

BURROUGH, P. *Principles of Geographical Information Systems for Land Resources Assessment*. Oxford, England, Oxford University Press, 1986.

CÂMARA, G.; MONTEIRO, Antônio Miguel; MEDEIROS, José Simeão de. Representações Computacionais do Espaço: Fundamentos Epistemológicos da Ciência da Geoinformação. *Geografia* (Rio Claro), Rio Claro, Brasil, v. 28, n.1, p. 83-96, 2003.

Campello, R. & E. Hruschka (2006). A fuzzy extension of the silhouette width criterion for cluster analysis. *Fuzzy Sets and Systems* 157 (21), 2858 - 2875.

Campello, R. J. G. B. Generalized External Indexes for Comparing Data Partitions with Overlapping Categories. *Pattern Recognition Letters*, Vol. 31, p. 966-975, 2010

CRUZ, A. J. de Oliveira. *Lógica Nebulosa*. UFRJ, 17 de junho, 2004

Dave, R. N. (1996). Validating fuzzy partitions obtained through c-shells clustering. *Pattern Recognition Letters* 17 (6), 613 - 623.

Desgraupes, Bernard. An R Package for Computing Clustering Quality Indices. University Paris Ouest, Lab Modal'X. April 2013, Vignettes do pacote clusterCrit disponível em: <http://cran.rproject.org/web/packages/clusterCrit/index.html>

Dunn, J. (1974). A Fuzzy Relative of the ISODATA Process and Its Use in Detecting Compact Well-Separated Clusters. *Journal of Cybernetics*, 3:32–57.

Faissol, Speridião. "Análise Fatorial: problemas e aplicações na geografia, especialmente nos estudos urbanos." *Revista Brasileira de Geografia* 34.4 (1972): 7-100.

Ferrán Casas. Calidad de vida y calidad humana. *Papeles del Psicólogo*, Noviembre, nº 74, 1999, disponível in: <http://www.papelesdelpsicologo.es/vernumero.asp?id=812>

FERREIRA, M. L.; LIMA, O.M.B. Processo de Classificação. in: Faissol, Speridião. *Tendências atuais na geografia urbano/regional: teorização e quantificação*. Fundação Instituto Brasileiro de Geografia e Estatística, 1978.

Furtado, Celso. *Em Busca de Novo Modelo*. São Paulo: Editora Paz e Terra, 2002.

GADREY, J.; JANY-CATRICE, F. *Os novos indicadores de riqueza*. São Paulo: Senac, 2006.

GARCIA, R. A.; MATOS, R. E. S. Mensurando a inserção social dos migrantes brasileiros. In *Anais do V encontro nacional sobre migrações*, Campinas, 2007.

Gath, I.; Geva, A. B. Unsupervised Optimal Fuzzy Clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 11, n. 7, p. 773-781, 1989.

Giordani, P, Ferraro, M. B. Package 'fclust', August 29, 2013, disponível em: <http://cran.rproject.org/web/packages/fclust/index.html>

GUIMARÃES, José Ribeiro Soares; JANNUZZI, Paulo M. IDH, Indicadores Sintéticos e suas aplicações em políticas públicas: uma análise crítica. *Revista Brasileira de Estudos Urbanos e Regionais*, v. 7, n. 1, p.73-90, 2011.

Hair Jr., J.F., Anderson, R.E., Tatham, R. L., Black, W.C. *Análise multivariada de dados*. trad. Adonai Schlup Sant'Anna e Anselmo Chaves Neto.- 5. ed.- Porto Alegre: Bookman, 2005.

IPEA - Instituto de Pesquisa Econômica Aplicada. *Objetivos de desenvolvimento do milênio: relatório nacional de acompanhamento*. Coordenação: IPEA e Secretaria de Planejamento e Investimentos Estratégicos. Supervisão: Grupo técnico para o acompanhamento dos ODM. Brasília: Ipea, MP, SPI, 2007. 152p

IPEA - Instituto de Pesquisa Econômica Aplicada. *PNAD 2009 – Primeiras Análises: Distribuição de Renda entre 1995 e 2009*. Comunicados do Ipea Nº 63, 05 de outubro de 2010, Disponível em: http://www.ipea.gov.br/portal/images/stories/PDFs/comunicado/101005_comunicadoipea63.pdf

IPEA - Instituto de Pesquisa Econômica Aplicada. 2012: Desenvolvimento Inclusivo Sustentável? Comunicados do Ipea Nº 158, 18 de dezembro de 2012, Disponível em: http://www.ipea.gov.br/portal/images/stories/PDFs/comunicado/121218_comunicadoipea158.pdf

Jain, A. K., Murty, M. N., and Flynn, P. J. (1999). Data Clustering: a Review. *ACM Comput. Surv.*,31(3):264–323.

JANÉ, D. de Almeida. Uma introdução ao estudo da lógica Fuzzy. *Hórus – Revista de Humanidades e Ciências Sociais Aplicadas, Ourinhos/SP*, no. 02, 2004.

Kaufman, L. and Rousseeuw, P.J. (1990). *Finding Groups in Data: An Introduction to Cluster Analysis*. New York: John Wiley & Sons, Inc.

Krishnapuram R., Joshi A., Nasraoui O., Yi L., 2001. Low-complexity fuzzy relational clustering algorithms for web mining. *IEEE Transactions on Fuzzy Systems*, 9, 595-607.

Krishnapuram, R., Joshi, A., & Yi, L. (1999). A Fuzzy relative of the k-medoids algorithm with application to web document and snippet clustering. *IEEE International Fuzzy Systems* (pp.1281–1286).

Maechler, Martin. Package ‘cluster’, Version 1.14.4, August 29, 2013, disponível em: <http://cran.rproject.org/web/packages/cluster/index.html>
MARICATO, Ermínia. É a questão urbana, estúpido!. *Le Monde Diplomatique Brasil*, ano 7, n. 73, p. 6-7, agosto de 2013.

MATOS, Ralfo; FERREIRA, Rodrigo Nunes. Espaço, população e economia; dos subespaços proeminentes ao transporte público premente. *Anais do II CONINTER - Congresso Internacional Interdisciplinar em Sociais e Humanidades*. Belo Horizonte, de 8 a 11 de outubro de 2013.

Meyer, David. Package ‘e1071’. August 29, 2013, disponível em: <http://cran.rproject.org/web/packages/e1071/index.html>

Milligan, G. W.; Cooper, M. C. An examination of procedures for determining the number of clusters in a data set, *Psychometrika* 50(2) (1985), 159–179.

Mingoti, S.A. *Análise de dados através de métodos de estatística multivariada – uma abordagem aplicada*. Belo Horizonte: Editora: UFMG, 2005. 295p.

MIRANDA-RIBEIRO, A. ; GARCIA, R. A. . Segregação socioespacial em Belo Horizonte: uma aplicação dos modelos difusos. *Geografias (UFMG)*, v. 1, p. 86-97, 2005.

NOGUEIRA, Marly; GARCIA, Ricardo Alexandrino; HORTA, Célio Augusto da Cunha; MOREIRA, Kely Cristina. A urbanização na América Latina e Caribe: um indicador de âmbito regional. *Revista Geográfica de América Central, Universidad de Costa Rica*, 2011. p. 1-18.

OLIVEIRA, Suzana Abreu de Souza. Alguns comentários sobre a teoria Fuzzy. *Exacta*, n. 1, abril, 2003, p. 139-147, Universidade Nove de Julho.

Pal, Nikhil R., and James C. Bezdek. "On cluster validity for the fuzzy c-means model." *Fuzzy Systems, IEEE Transactions on* 3.3 (1995): 370-379.

PBH - Prefeitura de Belo Horizonte (MG). Relatório de acompanhamento dos Objetivos de Desenvolvimento do Milênio Belo Horizonte – 2008. Secretaria Municipal de Planejamento, Orçamento e Informação -SMPL. Belo Horizonte: SMPL, 2008

Ribeiro, Luiz Cesar de Queiroz, Ribeiro, Marcelo Gomes (Orgs). IBEU: índice de bem-estar urbano. 1.ed. - Rio de Janeiro: Letra Capital, 2013.

Rousseeuw, Peter J. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics* 20 (1987): 53-65.

Ruspini, E.H., (1969) A new approach to clustering, *Inform. and Control* 15, 22-32. Sen, Amartya. (1995) Rationality and social choice. *American Economic Review* 85 (1).

Sen, Amartya. Evaluación del Desarrollo Humano Contribución Especial. PNUD - Programa de las Naciones Unidas para el Desarrollo. Informe sobre Desarrollo Humano 1999. Madrid: Ediciones Mundi-Prensa, 1999.

SIMÕES, R. F. Complexos industriais no espaço: uma análise de fuzzy cluster. TD 209. Belo Horizonte: UFMG/Cedeplar, 2003.

STIGLITZ, J.; SEN, A e FITOUI, Report by the Commission on the Measurement of Economic Performance and Social Progress, 2009. September, 2009.

VEIGA, José Eli da. Indicadores de sustentabilidade. *Estudos Avançados*, São Paulo, v. 24, n. 68, 2010.

Vendramin, L., Campello, R. J. G. B. & Hruschka, E. R. "Relative Clustering Validity Criteria: A Comparative Overview" *Statistical Analysis and Data Mining*, Vol. 3, p. 209-235, 2010

Vendramin, Lucas. Estudo e desenvolvimento de algoritmos para agrupamento fuzzy de dados em cenários centralizados e distribuídos. Dissertação (Mestrado - Programa de Pós-Graduação em Ciências de Computação e Matemática Computacional) Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo. São Carlos, 2012 138p.

YANG, M.S. Survey of Fuzzy Clustering. *Mathl. Comput. Modelling* Vol. 18, No. 11, pp. 1-16, 1993

Xie, L. X.; Beni, G.. Validity measure for fuzzy clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 3, n. 8, p. 841-847, 1991.

Zadeh, L.A. (1965) Fuzzy sets, *Inform. and Control* 8, 338-353.