

MINISTÉRIO DA EDUCAÇÃO  
Universidade Federal de Alfenas - UNIFAL-MG

## Projeto de Pesquisa

# **Classificação dos municípios brasileiros segundo gastos previdenciários utilizando componentes principais e análise de agrupamento**

Autor: Larissa Gonçalves Souza

Orientadora: Patrícia de Siqueira Ramos

Coorientador: Reinaldo Antônio Gomes Marques

Alfenas - MG

2015

## Sumário

1	Introdução . . . . .	3
2	Objetivos . . . . .	4
2.1	Objetivo primário . . . . .	4
2.2	Objetivos secundários . . . . .	4
3	Revisão de Literatura . . . . .	4
3.1	A transição demográfica no Brasil . . . . .	4
3.2	A Previdência Social . . . . .	5
3.3	Análise multivariada . . . . .	6
3.4	Análise de componentes principais . . . . .	9
3.5	Análise de agrupamento ( <i>Cluster Analysis</i> ) . . . . .	12
3.5.1	Número de grupos . . . . .	15
4	Metodologia . . . . .	17
4.1	Bases de dados e variáveis do estudo . . . . .	17
5	Resultados Esperados . . . . .	19
6	Cronograma . . . . .	19
7	Disciplinas necessárias . . . . .	20
8	Referências Bibliográficas . . . . .	20

## 1 Introdução

As aposentadorias e os outros benefícios concedidos pela previdência social formam grande parte da renda dos municípios brasileiros, influenciando fortemente sua economia e reduzindo desigualdades sociais (FRANÇA, 2004). Entretanto, o impacto dessa política pública não é uniforme, pois eles são bastante diversos entre si em relação à estrutura etária, à taxa de fecundidade, às atividades econômicas predominantes etc. O objetivo da presente pesquisa é construir uma tipologia dos municípios de acordo com a sua relação com a previdência social.

Em 3.546 dos 5.561 municípios, em 2004, o pagamento de benefícios superou o FPM (Fundo de Participação dos Municípios), como era o caso das quatro melhores cidades em desenvolvimento humano no Brasil: São Caetano do Sul - SP, em que o pagamento de benefícios foi 27,52 vezes maior do que o FPM; Águas de São Pedro - SP, 2,61 vezes; Niterói - RJ, 38,39; e Florianópolis - SC, 13,95. Além disso, o recebimento desses benefícios evita o êxodo para grandes cidades, fixando as pessoas em seus municípios de origem (FRANÇA, 2004).

Ainda segundo França (2004), nos municípios de até 5 mil habitantes, os benefícios representavam 20,3% da renda das famílias, sendo que a média brasileira era de 7,2%. Para a população rural, a previdência social é uma forma de redistribuição de renda, já que essas pessoas dificilmente contribuíram diretamente para a Previdência Social, fazendo com que este seja como um programa de renda mínima para os idosos no país. O acesso a esses benefícios melhora de forma significativa a qualidade de vida dos domicílios.

Atualmente há razoável disponibilidade de dados municipais no Brasil, seja através das pesquisas do Instituto Brasileiro de Geografia e Estatística (IBGE) (em especial, os Censos Demográficos e o Perfil dos Estados e dos Municípios Brasileiros) seja através dos registros administrativos de diversos órgãos públicos (Ministério da Saúde, Ministério da Educação, Ministério da Previdência Social etc.). Porém, na maior parte das vezes, as variáveis disponíveis nesses bancos de dados são analisadas separadamente e uma visão geral fica comprometida. Por isso, a análise multivariada é fundamental na análise dos dados socioeconômicos municipais, uma vez que a realidade municipal é multidimensional e muitas dessas variáveis estão intimamente relacionadas.

## **2 Objetivos**

### **2.1 Objetivo primário**

O objetivo deste trabalho é propor uma classificação dos municípios brasileiros em relação aos seus gastos com previdência social e a algumas variáveis demográficas.

### **2.2 Objetivos secundários**

- Primeiramente, será utilizada a análise de componentes principais para reduzir a dimensionalidade dos dados.
- Em seguida, será realizada a análise de agrupamento. O que se espera é que os grupos de municípios obtidos apresentem grande homogeneidade interna e grande heterogeneidade externa em relação às variáveis analisadas. Serão utilizados diferentes métodos de agrupamento e os resultados obtidos serão comparados.
- Serão analisados diferentes critérios para definição do número de grupos na partição final da análise de agrupamento.

## **3 Revisão de Literatura**

### **3.1 A transição demográfica no Brasil**

O Brasil vivenciou grandes mudanças nas últimas décadas no que diz respeito à sua dinâmica demográfica. Em torno de 1960, a população crescia num ritmo acelerado, sendo um país jovem, enquanto que, cerca de cinquenta anos depois, vivencia uma desaceleração desse crescimento. A idade mediana era 18 anos em 1960 e passou para 27 anos em 2010. Isso é apenas uma indicação de fenômenos mais gerais que têm sido observados: uma brusca diminuição das taxas de fecundidade e mortalidade em todas as idades, o envelhecimento da

população, novos arranjos familiares se formando, modificações na magnitude e limites etários da população economicamente ativa, além de outras mudanças (VASCONCELOS; GOMES, 2012).

A partir da década de 1970, houve uma rápida diminuição da fecundidade acompanhada da sustentada queda de mortalidade no país. A taxa de fecundidade total passou de 6,2 filhos/mulher em 1950 para 1,7 em 2012, atingindo níveis inferiores do que o que garantiria a reposição da população, que é de 2,1 filhos/mulher. Outra grande mudança ocorreu com a esperança de vida ao nascer, que era 45,4 anos em 1950, e hoje é 75,2 anos. Os dois movimentos ocorreram em um espaço de tempo muito curto no Brasil e em muitos dos outros países em desenvolvimento (CAMARANO, 2014).

Ainda segundo Camarano (2014), como consequências dessa transição demográfica que o Brasil vivencia, têm-se uma diminuição da taxa de crescimento da população e da força de trabalho, além do envelhecimento da população. Este último aspecto é visto como um problema, pois o crescimento rápido de um segmento populacional não produtivo e o menor crescimento do segmento produtivo podem desequilibrar a divisão de recursos na sociedade, gerando sérios problemas econômicos e previdenciários.

Entretanto, essas transformações na população brasileira não ocorreram de forma uniforme em todas as regiões do país. Em 1970, as regiões Norte e Nordeste ainda apresentavam valores altos de mortalidade infantil e de número médio de filhos por mulher, enquanto as regiões Sudeste, Sul e Centro-Oeste já apresentavam queda nesses índices. Apesar da queda da taxa de mortalidade infantil ter tido diferentes ritmos nas cinco regiões, em todas houve queda de cerca de 70%, entre 1980 e 2010, o que é espantoso (VASCONCELOS; GOMES, 2012).

Além da redução dos níveis de mortalidade, houve uma grande diminuição nos níveis de fecundidade. Em 2000, apenas a região Norte apresentava número médio superior a 3,0 filhos/mulher. Já em 2010, todas as outras regiões apresentavam níveis de fecundidade menores do que o nível de reposição de 2,1 filhos por mulher (VASCONCELOS, 2012).

### **3.2 A Previdência Social**

Uma das políticas públicas mais afetadas por essas mudanças demográficas é a Previdência Social. O sistema previdenciário brasileiro é formado por dois subsistemas, a Pre-

vidência Social básica e a Previdência Privada. A primeira é fornecida pelo poder público e formada pelos Regimes Próprios de Previdência Social (RPPS) - onde estão os trabalhadores do setor público - e pelo Regime Geral de Previdência Social (RGPS) - onde estão alocados os trabalhadores do setor privado. A Previdência Privada, de caráter facultativo e complementar, é composta pelas Empresas Abertas de Previdência Complementar (EAPC) e pelas Empresas Fechadas de Previdência Complementar (EFPC) (REIS; SILVEIRA; BRAGA, 2013).

O RGPS, composto por um número muito maior de pessoas, é o principal regime previdenciário do Brasil. Seus benefícios são um tipo de seguro social para o trabalhador e sua família, repondo sua renda quando este perde a capacidade de trabalho, devido a doença, invalidez, idade avançada, morte, desemprego involuntário, maternidade e reclusão. Há três grupos de benefícios concedidos pelo RGPS: previdenciários, acidentários e assistenciais. De acordo com o Ministério da Previdência Social (MPS), no ano de 2009, a Previdência Social emitiu 224,8 bilhões de reais em benefícios para 28,1 milhões de pessoas, o que representou 7,15% do PIB nacional (REIS; SILVEIRA; BRAGA, 2013).

### **3.3 Análise multivariada**

Os dados levantados em uma pesquisa são considerados multivariados quando os valores referentes a cada unidade amostral ou observação se referem a diversas variáveis aleatórias ao mesmo tempo, levando cada observação a ser multidimensional. Na maioria das pesquisas, os dados são multivariados mas, muitas vezes, o pesquisador opta por analisar cada variável separadamente. Porém, em geral, as variáveis são correlacionadas entre si e, quanto maior o número de variáveis, mais complexa se torna a análise univariada. Ao se utilizar a análise multivariada, as variáveis são analisadas ao mesmo tempo, fornecendo uma avaliação muito mais ampla do conjunto de dados, encontrando-se padrões e levando-se em conta a correlação entre as variáveis (MINGOTI, 2005).

A representação de dados multivariados se dá como em planilhas eletrônicas. Se há uma amostra aleatória de tamanho  $n$  e, para cada unidade amostral ou observação, os valores de  $p$  variáveis foram observados, cria-se uma matriz de dados  $X$  com dimensão  $n$  (linhas) por  $p$  (colunas):

$$\mathbf{X}_{n \times p} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}, \quad (3.1)$$

em que cada unidade amostral é representada por uma linha da matriz de dados  $\mathbf{X}$ , sendo um vetor com  $p$  elementos (variáveis), e cada variável é representada por uma coluna de  $\mathbf{X}$ , sendo um vetor com  $n$  elementos, as observações (EVERITT; HOTHORN, 2011).

A obtenção da matriz de dados a partir de uma amostra aleatória, como expressa na forma da equação (3.1), pode não ser muito informativa, principalmente se o tamanho amostral  $n$  for grande e houver um número excessivo de variáveis  $p$ . Torna-se interessante utilizar medidas resumo dos dados amostrais, da mesma forma que é feito no caso univariado, calculando-se a média, mediana, desvio padrão etc., de forma a sintetizar os dados da amostra obtida (FERREIRA, 2011).

Uma medida de tendência central muito utilizada é a média amostral que, no caso multivariado, torna-se o vetor de médias amostral de dimensão  $p \times 1$ , em que cada elemento é a média de cada variável:

$$\bar{\mathbf{X}} = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \\ \vdots \\ \bar{X}_p \end{bmatrix}.$$

Para medir a dispersão dos dados, no lugar da variância amostral, utiliza-se a matriz de covariâncias amostral  $\mathbf{S}$  de dimensão  $p \times p$ . Sua diagonal principal é composta pelas variâncias das  $p$  variáveis e os elementos fora da diagonal são as covariâncias entre as variáveis. Essa matriz é simétrica, ou seja,  $S_{ij} = S_{ji}$ .

$$\mathbf{S} = \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ S_{21} & S_{22} & \dots & S_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ S_{p1} & S_{p2} & \dots & S_{pp} \end{bmatrix}.$$

A correlação é também uma medida de covariação entre duas variáveis, porém em uma escala padronizada, ou seja, seus valores variam entre  $-1$  e  $+1$ . Valores próximos de  $+1$  indicam que as variáveis estão fortemente correlacionadas de forma positiva, grandes valores de uma estão associados a grandes valores da outra. Já valores próximos de  $-1$  indicam que as variáveis estão fortemente correlacionadas de forma negativa, indicando que grandes valores de uma estão associados a pequenos valores da outra. A matriz de correlações amostral é dada por

$$\mathbf{R} = \begin{bmatrix} 1 & r_{12} & \dots & r_{1p} \\ r_{21} & 2 & \dots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{p1} & r_{p2} & \dots & 1 \end{bmatrix}, \quad (3.2)$$

em que  $r_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}S_{jj}}}$  (FERREIRA, 2011).

Outras estatísticas descritivas, como a matriz de somas de quadrados e produtos, podem ser consideradas, dependendo do objetivo da pesquisa (FERREIRA, 2011).

Segundo Mingoti (2005), a análise multivariada se divide em dois grupos principais: técnicas exploratórias e técnicas de inferência estatística, como também ocorre na análise univariada. O primeiro possui um grande apelo prático por não dependerem do conhecimento da forma matemática da distribuição de probabilidade que gerou os dados amostrais e permitem a detecção de padrões. Exemplos de técnicas desse tipo são análise de componentes principais, análise fatorial exploratória, análise de agrupamento (*clusters*), entre outras. O foco do segundo grupo de técnicas é a estimação de parâmetros, testes de hipóteses, análise de regressão multivariada etc., cujo objetivo é utilizar a amostra para realizar inferências sobre a população de onde essa amostra foi extraída.

As técnicas exploratórias são muitas vezes denominadas técnicas de sintetização por se concentrarem em condensar uma grande massa de dados em uma forma mais simples. Assim, há uma redução significativa do volume de dados envolvido na análise ou uma redução da dimensionalidade (BARTHOLOMEW et al., 2008).

A presente proposta empregará técnicas exploratórias. A análise de componentes principais será usada como forma de reduzir a dimensionalidade dos dados, simplificando a sua estrutura de covariâncias antes de aplicar a análise de agrupamento (*clusters*) que ajudará a identificar os grupos de municípios com perfis similares quanto à presença da previdência social. As duas técnicas são apresentadas a seguir.



### 3.4 Análise de componentes principais

O objetivo da técnica de análise de componentes principais (ACP) é explicar a estrutura de covariâncias das  $p$  variáveis,  $\mathbf{X}^T = [X_1 \ X_2 \ \dots \ X_p]$ , por meio da construção de combinações lineares das variáveis originais. Os componentes principais são as  $p$  combinações lineares obtidas,  $\mathbf{Y}^T = [Y_1 \ Y_2 \ \dots \ Y_p]$ , e são não correlacionados entre si. Entretanto, como a intenção é reduzir o número de variáveis, a informação contida nelas é substituída pela informação contida em  $k$  componentes principais, em que  $k < p$ . Os  $k$  componentes são ordenados de forma que os primeiros deles já contabilizem a maior parte da variação presente em todas as variáveis originais (MINGOTI, 2005; EVERITT; HOTHORN, 2011).

A análise de componentes principais é uma técnica principalmente exploratória. Há métodos inferenciais para se testar hipóteses sobre componentes principais populacionais a partir de uma amostra aleatória de observações, mas eles são menos frequentes na literatura especializada (EVERITT; HOTHORN, 2011).

É preciso adotar um critério para reter apenas parte dos componentes, de maneira que grande parte da variância total seja explicada pelo conjunto pequeno de novas variáveis. Se o valor de  $k$  for pequeno e a quantidade de variação explicada pelos  $k$  componentes for grande, haverá uma simplificação da estrutura de covariâncias das variáveis originais. Essa técnica pode, então, ser utilizada como uma etapa intermediária para auxiliar em outras técnicas, como em problemas de multicolinearidade em regressão linear, por exemplo (FERREIRA, 2011).

A suposição de normalidade das  $p$  variáveis não é imprescindível para a aplicação da técnica, mas, se ocorrer, os componentes principais obtidos são, além de não correlacionados, independentes e normais. Os componentes podem ser obtidos a partir da matriz de covariâncias ou a partir da matriz de correlações das variáveis originais. Essa é uma questão discutida por alguns autores. Em geral, recomenda-se obter os componentes a partir da matriz de covariâncias amostral quando as variáveis estão na mesma escala e a partir da matriz de correlações amostral nos outros casos, que é o que ocorre mais frequentemente em situações práticas (EVERITT; HOTHORN, 2011). Já outros autores, como Khatree e Naik (2000) questionam essa escolha e argumentam que é preciso levar outras questões em conta.

O primeiro componente principal  $Y_1$  é a combinação linear

$$Y_1 = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p,$$

cuja variância amostral é a maior dentre todas as outras combinações lineares. É importante usar uma restrição nos valores desses coeficientes, geralmente  $\mathbf{a}_1^T \mathbf{a}_1 = 1$ , ou seja, a soma dos quadrados desses valores deve ser igual a 1. Isso deve ser feito porque a variância de  $Y_1$  poderia crescer de forma ilimitada apenas aumentando os coeficientes  $\mathbf{a}_1^T = [a_{11} \ a_{12} \ \dots \ a_{1p}]$ . A variância amostral de  $Y_1$  é dada por  $\mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$ , sendo  $\mathbf{S}$  a matriz de covariâncias amostral das  $X$  variáveis e  $\mathbf{a}_1$  é o autovetor da matriz  $\mathbf{S}$  associado ao maior autovetor  $\lambda$  dessa matriz (EVERITT; HOTHORN, 2011). A obtenção de autovalores  $\lambda$  e autovetores  $\mathbf{e}$  de uma matriz quadrada  $p \times p$  são tais que  $\mathbf{A}\mathbf{e} = \lambda\mathbf{e}$ . Para maiores detalhes, consultar, por exemplo, Ferreira (2011).

Ainda de acordo com Everitt e Hothorn (2011), o segundo componente principal,  $Y_2$  é definido como a combinação linear

$$Y_2 = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p,$$

ou seja,  $Y_2 = \mathbf{a}_2^T \mathbf{X}$ , em que  $\mathbf{a}_2^T = [a_{21} \ a_{22} \ \dots \ a_{2p}]$  e  $\mathbf{X}^T = [X_1 \ X_2 \ \dots \ X_p]$ , que possui a maior variância sujeito às condições

$$\mathbf{a}_2^T \mathbf{a}_2 = 1,$$

$$\mathbf{a}_2^T \mathbf{a}_1 = 0,$$

em que a segunda condição garante que  $Y_1$  e  $Y_2$  são não correlacionados. De forma similar, todos os outros componentes serão obtidos.

O vetor de coeficientes que define o  $i$ -ésimo componente principal,  $\mathbf{a}_i$  é o autovetor de  $\mathbf{S}$  associado com o seu  $i$ -ésimo maior autovetor. A variância do  $i$ -ésimo componente principal é dada por  $\lambda_i$ , sendo os  $\lambda_1, \lambda_2, \dots, \lambda_p$  os autovalores de  $\mathbf{S}$  sujeitos à restrição  $\mathbf{a}_i^T \mathbf{a}_i = 1$  (EVERITT; HOTHORN, 2011).

A proporção da variância total de  $\mathbf{X}$  explicada pelo  $i$ -ésimo componente principal é definida por

$$\frac{Var(Y_i)}{\text{Variância total de } \mathbf{X}} = \frac{\lambda_i}{\text{traço}(\mathbf{S})} = \frac{\lambda_i}{\sum_{j=1}^p \lambda_j}.$$

Além disso, as variâncias total e generalizada de  $\mathbf{X}$  podem ser descritas pelas variâncias total e

generalizada de  $\mathbf{Y}$ :

$$\text{traço}(\mathbf{S}) = \sum_{j=1}^p \lambda_j = S_1^2 + S_2^2 + \cdots + S_p^2 \quad \text{e} \quad |\mathbf{S}| = \prod_{j=1}^p \lambda_j.$$

Dessa forma, os vetores  $\mathbf{X}$  e  $\mathbf{Y}$  são equivalentes em relação a essas duas medidas de variação. Além disso, sempre o primeiro componente principal tem a maior proporção de explicação da variância total de  $\mathbf{X}$  (MINGOTI, 2005).

De acordo com Everitt e Hothorn (2011), os primeiros  $k$  componentes, em que  $k < p$ , explicam uma proporção da variância total,

$$\frac{\sum_{i=1}^k \text{Var}(Y_i)}{\text{Variância total de } X} = \frac{\sum_{i=1}^k \lambda_i}{\text{traço}(\mathbf{S})} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{j=1}^p \lambda_j}. \quad (3.3)$$

Os componentes principais podem ser obtidos a partir da matriz de covariâncias amostrais  $\mathbf{S}$  ou a partir da matriz de correlações amostrais  $\mathbf{R}$ . Extrair os componentes da matriz de covariâncias deve ser preferido quando as variáveis originais estão na mesma escala, o que é raro ocorrer. Extrair os componentes como os autovetores de  $\mathbf{R}$  é equivalente a calcular os componentes das variáveis originais para depois padronizar cada um para ter variância igual a 1 (EVERITT; HOTHORN, 2011).

Um passo importante da aplicação da técnica de ACP é a escolha de quantos componentes serão retidos. Um critério muito utilizado é avaliar a representatividade dos  $k$  primeiros componentes, de acordo com a equação (3.3). Define-se qual o valor de porcentagem da variação é pretendido (mínimo de 70%, por exemplo) e escolhem-se quantos componentes forem necessários para atingir essa representatividade. Porém, é necessário ter cautela com a escolha do número  $k$  pois a utilidade prática dos componentes principais diminui com o aumento desse valor (MINGOTI, 2005).

Um método gráfico que pode auxiliar na escolha do valor de  $k$  é o *scree plot*, em que é representado o valor  $k$  no eixo  $x$  e a porcentagem da variação explicada no eixo  $y$ . Assim, busca-se o ponto em que não há grande variação no eixo  $y$ , indicando que a inclusão de mais componentes não auxiliará muito na interpretação (EVERITT; HOTHORN, 2011). Mais detalhes sobre critérios podem ser vistos também em Khatree e Naik (2000).

Os valores numéricos dos componentes, denominados escores, podem ser calculados

para cada elemento amostral e, em seguida, esses valores podem ser analisados utilizando outras técnicas como análise de variância e análise de regressão (MINGOTI, 2005). Os escores dos primeiros dois componentes principais podem ser plotados em um diagrama de dispersão para identificar agrupamentos ou outros tipos de padrão existente nos dados.

Para calcular os escores dos componentes de cada observação  $i$ , se os componentes foram obtidos a partir da matriz de covariâncias amostral  $S$ , deve-se obter

$$Y_{i1} = \mathbf{a}_1^T \mathbf{X}_i, \quad Y_{i2} = \mathbf{a}_2^T \mathbf{X}_i, \quad \dots, \quad Y_{ik} = \mathbf{a}_k^T \mathbf{X}_i,$$

em que  $k$  é o número de componentes retidos e  $\mathbf{X}_i$  é o vetor de variáveis  $p \times 1$  para a observação  $i$ .

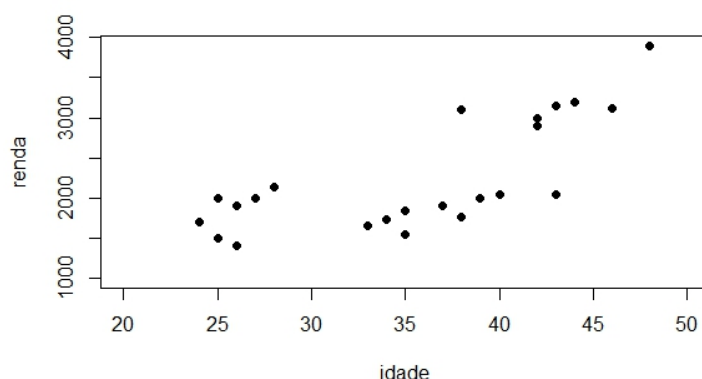
### 3.5 Análise de agrupamento (*Cluster Analysis*)

A técnica de análise de agrupamento (AA), também conhecida como análise de conglomerados, classificação ou *cluster analysis*, objetiva agrupar elementos da amostra de tal forma que (a) as observações de um mesmo grupo sejam similares entre si em relação às variáveis medidas e que (b) as observações de grupos diferentes sejam heterogêneas em relação a essas mesmas variáveis (MINGOTI, 2005).

“Análise de agrupamento” é uma expressão genérica que compreende vários métodos numéricos que pretendem descobrir grupos de observações homogêneas. O que esses métodos tentam é formalizar o que os seres humanos conseguem fazer bem em duas ou três dimensões, por exemplo, por meio de diagramas de dispersão. Os grupos são identificados pela avaliação das distâncias entre os pontos (EVERITT; HOTHORN, 2011).

Como forma de ilustração, considere um conjunto de dados fictícios em que há  $n = 23$  observações (pessoas) e registros sobre suas idades e rendas mensais, ou seja, há  $p = 2$  variáveis, ambas medidas na escala contínua. O objetivo é agrupar as pessoas em relação às duas variáveis. O diagrama de dispersão obtido se encontra na Figura 1 e é possível identificar três agrupamentos, apenas pela análise visual. Se houvesse mais uma variável e, conseqüentemente mais uma dimensão, ainda seria possível imaginar um gráfico com os pontos. Porém, com mais de três dimensões, torna-se mais difícil a visualização dos dados.

Figura 1 Diagrama de dispersão de dados fictícios de idade e renda de várias pessoas.



Fonte: modificado a partir de (BARTHOLOMEW et al., 2008, p.18)

Há dois objetivos possíveis de um agrupamento: agrupar as  $n$  observações em um número desconhecido de grupos ou classificar as observações em um conjunto predefinido de grupos. A análise de agrupamento deve ser utilizada no primeiro caso e análise discriminante no segundo. Na análise de agrupamento, geralmente, o número de grupos não é conhecido a princípio e encontrar o melhor agrupamento não é uma tarefa simples (FERREIRA, 2011).

Segundo Ferreira (2011), os métodos de agrupamento são divididos em hierárquicos (aglomerativos ou divisivos) e não hierárquicos. No primeiro tipo, há  $n$  grupos no início, cada um com uma observação, e no fim há um único grupo com todas as observações. A cada passo, cada observação ou grupo é unido a outra observação ou grupo (método hierárquico aglomerativo). No método hierárquico divisivo, há um único grupo com as  $n$  observações no início e, ao final, há  $n$  grupos. Nos métodos que não são hierárquicos é preciso definir o número  $k$  de grupos inicialmente para, em seguida, atribuir as  $n$  observações aos  $k$  grupos da melhor maneira possível. Sempre é preciso usar uma alocação arbitrária no início do processo e, iterativamente, buscar a alocação ótima.

Ao se utilizar o método hierárquico aglomerativo, depois que uma fusão é realizada, ela não será mais desfeita. Assim, quando o método coloca dois elementos em um mesmo grupo, eles não mais aparecerão em grupos diferentes. Para o pesquisador encontrar a melhor solução com o número de agrupamentos ótimo, ele deverá adotar algum critério de divisão (EVERITT; HOTHORN, 2011).

Agrupamentos obtidos a partir de métodos hierárquicos, aglomerativos ou divisivos, podem ser representados por um diagrama bidimensional muito utilizado, o dendrograma. Esse

gráfico ilustra as fusões ou divisões realizadas a cada passo do processo. Maiores detalhes podem ser obtidos em Everitt et al. (2011).

Ao se tentar realizar um agrupamento é muito importante saber quão próximas ou distantes estão as observações. Muitos métodos de agrupamento iniciam com uma matriz de distâncias  $n \times n$  que refletem uma medida de similaridade ou dissimilaridade entre os  $n$  elementos da amostra. Dois elementos são considerados próximos quando sua distância é pequena ou sua similaridade é grande Everitt et al. (2011).

A distância entre as observações  $i$  e  $j$  aparece na  $i$ -ésima linha e  $j$ -ésima coluna da matriz de distâncias. Por exemplo, se há  $n = 4$  elementos na amostra, a matriz de distâncias terá dimensão  $4 \times 4$  e será da forma

$$\begin{bmatrix} - & d_{12} & d_{13} & d_{14} \\ d_{21} & - & d_{23} & d_{24} \\ d_{31} & d_{32} & - & d_{34} \\ d_{41} & d_{42} & d_{43} & - \end{bmatrix},$$

em que  $d_{ij}$  é a distância entre os elementos  $i$  e  $j$ . Geralmente, essa matriz é simétrica, ou seja,  $d_{12} = d_{21}$ ,  $d_{13} = d_{31}$ , e assim por diante (BARTHOLOMEW et al., 2008).

Há muitos tipos de distâncias que podem ser calculadas entre pares de observações, mas um tipo muito simples e comum é a distância Euclidiana, calculada por

$$d_{ij} = \sqrt{\sum_{k=1}^p (X_{ik} - X_{jk})^2}$$

em que  $d_{ij}$  é a distância Euclidiana entre os elementos  $i$ , com os valores  $X_{i1}, X_{i2}, \dots, X_{ip}$ , e  $j$ , com os valores  $X_{j1}, X_{j2}, \dots, X_{jp}$ .

Uma técnica de agrupamento que pode ser utilizada é *k-means*, não hierárquica, que procura uma partição das  $n$  observações em  $k$  agrupamentos ( $G_1, G_2, \dots, G_k$ ), em que  $G_i$  denota o conjunto de observações que está no  $i$ -ésimo grupo e  $k$  é dado por algum critério numérico de minimização. A implementação mais usada do método tenta encontrar a partição dos  $n$  elementos em  $k$  grupos que minimize a soma de quadrados dentro dos grupos (*SQDG*)

em relação a todas as variáveis. Esse critério pode ser escrito como

$$SQDG = \sum_{j=1}^p \sum_{l=1}^k \sum_{i \in G_l} (X_{ij} - \bar{X}_j^{(l)})^2,$$

em que  $\bar{X}_j^{(l)} = \frac{1}{n_l} \sum_{i \in G_l} X_{ij}$  é a média dos indivíduos no grupo  $G_l$  em relação à variável  $j$  (EVERITT; HOTHORN, 2011).

Ainda de acordo com os autores, apesar de o problema parecer relativamente simples, ele não é tão direto. A tarefa de selecionar a partição com a menor soma de quadrados dentro dos grupos se torna complexa porque o número de partições possíveis se torna muito grande mesmo com um tamanho amostral não tão grande. Por exemplo, para  $n = 100$  e  $k = 5$ , o número de partições é da ordem de  $10^{68}$ . Esse problema levou ao desenvolvimento de algoritmos que não garantem encontrar a solução ótima, mas que levam a soluções satisfatórias.

### 3.5.1 Número de grupos

A etapa final do processo de agrupamento é definir a partição do conjunto de dados. Essa não é uma tarefa simples e existem vários métodos propostos para definir o número  $k$  de agrupamentos ou em qual passo o algoritmo de agrupamento deve ser interrompido. Apesar de não haver um consenso, alguns critérios podem ser utilizados para auxiliar na decisão final (MINGOTI, 2005).

Quando métodos hierárquicos de agrupamento são utilizados, um dendrograma é obtido e deve-se decidir em qual altura o corte deve ser realizado, o que vai gerar um determinado número de grupos. A questão é decidir o ponto de corte. Uma maneira informal de tomar essa decisão é avaliar as mudanças na altura do dendrograma nos diferentes passos e escolher a maior mudança observada. Porém, mesmo com um número de observações não muito grande (como 15 ou 20), não é simples decidir onde está essa maior mudança (EVERITT; HOTHORN, 2011).

Outros critérios informais utilizam um gráfico da medida *versus* o número de grupos em que se busca identificar grandes mudanças no gráfico para um determinado número  $k$  e um ponto de parada é sugerido. Porém, esses critérios são subjetivos. Várias técnicas mais formais têm sido propostas e alguns trabalhos avaliaram suas propriedades, tais como Milligan e Cooper

(1985) e Dimitriadou, Dolničar e Weingessel (2002).

O trabalho de Milligan e Cooper (1985) identificou como as duas melhores técnicas as propostas por Caliński e Harabasz (1974), também conhecida como pseudo  $F$ , e Duda e Hart (1973), denominada pseudo  $T^2$ . Everitt et al. (2011) apresentam um resumo desses e de outros critérios que podem auxiliar nessa decisão.

Um método que utiliza a matriz de dissimilaridade é o *silhouette plot* proposto por Kaufman e Rousseeuw (2009). Neste método, para cada observação  $i$ , é definido um índice  $s(i)$  entre  $-1$  e  $1$ . Quando este valor é próximo de  $1$  indica que a observação foi bem classificada no grupo, se próximo de  $-1$  indica o contrário. Quando o valor assumido é próximo de  $0$  não está claro se a observação deveria estar no seu grupo ou em outro. O gráfico mostra os valores de  $s(i)$  em forma de barras horizontais, ordenadas de forma decrescente para cada agrupamento. Comparar *silhouette plots* para soluções obtidas com diferentes números de grupos pode auxiliar na escolha deste número, levando a melhores agrupamentos.

A estatística GAP foi criada com o mesmo propósito por Tibshirani, Walther e Hastie (2001). O método formaliza a ideia de procurar pelo “*elbow*” no gráfico do número de grupos *versus* algum critério de otimização, ou seja, procurar por uma grande mudança na inclinação do gráfico a partir da qual não há grandes ganhos no critério.

Diante dessa diversidade, é crucial não utilizar apenas um método para definir o número de grupos, mas avaliar os resultados obtidos com diferentes critérios. Além disso, alguns deles fazem suposições sobre a estrutura dos grupos e terão bom desempenho apenas se as suposições forem atendidas (EVERITT; HOTHORN, 2011).

Além disso, quando estudos socioeconômicos estão em questão, adicionalmente aos métodos estatísticos, as especificidades do problema analisado devem ser levadas em conta para que se decida qual critério fornece grupos cuja interpretação seja mais útil (CARVALHO; MATA; RESENDE, 2007).



## 4 Metodologia

### 4.1 Bases de dados e variáveis do estudo

As bases de dados utilizadas neste trabalho são provenientes do Ministério da Previdência Social (disponível em [www.previdencia.gov.br](http://www.previdencia.gov.br)) e do Atlas do Desenvolvimento Humano no Brasil 2013 (disponível em [www.atlasbrasil.org.br](http://www.atlasbrasil.org.br)), que utiliza os censos demográficos realizados pelo IBGE em 1991, 2000 e 2010 para calcular cerca de 230 variáveis para os 5.565 municípios brasileiros. Os dados estão tabulados em formato de planilhas .xls, o que facilita seu tratamento. A partir das variáveis demográficas presentes no Atlas, inicialmente três foram escolhidas - outras poderão ser acrescentadas de acordo com o que for sugerido pela literatura a ser analisada. A Tabela 1 apresenta a lista provisória das variáveis .

Tabela 1 Lista provisória das variáveis

Sigla	Descrição	Fonte
QUANT	quantidade de benefícios em dezembro	MPS
ARREC	valor arrecadado	MPS
VAB	valor anual dos benefícios	MPS
VBD	valor dos benefícios em dezembro	MPS
POP	população residente total no município	Atlas
T_ENV	razão entre a população de 65 anos ou mais de idade e a população total multiplicado por 100	Atlas
P_FORMAL	razão entre o número de pessoas de 18 anos ou mais formalmente ocupadas e o número total de pessoas ocupadas nessa faixa etária multiplicado por 100	Atlas

A escolha das variáveis do Atlas se deu devido à suposição que os valores do benefícios dos municípios devem ter relação com o número de habitantes, a idade de seus habitantes e a porcentagem de trabalho formalizado. Outras variáveis poderiam ser escolhidas para este fim, porém, posteriormente, mais algumas podem ser incluídas no estudo.

Como as variáveis de previdência estão expressas em diferentes unidades/medidas, torna-se necessária uma padronização. Será utilizado o valor por habitante, ou seja, as variáveis serão divididas pelo número de habitantes do município (POP). Outro tipo de transformação possível seria fazer com que todas variassem no intervalo de 0 a 1 ou de 0 a 100.

O primeiro passo da análise será a obtenção da matriz de correlações entre as variáveis, duas a duas, de forma a identificar os pares de variáveis mais associadas entre si. Tal matriz terá a forma apresentada na equação (3.2). Essa etapa compreende uma análise exploratória dos dados que auxiliará na aplicação das técnicas multivariadas posteriores.

A ACP será utilizada para reduzir a dimensionalidade dos dados, conforme explicitado na seção 3.4. Pretende-se reduzir o conjunto das variáveis originais correlacionadas entre si a um novo conjunto de variáveis, os componentes principais, não correlacionadas. O que se espera é que um pequeno número dessas novas variáveis expliquem boa parte da variação presente nos dados originais. Se apenas dois componentes,  $Y_1$  e  $Y_2$ , já explicarem boa parte da variação presente nos dados, será possível obter um gráfico bidimensional com os valores das observações, os escores, de  $Y_1$  e  $Y_2$ .

Além disso, é interessante que os componentes principais tenham uma interpretação prática. Para isso, após a definição de quantos serão utilizados, as correlações entre cada variável original e cada componente serão calculadas, os chamados *loadings*. Os valores dessas correlações, bem como seus sinais, indicarão como cada componente poderá ser interpretado.

Após a aplicação da técnica de componentes principais será utilizada a análise de agrupamento para identificar os grupos de municípios com características semelhantes. Como foi mostrado na seção 3.5, a AA oferece várias opções para a escolha da medida de distância, do método de agrupamento e do número de grupos.

Como medida de similaridade será utilizada a distância euclidiana, a mais aplicada em análise de agrupamento. Um dos métodos de agrupamento será o método hierárquico aglomerativo de Ward, o mais indicado quando as variáveis medidas estão na escala contínua.

Além do método aglomerativo também será aplicado o não hierárquico das  $k$ -médias, também muito popular. Este método requer que o número de grupos seja definido antes de sua aplicação. Assim, os métodos para definição do número de grupos mostrados na seção 3.5.1 serão utilizados e verificar-se-á se, e em quais condições, há convergência entre eles no caso dos dados analisados. Ao final, os resultados dos agrupamentos obtidos com os métodos Ward e  $k$ -médias serão comparados entre si.

Todas as rotinas necessárias para a análise dos dados serão realizadas utilizando o programa *R* em sua versão 3.2.0 (??). (R CORE TEAM, 2014).

## 5 Resultados Esperados

O que se espera é que os municípios brasileiros sejam agrupados de forma satisfatória de acordo com as variáveis demográficas e previdenciárias escolhidas. Pretende-se identificar padrões nos dados e que semelhanças e diferenças entre os municípios sejam detectadas. O resultado pode ajudar na análise dos impactos e dos problemas da previdência social, contribuindo para novas pesquisas e para a melhoria dessa política pública.

## 6 Cronograma

O projeto será concluído em 24 meses, no período de março de 2015 a fevereiro de 2017. As atividades mensais a serem desenvolvidas compõem-se das seguintes etapas:

Atividades	Meses																							
	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
1	X	X	X	X																				
2					X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
3										X	X	X	X	X	X	X								
4													X	X	X	X	X	X						
5																		X						
6															X	X	X	X	X	X	X	X	X	X
7																								X

- 1: Definição do Tema;
- 2: Revisão de literatura;
- 3: Implementação do trabalho;
- 4: Obtenção dos resultados;
- 5: Qualificação;
- 6: Redação da dissertação;
- 7: Finalização do trabalho/defesa da dissertação.

## 7 Disciplinas necessárias

As disciplinas necessárias para realizar este projeto são: Álgebra Linear Aplicada; Probabilidade; Inferência Estatística; Inglês Instrumental em Estatística Aplicada e Biometria; Estatística Computacional; Análise Multivariada.

## 8 Referências Bibliográficas

- BARTHOLOMEW, D. J.; STEELE, F.; GALBRAITH, J.; MOUSTAKI, I. **Analysis of multivariate social science data**. Boca Raton: CRC press, 2008.
- CALIŃSKI, T.; HARABASZ, J. A dendrite method for cluster analysis. **Communications in Statistics - theory and methods**, Taylor and Francis, v. 3, n. 1, p. 1–27, 1974.
- CAMARANO, A. A. O. Novo regime demográfico: uma nova relação entre população e desenvolvimento? Instituto de Pesquisa Econômica Aplicada (Ipea), 2014.
- CARVALHO, A.; MATA, D. D.; RESENDE, G. M. Clusterização dos municípios brasileiros. **Dinâmica dos Municípios, Brasília: IPEA**, 2007.
- DIMITRIADOU, E.; DOLNIČAR, S.; WEINGESSEL, A. An examination of indexes for determining the number of clusters in binary data sets. **Psychometrika**, Springer, v. 67, n. 1, p. 137–159, 2002.
- DUDA, R. O.; HART, P. E. **Pattern classification and scene analysis**. [S.l.: s.n.], 1973.
- EVERITT, B.; HOTHORN, T. **An introduction to applied multivariate analysis with R**. NY: Springer Science & Business Media, 2011.
- EVERITT, B. S.; LANDAU, S.; LEESE, M.; STAHL, D. **Cluster analysis**. 5. ed. UK: John Wiley and Sons, 2011.
- FERREIRA, D. F. **Estatística multivariada**. 2. ed. Lavras: Editora UFLA, 2011.
- FRANÇA, Á. S. de. **Previdência social e a economia dos municípios**. 5. ed. Brasília: Anfip, 2004.
- KAUFMAN, L.; ROUSSEEUW, P. J. **Finding groups in data: an introduction to cluster analysis**. [S.l.]: John Wiley & Sons, 2009.
- KHATREE, R.; NAIK, D. N. **Multivariate data reduction and discrimination with SAS software**. NC: SAS Institute Inc., 2000.

MILLIGAN, G. W.; COOPER, M. C. An examination of procedures for determining the number of clusters in a data set. **Psychometrika**, Springer, v. 50, n. 2, p. 159–179, 1985.

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

R CORE TEAM. **R: A Language and Environment for Statistical Computing**. Vienna, Austria, 2014. Disponível em: <<http://www.R-project.org/>>.

REIS, P. R. da C.; SILVEIRA, S. d. F. R.; BRAGA, M. J. Previdência social e desenvolvimento socioeconômico: impactos nos municípios de pequeno porte de minas gerais. **RAP: Revista Brasileira de Administração Pública**, SciELO Brasil, v. 47, n. 3, 2013.

TIBSHIRANI, R.; WALTHER, G.; HASTIE, T. Estimating the number of clusters in a data set via the gap statistic. **Journal of the Royal Statistical Society: Series B (Statistical Methodology)**, Wiley Online Library, v. 63, n. 2, p. 411–423, 2001.

VASCONCELOS, A. M. N.; GOMES, M. M. F. Transição demográfica: a experiência brasileira. **Epidemiologia e Serviços de Saúde**, Coordenação-Geral de Desenvolvimento da Epidemiologia em Serviços/Secretaria de Vigilância em Saúde/Ministério da Saúde, v. 21, n. 4, p. 539–548, 2012.