

UNIVERSIDADE FEDERAL DE MINAS GERAIS
DEPARTAMENTO DE ESTATÍSTICA
PÓS-GRADUAÇÃO LATO SENSU –
ESPECIALIZAÇÃO EM ESTATÍSTICA COMPUTACIONAL APLICADA

Pedro Mateus Moraes de Almeida

Trabalho final da disciplina Estatística Multivariada Computacional
Aplicação das técnicas de PCA e EFA na base de dados SES do Rio de Janeiro

Belo Horizonte

2023

Sumário

1. Elabore uma breve descritiva (tamanho populacional, vetor de médias e matriz de correlações dos indicadores em cada dimensão) das variáveis disponíveis para a cidade escolhida. Discorra sucintamente sobre os resultados.	2
Descrição da Base de Dados SES_Rio_de_Janeiro.csv	2
Distribuição das Dimensões:	4
Saneamento	4
Moradia	5
Emprego	6
Educação	7
Médias dos Indicadores.....	8
Desvio Padrão dos Indicadores	9
Correlação	10
2. Construa as componentes principais para cada dimensão. Faça uma breve análise qualitativa dos pesos obtidos em termos das variáveis originais.	13
Dimensão Saneamento	16
Dimensão Moradia	18
Dimensão Emprego	21
Dimensão Educação	23
Todas as dimensões.....	25
3. Estime (pelo menos um) modelo EFA para cada dimensão. Faça uma breve análise qualitativa das cargas fatoriais estimadas em termos das variáveis originais.	29
Dimensão Saneamento	32
Dimensão Moradia	35
Dimensão Emprego	37
Dimensão Educação	39
Todas as dimensões.....	41
4. Exemplifique como as diferentes regras de seleção podem afetar a escolha do número de componentes e/ou fatores em cada dimensão. Comente sobre como o número de componentes /fatores escolhido é afetado pela correlação entre as variáveis originais. Opcional: faça também uma conexão com a regra de aceitação/rejeição de H_0 no teste da qualidade de ajuste do(s) modelo(s) EFA ajustado.	44
5. Construa um “índice de status socioeconômico” utilizando a PCA e/ou a EFA (utilizando pelo menos um indicador de cada dimensão). Interprete os resultados em termos da dimensão correspondente e do padrão de “status socioeconômico” para a cidade escolhida.....	47

1. Elabore uma breve descritiva (tamanho populacional, vetor de médias e matriz de correlações dos indicadores em cada dimensão) das variáveis disponíveis para a cidade escolhida. Discorra sucintamente sobre os resultados.

Descrição da Base de Dados SES_Rio_de_Janeiro.csv

Cidade: Rio de Janeiro

População: 2.177.297

Granularidade: cada linha representa uma região da cidade Rio de Janeiro

Conjunto de Dados:

- Número de Linhas: 200
- Número de Colunas: 18

2. Dados Faltantes:

- **Valores Nulos/Em Branco:** 0

3. Dimensões dos dados:

- Saneamento
- Moradia
- Emprego
- Educação.

4. Nomes das Colunas:

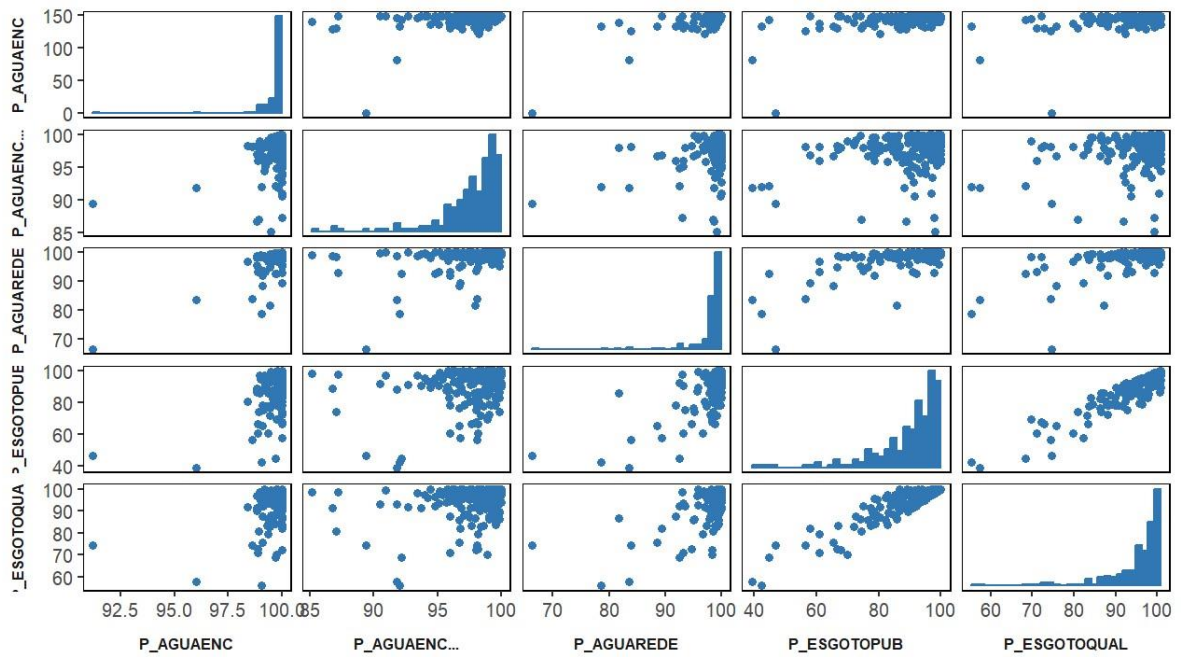
- UF: Unidade Federativa (estado).
- municipality: Identificação do município.
- code: Código do município.

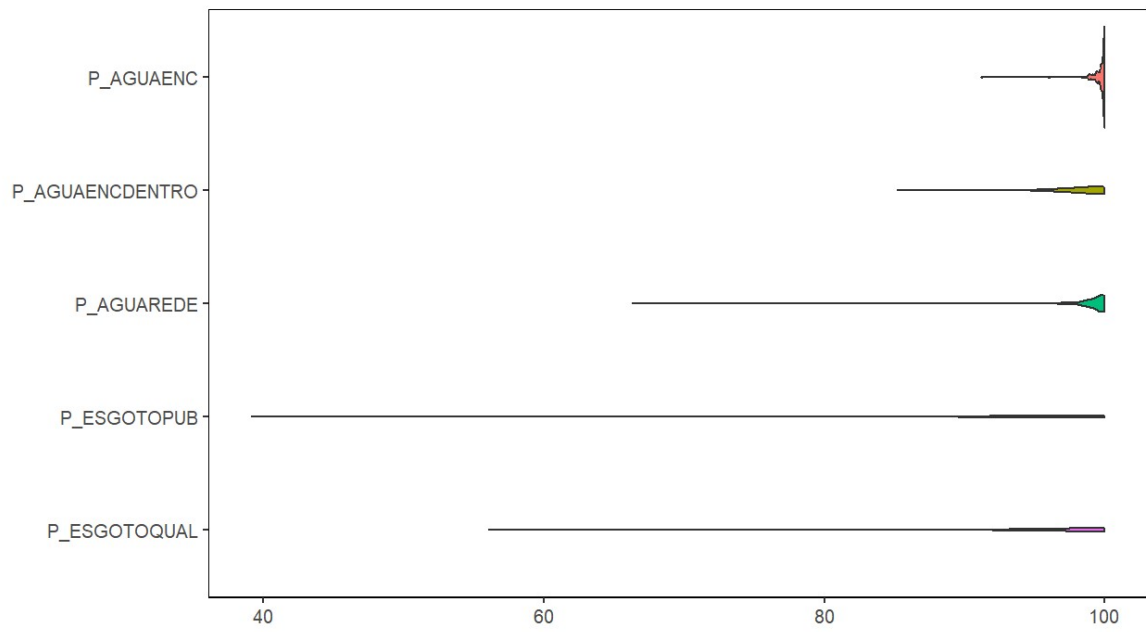
- area_de_ponderacao: unidade geográfica
- N: População estimada da área.
- **Dimensão Saneamento**
 1. **P_AGUAENC**: Proporção de domicílios com acesso à água encanada.
 2. **P_AGUAENC DENTRO**: Proporção de domicílios com acesso à água encanada dentro do domicílio.
 3. **P_AGUAREDE**: Proporção de domicílios com acesso à água de uma rede pública.
 4. **P_ESGOTOPUB**: Proporção de domicílios conectados à rede pública de esgoto.
 5. **P_ESGOTOQUAL**: Proporção de domicílios conectados a algum tipo de sistema de esgoto.
- **Dimensão Moradia**
 1. **P_MATPAREDES**: Proporção de domicílios com paredes externas feitas em sua maioria de materiais duráveis.
 2. **P_OVERCROWDING**: Proporção de domicílios com mais de 3 moradores por quarto/dormitório (superlotação).
- **Dimensão Emprego**
 1. **P_DESEMP**: Taxa de desemprego entre a população com 15 anos de idade ou mais.
 2. **P_FORTRAB**: Participação da força de trabalho entre a população com 15 anos de idade ou mais.
- **Dimensão Educação**
 1. **P_FREQESCOLA**: Proporção da população entre 15 e 17 anos frequentando a escola.
 2. **P_ENSFUND**: Proporção da população com 25 anos ou mais com o Ensino Fundamental completo.

3. **P_ENSMED**: Proporção da população com 25 anos ou mais com o Ensino Médio completo.
4. **P_ENSSUP**: Proporção da população com 25 anos ou mais com o Ensino Superior completo.

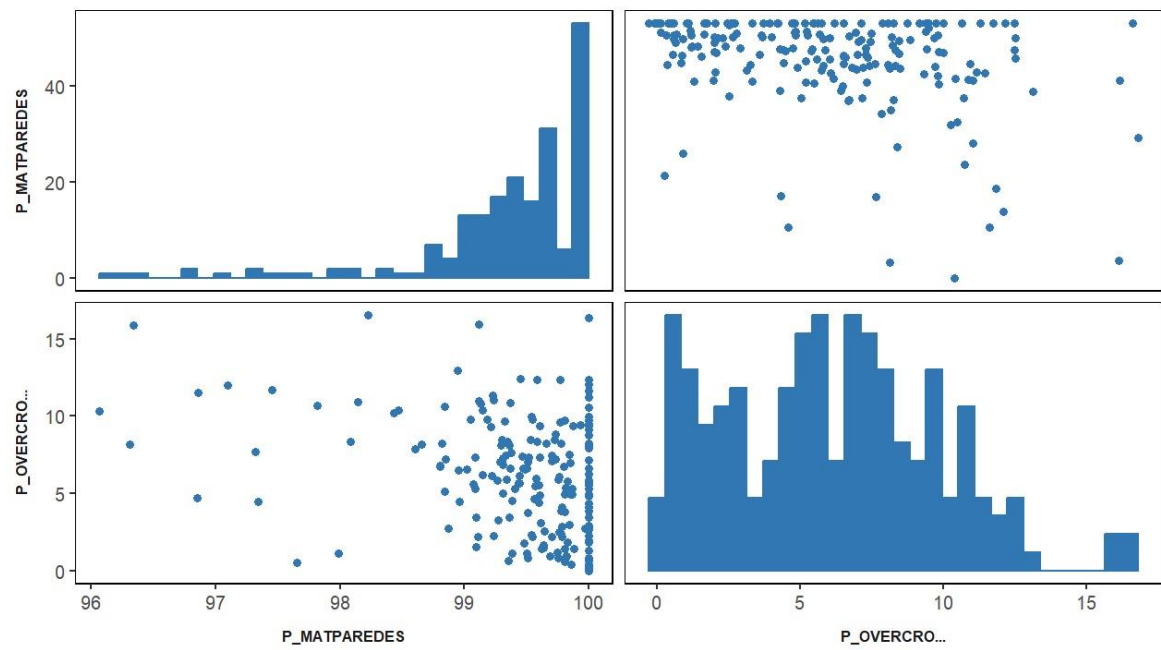
Distribuição das Dimensões:

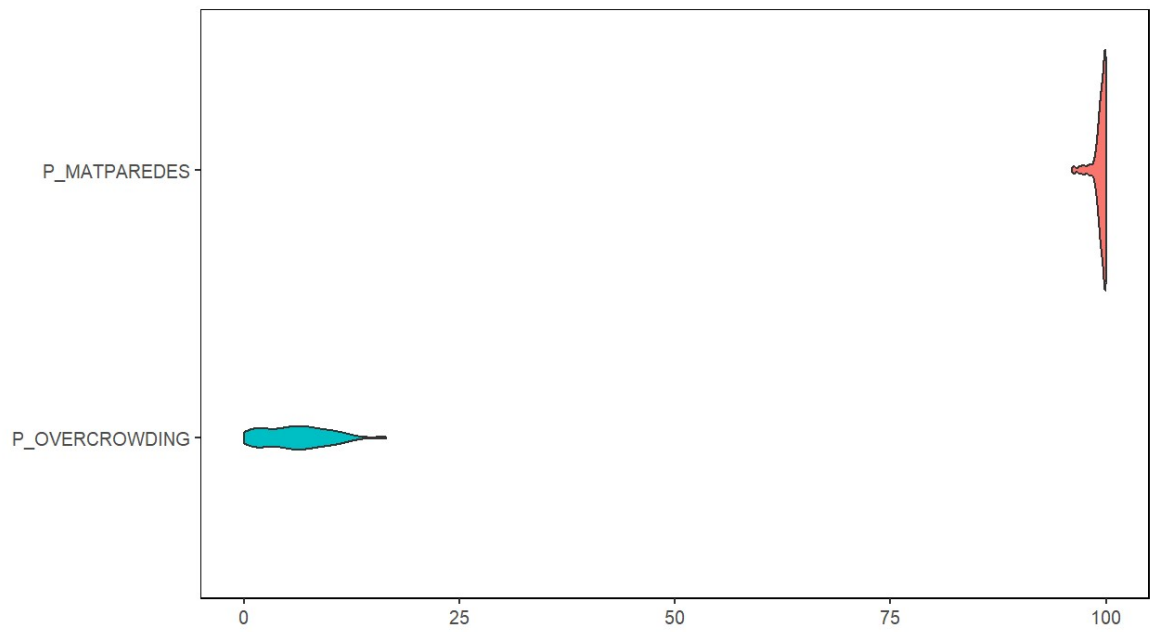
Saneamento



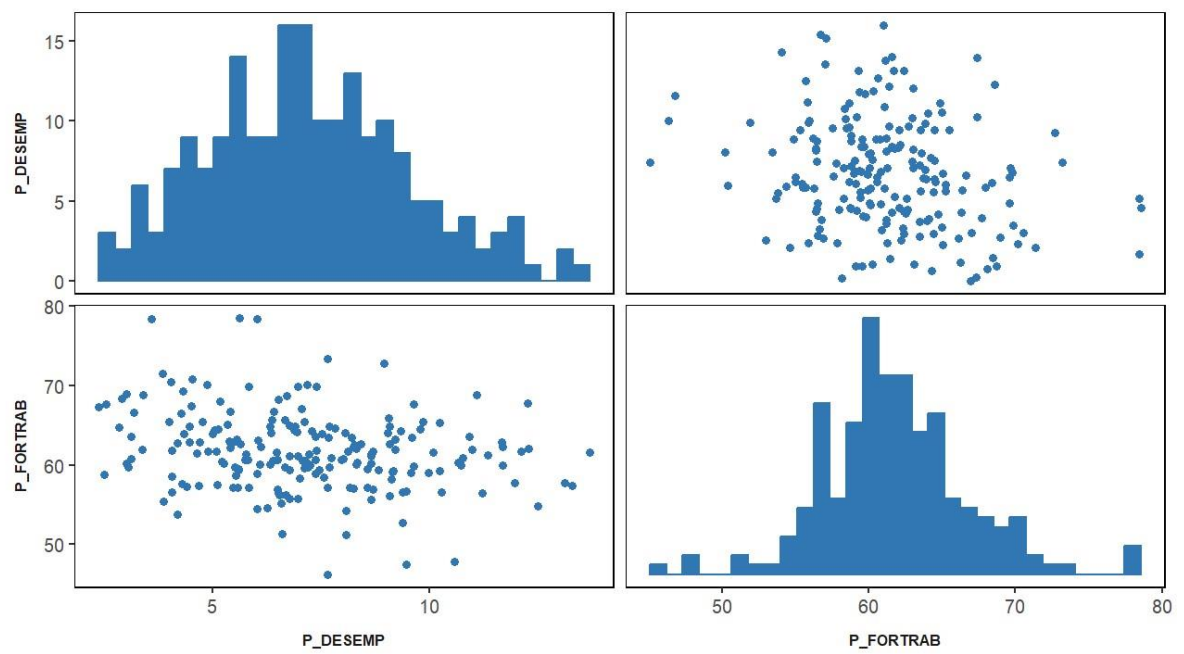


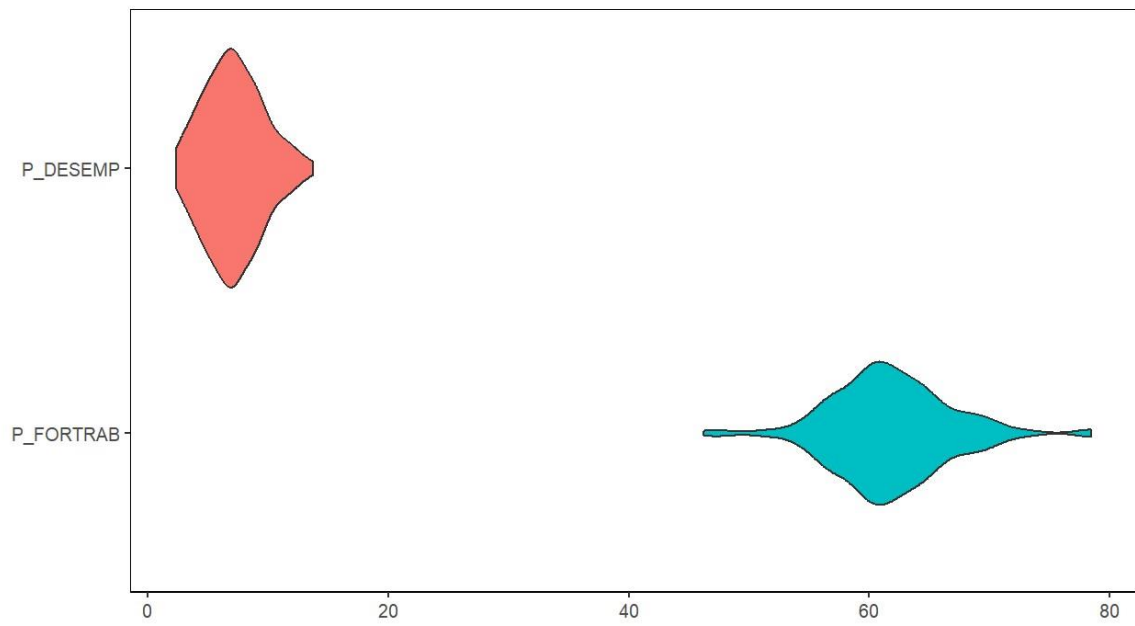
Moradia



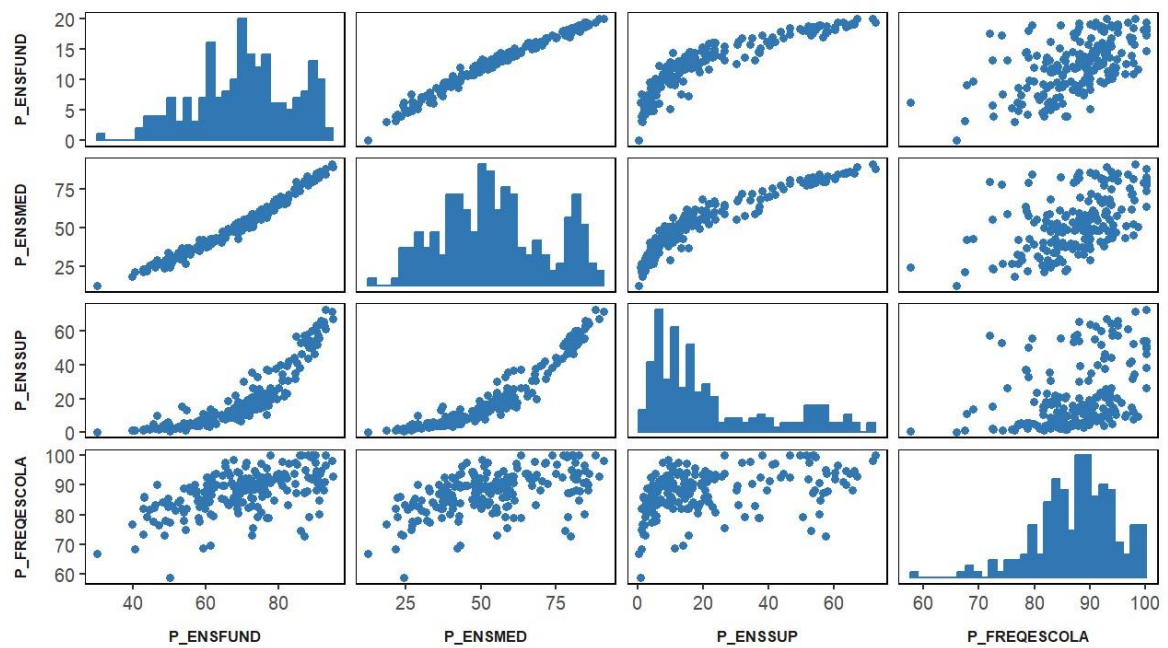


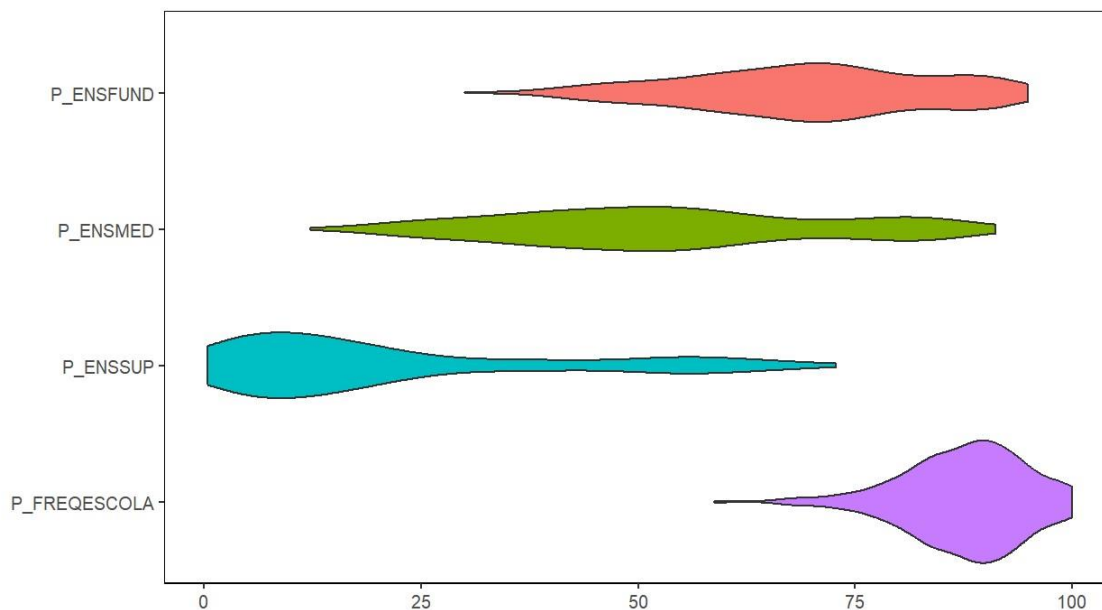
Emprego





Educação





Médias dos Indicadores

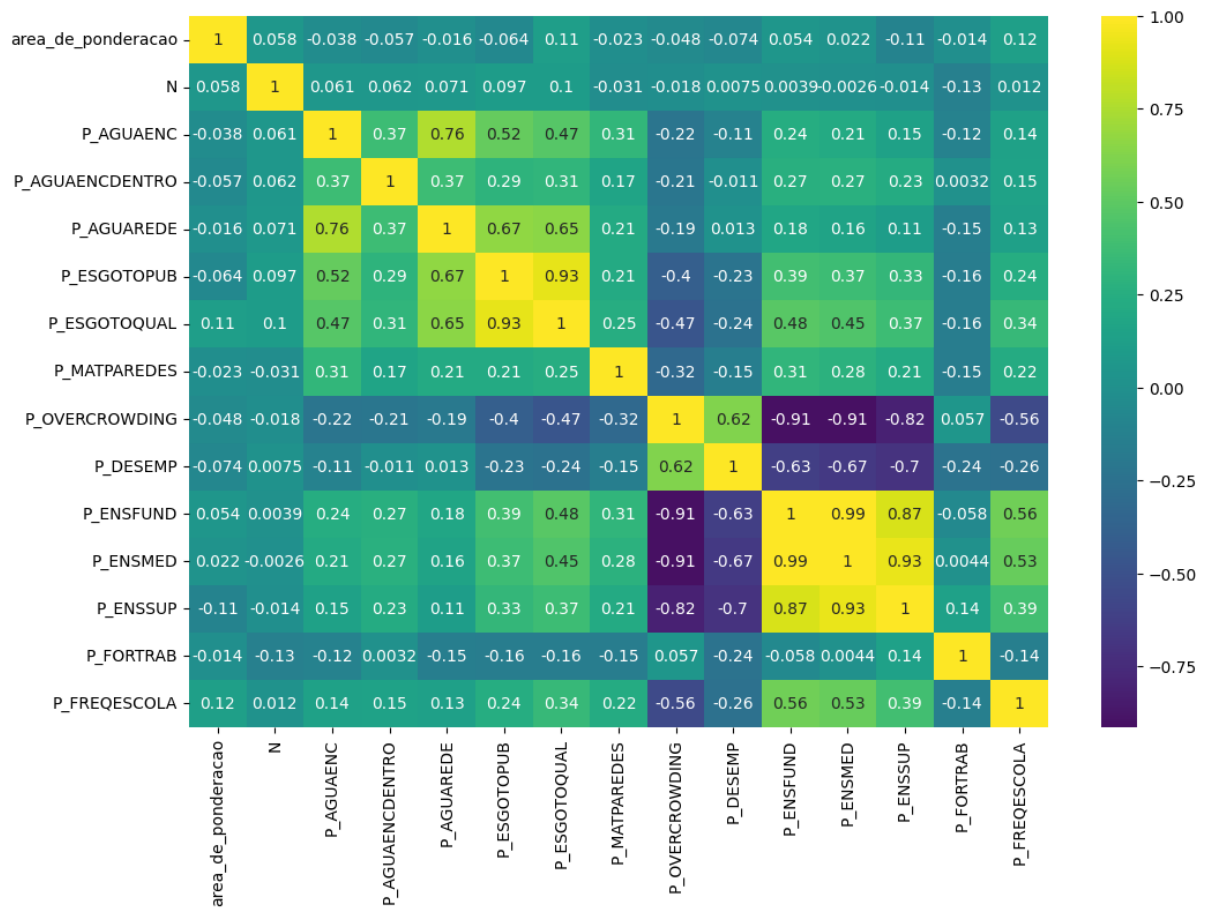
- População estimada (**N**): 10.887 habitantes (em média por área de ponderação)
- Acesso à água encanada (**P_AGUAENC**): 99,73%
- Acesso à água encanada dentro do domicílio (**P_AGUAENCDENTRO**): 97,81%
- Acesso à água de uma rede pública (**P_AGUAREDE**): 98,36%
- Conexão à rede pública de esgoto (**P_ESGOTOPUB**): 90,64%
- Sistema de esgoto qualificado (**P_ESGOTOQUAL**): 94,72%
- Paredes externas feitas de materiais duráveis (**P_MATPAREDES**): 99,44%
- Superlotação (**P_OVERCROWDING**): 5,99%
- Taxa de desemprego (**P_DESEMP**): 7,14%
- Ensino Fundamental completo (**P_ENSFUND**): 69,82%
- Ensino Médio completo (**P_ENSMED**): 53,36%
- Ensino Superior completo (**P_ENSSUP**): 21,22%
- Participação da força de trabalho (**P_FORTRAB**): 61,77%

- Frequência escolar para população entre 15 e 17 anos (**P_FREQESCOLA**): 87,66%

Desvio Padrão dos Indicadores

- **População estimada (N)**: O desvio padrão é de 2774.20 habitantes.
- **Acesso à água encanada (P_AGUAENC)**: O desvio padrão é de 0.74%.
- **Acesso à água encanada dentro do domicílio (P_AGUAENCIENTRO)**: O desvio padrão é de 2.58%.
- **Acesso à água de uma rede pública (P_AGUAENREDE)**: O desvio padrão é de 3.78%.
- **Conexão à rede pública de esgoto (P_ESGOTOPUB)**: O desvio padrão é de 11.50%.
- **Sistema de esgoto qualificado (P_ESGOTOQUAL)**: O desvio padrão é de 7.52%.
- **Paredes externas feitas de materiais duráveis (P_MATPAREDES)**: O desvio padrão é de 0.73%.
- **Superlotação (P_OVERCROWDING)**: O desvio padrão é de 3.70%.
- **Taxa de desemprego (P_DESEMP)**: O desvio padrão é de 2.42%.
- **Ensino Fundamental completo (P_ENSFUND)**: O desvio padrão é de 13.87%.
- **Ensino Médio completo (P_ENSMED)**: O desvio padrão é de 18.29%.
- **Ensino Superior completo (P_ENSSUP)**: O desvio padrão é de 19.33%.
- **Participação da força de trabalho (P_FORTRAB)**: O desvio padrão é de 4.96%.
- **Frequência escolar para população entre 15 e 17 anos (P_FREQESCOLA)**: O desvio padrão é de 7.09%.

Correlação



Pontos Significativos da Matriz de Correlação (Foram considerados valores superiores a 0,45 ou inferiores a -0,45)

1. Saneamento (as correlações abaixo fazem sentido dado a relação entre as variáveis)

- **P_AGUAENC e P_AGUAEREDE: 0.7578**
- **P_AGUAENC e P_ESGOTOPUB: 0.5237**
- **P_AGUAENC e P_ESGOTOQUAL: 0.4718**
- **P_AGUAEREDE e P_ESGOTOPUB: 0.6705**
- **P_AGUAEREDE e P_ESGOTOQUAL: 0.6509**
- **P_ESGOTOPUB e P_ESGOTOQUAL: 0.9252**

2. Moradia e Educação

- **P_OVERCROWDING e P_ENSFUND: -0.9079**
- **P_OVERCROWDING e P_ENSMED: -0.9149**
- **P_OVERCROWDING e P_ENSSUP: -0.8156**
- **P_OVERCROWDING e P_FREQESCOLA: -0.5554**

3. Educação

- **P_ENSFUND e P_ENSMED: 0.9856 (justificável dado para fazer o médio precisa concluir o fundamenta)**
- **P_ENSFUND e P_ENSSUP: 0.8729 (justificável dado para fazer o médio precisa concluir o fundamenta)**
- **P_ENSMED e P_ENSSUP: 0.9272 (justificável dado para fazer o médio precisa concluir o fundamenta)**

4. Moradia

- **P_MATPAREDES e P_ESGOTOQUAL: 0.2497 (Baixa correlação)**

5. Outras Correlações Significativas

- **P_DESEMP e P_OVERCROWDING: 0.6223**
 - **P_DESEMP e P_ENSFUND: -0.6279**
 - **P_DESEMP e P_ENSMED: -0.6661**
 - **P_DESEMP e P_ENSSUP: -0.7023**
 - **P_FREQESCOLA e P_ENSFUND: 0.5622**
 - **P_FREQESCOLA e P_ENSMED: 0.5342**
 - **P_FREQESCOLA e P_ENSSUP: 0.3876 (Baixa correlação)**
-
- Alta correlação entre acesso à água encanada (**P_AGUAENC**) e acesso à água de uma rede pública (**P_AGUAREDE**): 0,76.

- Correlação significativa entre a taxa de desemprego (**P_DESEMP**) e a superlotação (**P_OVERCROWDING**): 0,62.
- Correlações negativas fortes entre superlotação (**P_OVERCROWDING**) e níveis de educação (**P_ENSFUND**, **P_ENSMED**, **P_ENSSUP**), variando de -0,81 a -0,91.
- Correlação negativa entre taxa de desemprego (**P_DESEMP**) e níveis de educação (**P_ENSFUND**, **P_ENSMED**, **P_ENSSUP**), variando de -0,63 a -0,70.

Essas correlações destacam relações importantes entre educação, emprego, condições de moradia e acesso a serviços básicos, refletindo aspectos críticos do bem-estar socioeconômico no Rio de Janeiro. Por exemplo, a forte correlação negativa entre superlotação e educação sugere que áreas com maior superlotação tendem a ter menores níveis de educação formal, o que pode ter implicações significativas para políticas de planejamento urbano e social.

2. Construa as componentes principais para cada dimensão. Faça uma breve análise qualitativa dos pesos obtidos em termos das variáveis originais.

Introdução e Fundamentos

A Análise de Componentes Principais (PCA) é uma técnica estatística crucial na Estatística Multivariada, que lida com a análise conjunta de múltiplas variáveis dependentes. Ela visa representar a variabilidade de um conjunto de dados através de novos vetores ortogonais chamados componentes principais. Estes vetores são construídos para maximizar a variabilidade capturada, com a primeira componente contendo a maior parte da variabilidade e as subsequentes decrescendo em termos de importância.

Construção das Componentes Principais

Cada componente principal no PCA é uma combinação linear das variáveis originais, com pesos específicos definidos para maximizar a variabilidade sob a restrição de que a norma ℓ_2 do vetor de pesos seja igual a 1. Este processo garante que a variância total não seja alterada pela transformação, mantendo a informação original dos dados.

A ortogonalidade entre os componentes, uma característica fundamental do PCA, assegura que eles sejam não correlacionados, simplificando a interpretação e análise dos dados.

Maximização da Variabilidade:

- O PCA busca direções (componentes principais) no espaço multidimensional que maximizam a variância dos dados projetados. Isso significa identificar os eixos nos quais os dados se espalham mais.
- A decomposição espectral da matriz de covariâncias dos dados é fundamental nesse processo. Os autovalores da matriz representam a variância capturada por cada componente principal, e os autovetores correspondem aos componentes principais.

Importância Relativa dos Componentes:

- A ordenação dos componentes principais por variância permite uma análise hierárquica da informação contida nos dados. O primeiro componente principal contém a maior variância, com cada componente subsequente capturando progressivamente menos variância.
- A análise da variância acumulada por um número selecionado de componentes principais ajuda a determinar quantos deles devem ser retidos para uma representação eficaz dos dados

Aplicações Práticas

1. Redução de Dimensionalidade:

- O PCA é frequentemente usado para reduzir o número de variáveis em um conjunto de dados, escolhendo um subconjunto de componentes principais que capturam a maior parte da variabilidade. Muito utilizado em modelos de Machine Learning, torna o modelo mais eficiente e menos sujeito ao overfitting, contribuindo para análises mais precisas e confiáveis. Esta técnica é particularmente útil em grandes conjuntos de dados, onde o número elevado de variáveis torna a análise complexa e difusa.

2. Feature Selection:

- A PCA é frequentemente usada como uma técnica de seleção de características (feature selection) em machine learning. Ao transformar um grande número de variáveis possivelmente correlacionadas em um conjunto menor de variáveis não correlacionadas (componentes principais), a PCA facilita a identificação e seleção das características mais significativas para um modelo.
- A análise qualitativa no PCA envolve a identificação das variáveis originais que são mais representativas em cada componente principal. Isso é feito observando **os pesos** atribuídos a cada variável na composição do componente.

3. Criação de Índices:

- O PCA é empregado para sintetizar variáveis complexas em componentes principais. Esta abordagem permite uma compreensão mais clara das dimensões essenciais que caracterizam o ambiente socioeconômico, simplificando a análise de dados que, de outra forma, seriam muito complexos.

4. Análise de Correlação Multivariada:

- O PCA pode ser utilizada para analisar a estrutura de correlação entre variáveis. Componentes principais não correlacionados indicam que a PCA transformou variáveis potencialmente correlacionadas em um novo conjunto de variáveis independentes. Esta técnica é particularmente útil em cenários onde a multicolinearidade (alta correlação entre variáveis explicativas) pode ser um problema, ajudando a melhorar a performance e a interpretabilidade dos modelos preditivos.

5. Análise Exploratória de Dados:

- A PCA ajuda na visualização e interpretação de dados multivariados, possibilitando uma compreensão mais abrangente do comportamento e interdependência das variáveis.
- Ferramentas visuais como biplots podem ser usadas para ilustrar as projeções ortogonais das variáveis originais nos componentes principais.

6. Visualização de Dados Complexos:

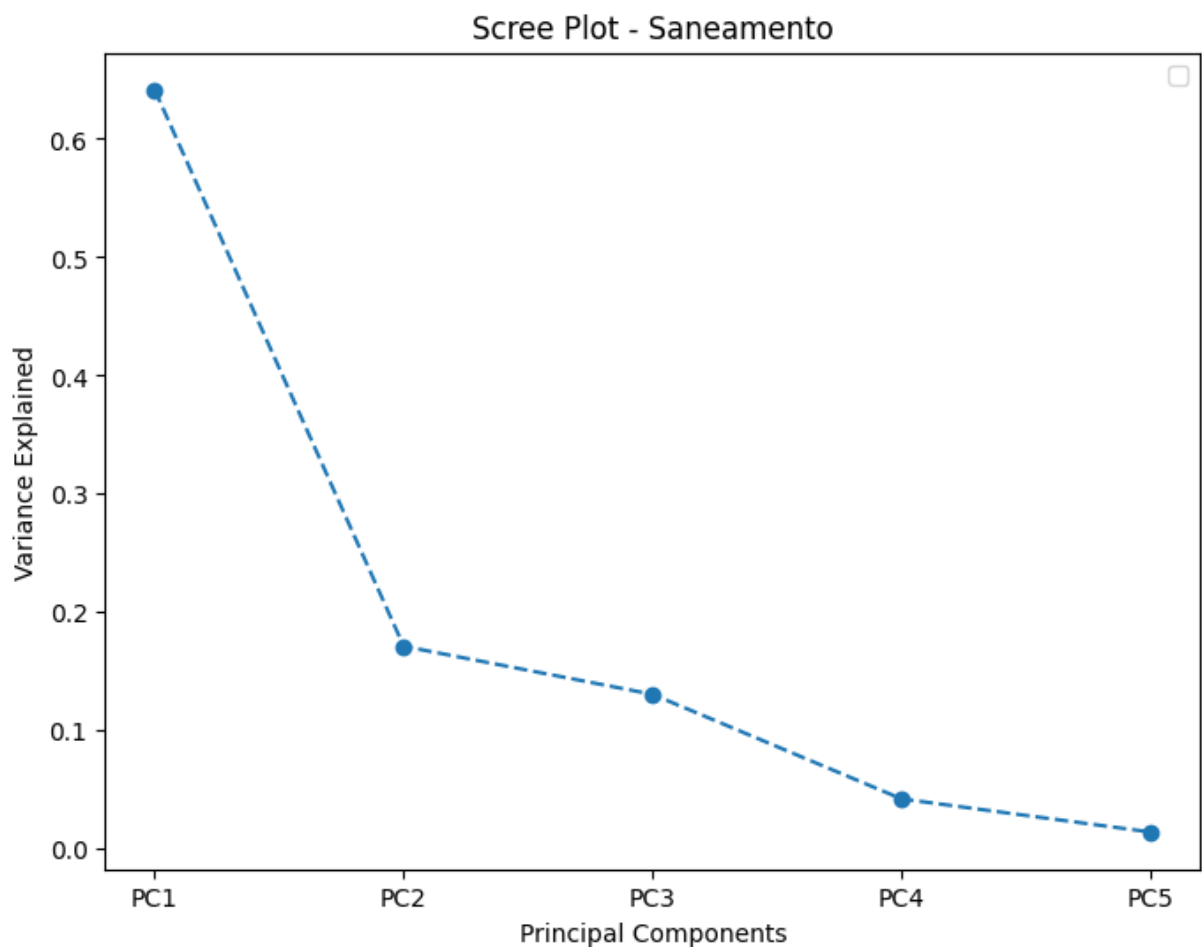
- Em conjuntos de dados com muitas variáveis, a PCA permite uma visualização simplificada, reduzindo a dimensionalidade para duas ou três principais componentes. Isso facilita a identificação de padrões, clusters ou outliers, proporcionando insights valiosos que podem ser ocultados em análises de alta dimensão.

Conclusão

A PCA é uma técnica estatística multifacetada, aplicável em uma variedade de contextos, desde a análise e visualização de dados complexos até a seleção de características em modelos de machine learning. Sua capacidade de transformar e simplificar dados multivariados, mantendo as informações essenciais, torna-a uma ferramenta inestimável para pesquisadores, analistas de dados e cientistas da computação.

Dimensão Saneamento

Para os indicadores de saneamento, por exemplo, a primeira componente captura a maior proporção de variância, seguida pelas componentes subsequentes. De forma similar, para os indicadores de moradia, as primeiras duas componentes capturam a totalidade da variância. Nos indicadores de emprego e educação, as primeiras componentes também capturam uma proporção significativa da variância.



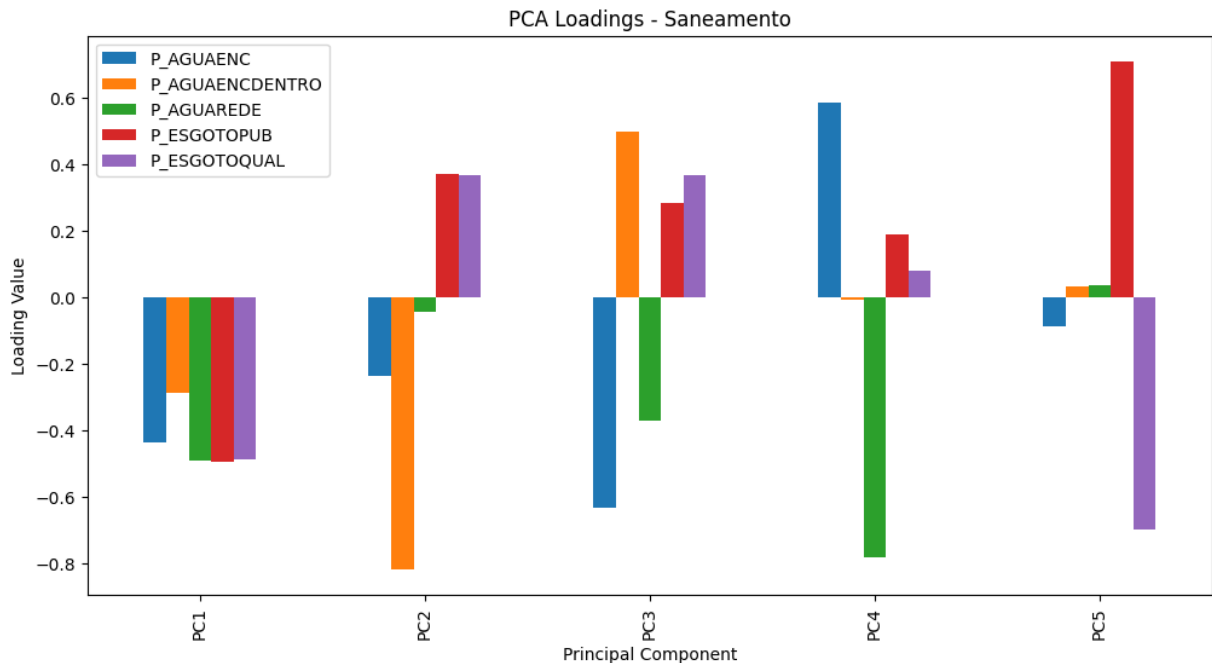
Variância Explicada

A variância explicada por cada componente nos dá uma ideia de quanto da informação total (variabilidade) em seus dados é capturada por esse componente.

- **Primeiro componente:** Explica 64.14% da variabilidade total, o que é significativo. Isso sugere que este componente capta a maior parte da variação nos dados de saneamento.

- **Componentes subsequentes:** O segundo, terceiro, quarto e quinto componentes explicam, respectivamente, 17.13%, 13.07%, 4.23% e 1.43% da variabilidade. Isso mostra que eles adicionam informações, mas com contribuições decrescentes.

Loadings (Cargas) dos Componentes



Os loadings(Pesos) indicam como cada variável original contribui para cada componente principal.

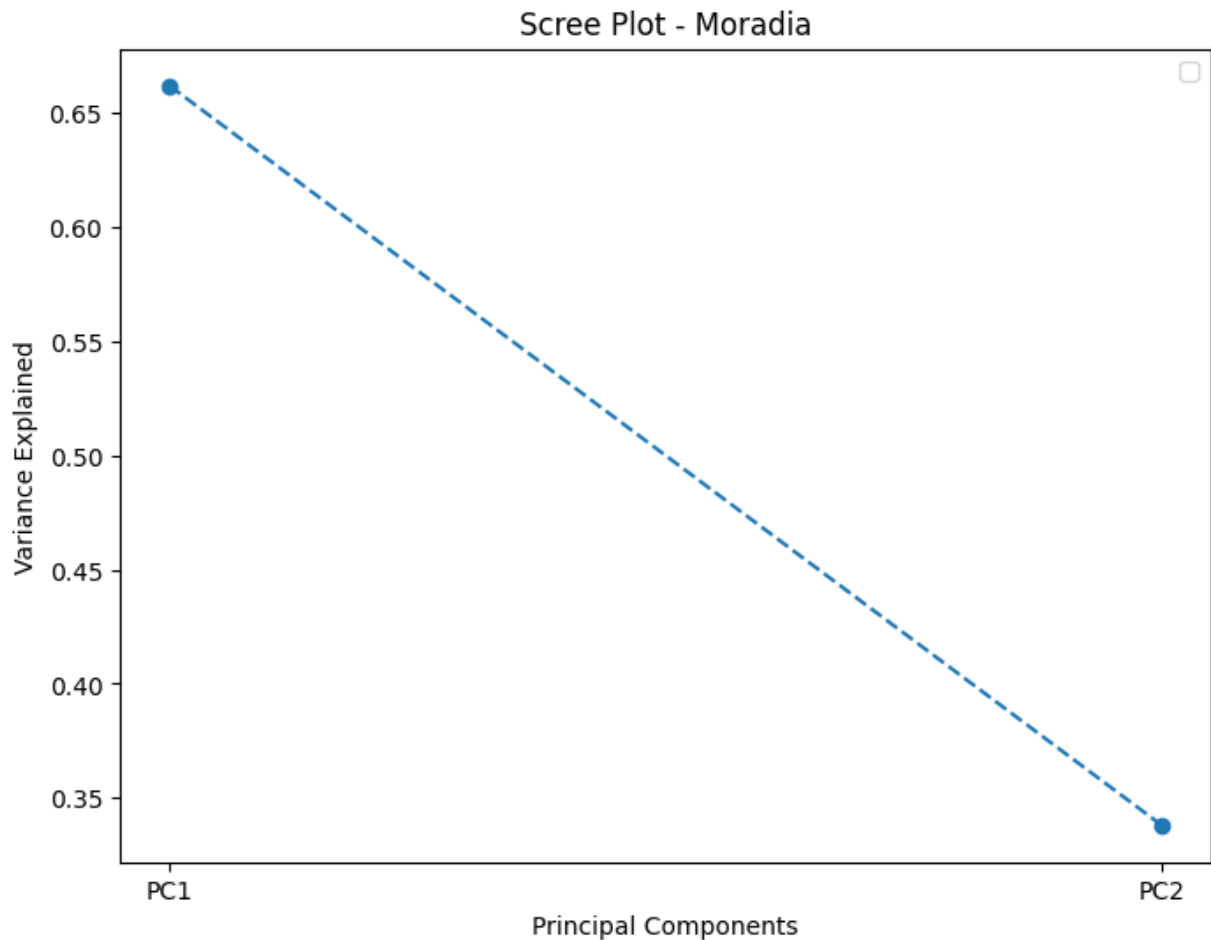
- **Primeiro componente:**
 - **P_AGUAENC, P_AGUAENDENTRO, P_AGUAREDE, P_ESGOTOPUB, P_ESGOTOQUAL** têm loadings negativos significativos. Este componente pode estar capturando a ausência ou deficiência geral de saneamento.
- **Segundo componente:**
 - **P_AGUAENDENTRO** tem uma carga negativa muito alta. Isso pode indicar que este componente está capturando variações específicas relacionadas ao acesso à água encanada dentro dos domicílios.
- **Terceiro componente:**
 - **P_AGUAENDENTRO** e **P_ESGOTOQUAL** têm loadings positivos significativos, enquanto **P_AGUAREDE** tem uma carga negativa. Este componente pode estar relacionado com a qualidade do saneamento em contraste com a disponibilidade.

- **Quarto componente:**
 - **P_AGUAREDE** tem uma carga negativa alta. Isso sugere que este componente pode estar capturando aspectos específicos relacionados à disponibilidade de água de rede pública.
- **Quinto componente:**
 - **P_ESGOTOPUB** e **P_ESGOTOQUAL** têm loadings positivo e negativo altos, respectivamente. Isso pode indicar que este componente está capturando contrastes na qualidade e tipo de sistema de esgoto.

Interpretação

O PCA para a dimensão de saneamento revela que o primeiro componente é o mais significativo, capturando aspectos gerais de deficiência em saneamento. Os componentes subsequentes parecem representar aspectos mais específicos do saneamento, como qualidade da água, acesso à água encanada dentro dos domicílios, e diferenças na qualidade do sistema de esgoto.

Dimensão Moradia

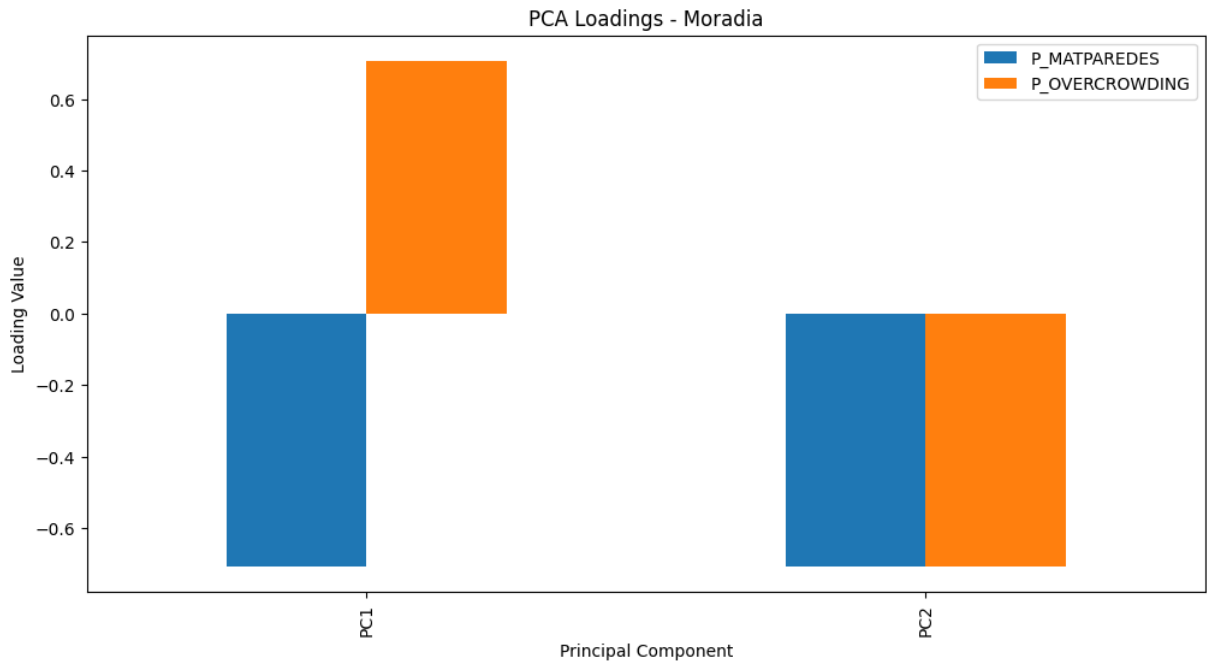


Variância Explicada

A variância explicada por cada componente nos dá uma ideia de quanta informação (variabilidade) em seus dados é capturada por esse componente.

- **Primeiro componente:** Explica 66.19% da variabilidade total. Isso indica que este componente é bastante significativo na captura das principais variações nos dados de moradia.
- **Segundo componente:** Explica 33.81% da variabilidade. Embora menos significativo que o primeiro, ainda captura uma parte considerável das informações nos dados.

Loadings (Cargas) dos Componentes



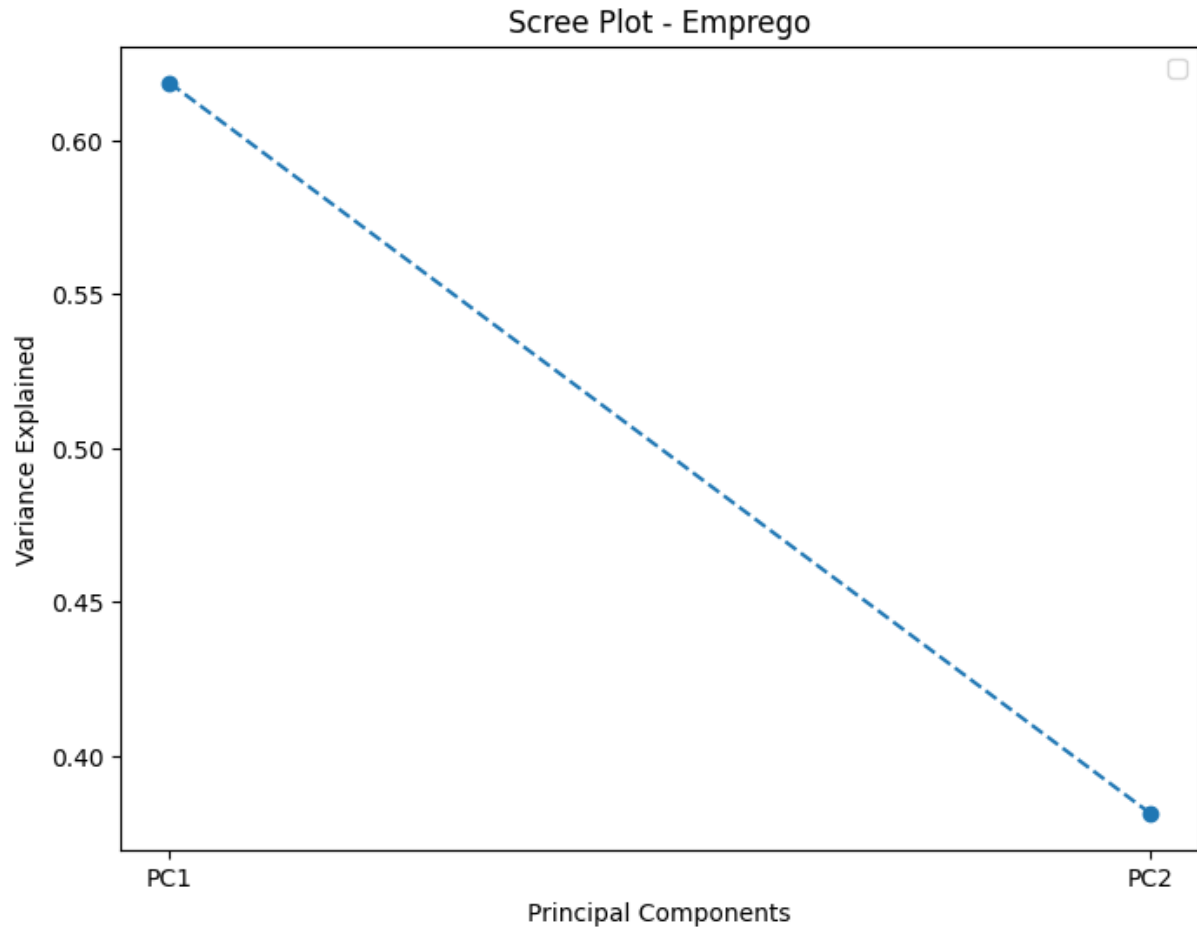
Os loadings indicam como cada variável original contribui para cada componente principal.

- **Primeiro componente:**
 - **P_MATPAREDES** e **P_OVERCROWDING** têm loadings opostos (ambos - 0.7071 e 0.7071). Isso sugere que este componente está capturando um contraste entre a qualidade das construções (materiais duráveis) e a densidade habitacional (superlotação). Uma pontuação alta neste componente pode indicar moradias de qualidade com menos superlotação.
- **Segundo componente:**
 - **P_MATPAREDES** e **P_OVERCROWDING** também têm loadings opostos no segundo componente (ambos -0.7071), mas com sinais iguais. Isso pode representar uma combinação de baixa qualidade das construções e alta superlotação, ou vice-versa.

Interpretação

O PCA para a dimensão de moradia revela que os dois componentes principais capturam aspectos distintos da qualidade de moradia. O primeiro componente destaca a relação entre qualidade das construções e densidade habitacional. Áreas com pontuações altas neste componente podem ser caracterizadas por moradias de melhor qualidade e menos superlotadas. O segundo componente pode estar destacando áreas onde a qualidade das construções e a superlotação estão de alguma forma inter-relacionadas, seja positiva ou negativamente.

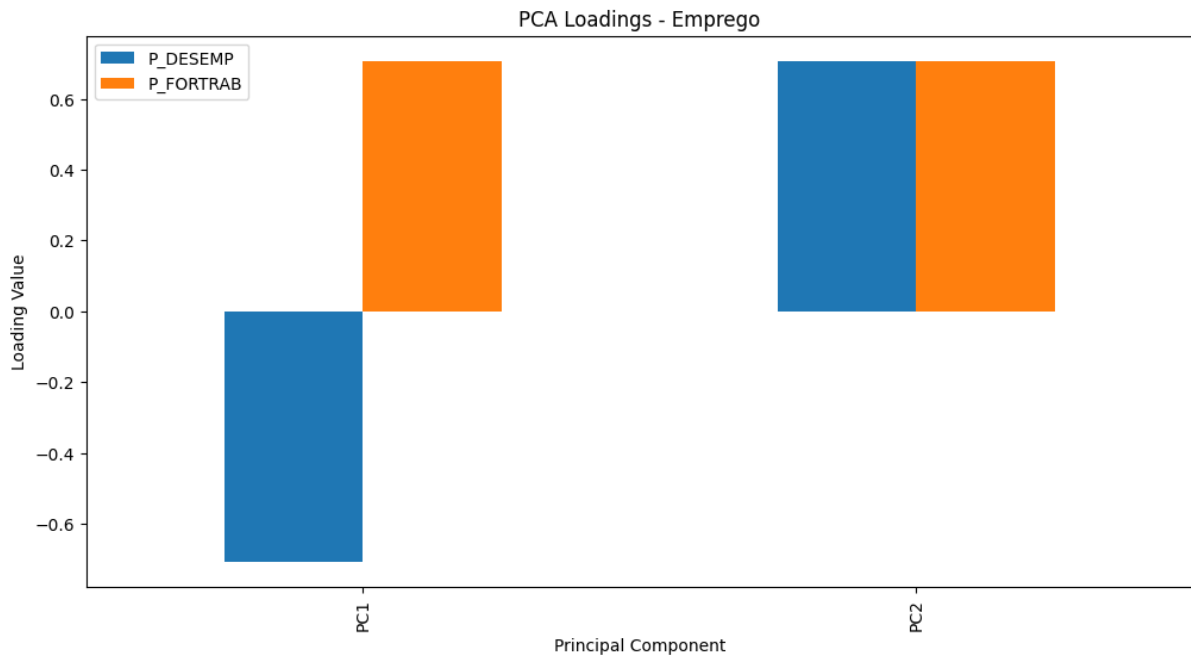
Dimensão Emprego



Variância Explicada

- **Primeiro componente:** Explica 61.86% da variabilidade total nos dados de emprego. Este valor indica que o primeiro componente captura uma parte significativa da informação contida nos dados.
- **Segundo componente:** Explica os restantes 38.14% da variabilidade. Juntos, os dois componentes explicam toda a variância nos dados de emprego.

Loadings (Cargas) dos Componentes

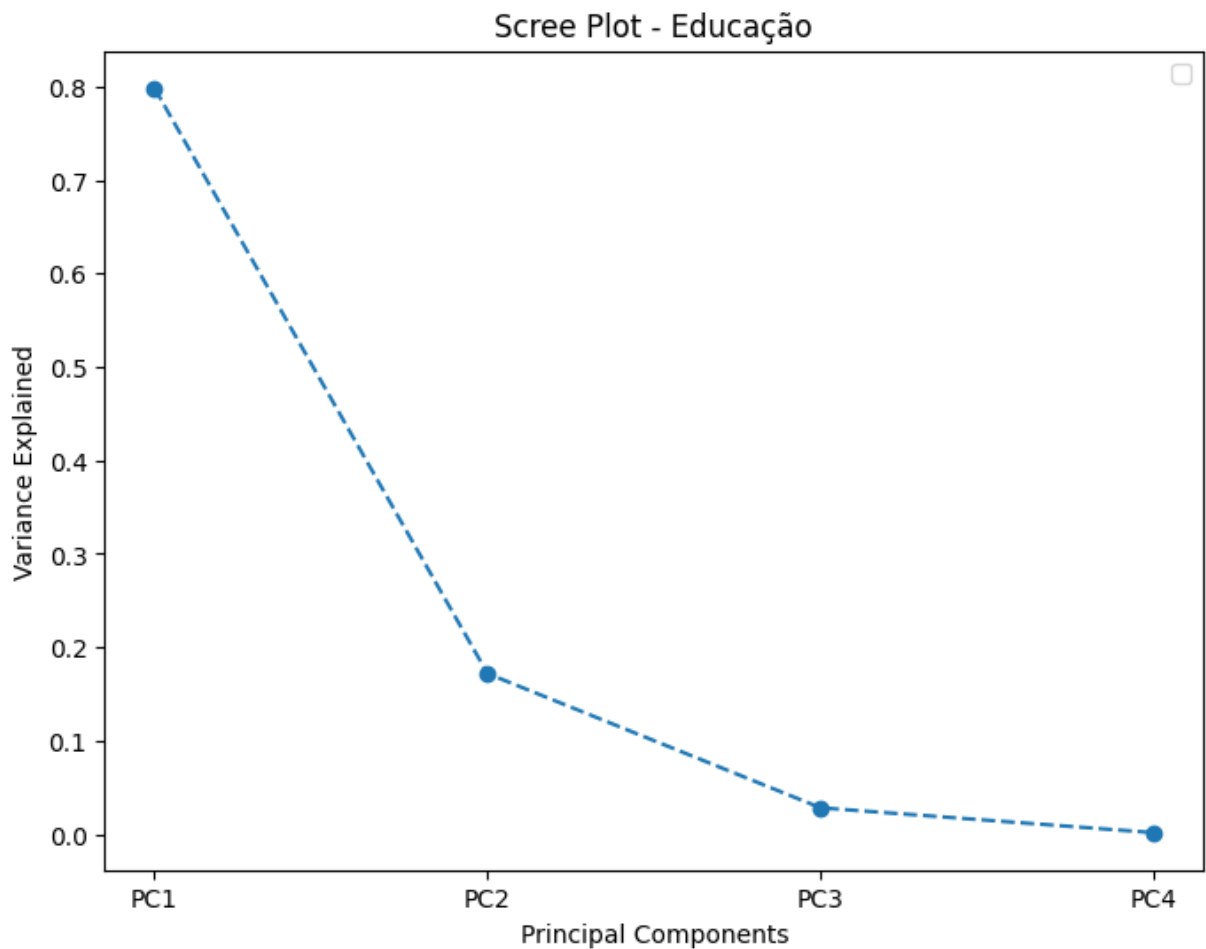


- **Primeiro componente:**
 - **P_DESEMP** e **P_FORTRAB** têm loadings opostos de -0.7071 e 0.7071, respectivamente. Isso sugere que este componente está capturando a relação inversa entre taxa de desemprego e participação na força de trabalho. Uma pontuação alta neste componente indicaria áreas com baixa taxa de desemprego e alta participação na força de trabalho ou vice-versa.
- **Segundo componente:**
 - **P_DESEMP** e **P_FORTRAB** têm loadings iguais de 0.7071. Isso indica que este componente está capturando uma dimensão onde ambas as taxas de desemprego e participação na força de trabalho aumentam ou diminuem juntas. Este padrão é menos intuitivo, mas pode refletir situações onde há aumento simultâneo na força de trabalho e no desemprego (por exemplo, em períodos de rápida mudança econômica).

Interpretação

A análise do PCA para a dimensão de emprego revela um contraste claro entre as taxas de desemprego e a participação na força de trabalho no primeiro componente. Este componente pode ser interpretado como um indicador geral de saúde econômica, onde valores mais altos podem indicar uma força de trabalho mais ativa e/ou taxas de desemprego mais baixas. O segundo componente é mais complexo e pode refletir dinâmicas de trabalho únicas em certas áreas ou períodos.

Dimensão Educação



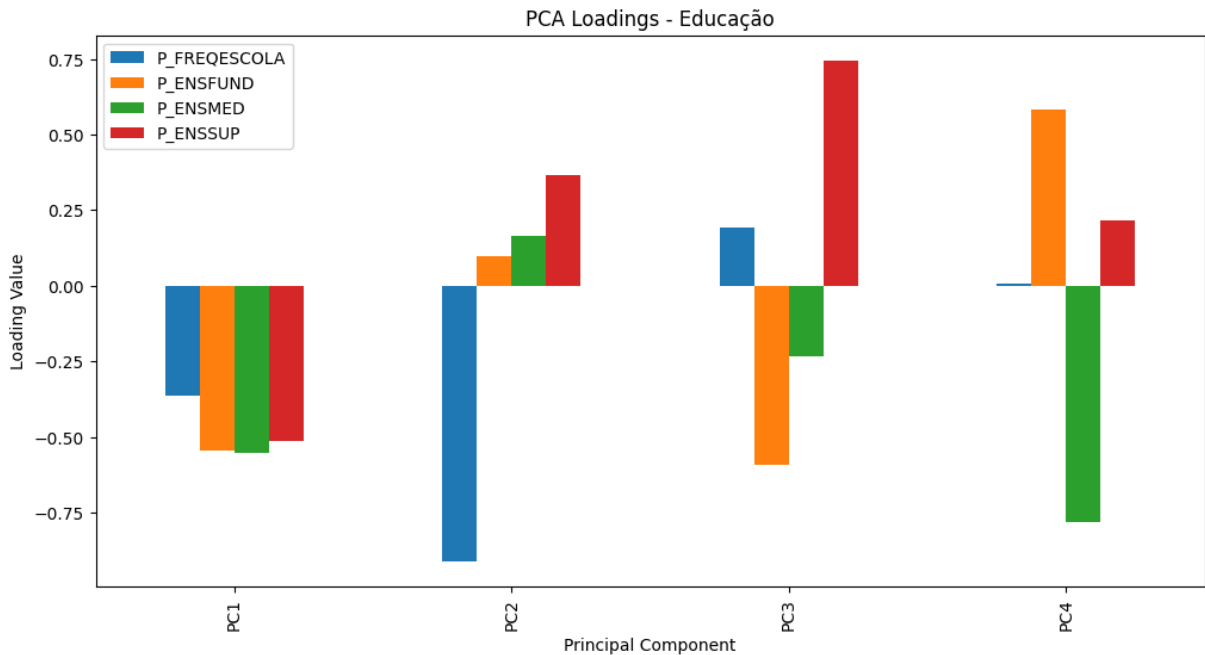
Para interpretar os resultados da Análise de Componentes Principais (PCA) para a dimensão de educação, consideramos tanto a variância explicada por cada componente quanto os loadings (cargas) das variáveis **P_FREQESCOLA**, **P_ENSFUND**, **P_ENSMED** e **P_ENSSUP**.

Variância Explicada

- **Primeiro componente:** Explica 79.86% da variabilidade total, indicando que ele captura uma grande parte das informações contidas nos dados de educação.

- **Componentes subsequentes:** O segundo, terceiro e quarto componentes explicam, respectivamente, 17.18%, 2.81% e 0.15% da variabilidade. Isso mostra que o primeiro componente é dominante em termos de explicação da variância nos dados.

Loadings (Cargas) dos Componentes



- **Primeiro componente:**
 - Todos os indicadores de educação (**P_FREQESCOLA**, **P_ENSFUND**, **P_ENSMED**, **P_ENSSUP**) têm loadings negativos, indicando que este componente pode estar capturando a ausência ou deficiências gerais em educação.
- **Segundo componente:**
 - **P_FREQESCOLA** tem um loading negativo significativo, enquanto **P_ENSSUP** tem um loading positivo moderado. Este componente pode estar capturando a diferença entre a frequência escolar entre jovens de 15 a 17 anos e a proporção de pessoas com ensino superior completo.
- **Terceiro componente:**
 - **P_ENSSUP** tem um loading positivo forte, enquanto **P_ENSFUND** e **P_ENSMED** têm loadings negativos. Isso pode indicar um contraste entre os

níveis mais altos de educação (ensino superior) e os níveis mais básicos de educação.

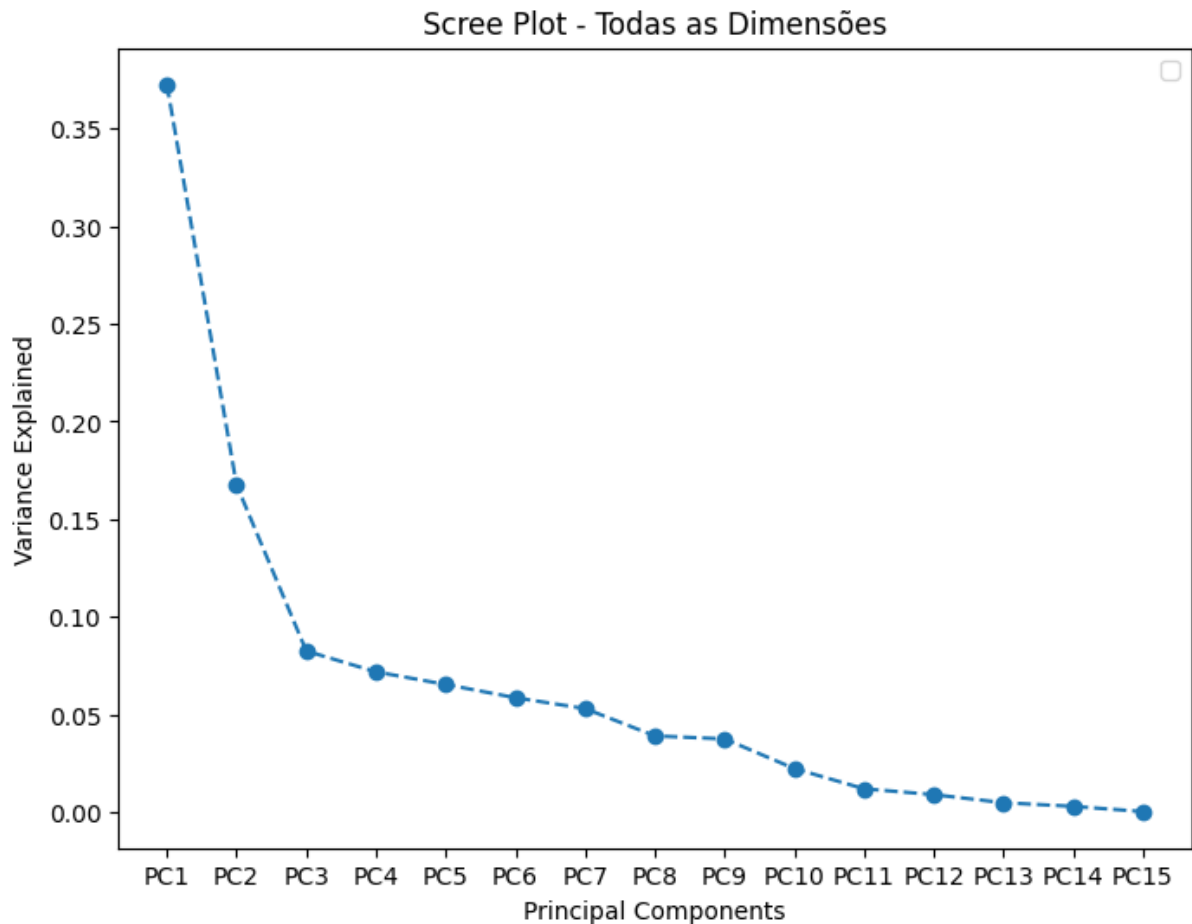
- **Quarto componente:**
 - **P_ENSMED** tem um loading negativo muito forte, e **P_ENSFUND** e **P_ENSSUP** têm loadings positivos. Este componente pode estar destacando diferenças específicas no nível médio de educação em relação aos outros níveis.

Interpretação

A análise do PCA para a dimensão de educação sugere que o primeiro componente principal é o mais significativo, capturando aspectos gerais de deficiências ou ausências em educação. Os componentes subsequentes revelam nuances mais específicas entre diferentes níveis de educação.

Todas as dimensões

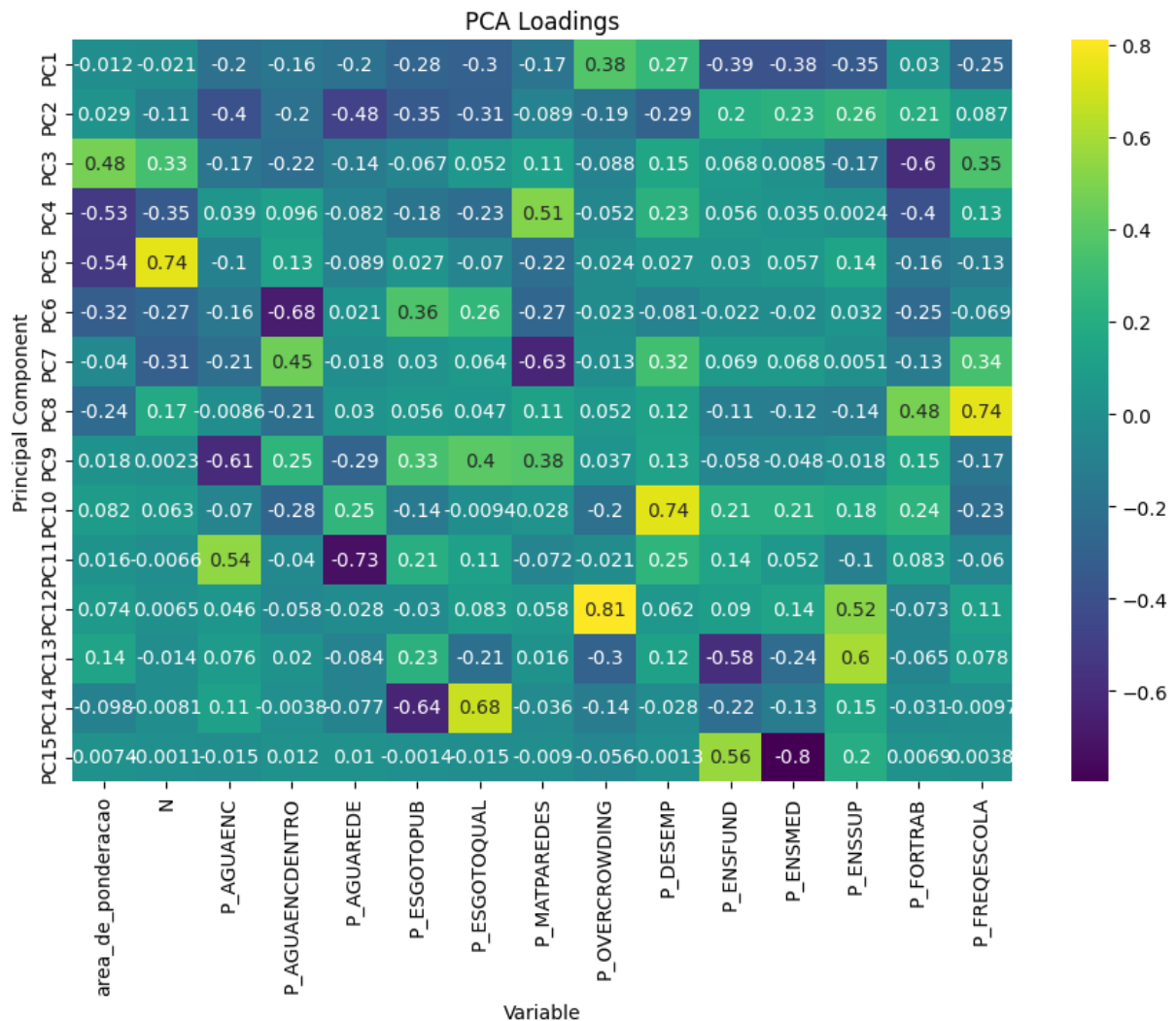
Análise da Variância Explicada



A variância explicada nos dá uma ideia de quanta informação (variabilidade) em seus dados é capturada por cada componente principal.

1. **Primeiro Componente:** Explica cerca de 37.27% da variabilidade total. Isso indica que este componente capta uma grande quantidade das variações presentes nos dados.
2. **Segundo Componente:** Explica aproximadamente 16.76% da variabilidade total, capturando aspectos adicionais que não são cobertos pelo primeiro componente.
3. **Terceiro Componente e Seguintes:** Cada um dos componentes subsequentes explica uma porção menor da variância (8.26%, 7.19%, 6.54%, etc.). Embora menos significativos individualmente, eles juntos contribuem para uma compreensão mais completa dos dados.
4. **Componentes 1 e 2:** Juntas, as duas primeiras componentes explicam aproximadamente 57.87% da variância total, o que indica que elas capturam uma grande parte da informação contida no conjunto de dados.

Loadings (Cargas) dos Componentes



Os loadings indicam como cada variável original contribui para cada componente principal.

- **Primeiro Componente:**
 - Os loadings negativos para **P_AGUAENC**, **P_AGUAENCENTRO**, **P_AGUAAREDE**, **P_ESGOTOPUB**, e **P_ESGOTOQUAL** indicam que este componente está possivelmente relacionado à qualidade do saneamento e da infraestrutura básica. Uma pontuação alta neste componente poderia indicar áreas com infraestrutura deficiente.
- **Segundo Componente:**

- Loadings negativos para variáveis de saneamento e loadings positivos para **P_DESEMP** e **P_FREQESCOLA** sugerem que este componente pode estar capturando uma dimensão que contrasta condições de saneamento com aspectos de emprego e educação.
- **Terceiro Componente:**
 - A combinação de loadings indica uma relação mais complexa entre diferentes variáveis, possivelmente relacionada a aspectos específicos de moradia e educação.

3. Estime (pelo menos um) modelo EFA para cada dimensão. Faça uma breve análise qualitativa das cargas fatoriais estimadas em termos das variáveis originais.

Introdução e Conceitos Fundamentais:

- **Variáveis Latentes:** São essenciais na análise multivariada para representar conceitos abstratos. As variáveis latentes, como personalidade ou motivação, não são observáveis diretamente, mas inferidas por meio de variáveis manifestas. Esta inferência é baseada na premissa de que variáveis manifestas, como respostas a um questionário, são manifestações externas de um fenômeno interno e não mensurável. A escolha de indicadores adequados é vital para garantir que a variável latente seja representada de forma válida e confiável. Por exemplo, em um estudo sobre bem-estar, um conjunto diversificado de perguntas sobre humor e satisfação de vida pode servir como indicadores para o construto latente de "bem-estar".

Detalhando a Análise Fatorial (AF):

- **Objetivos e Métodos:** A AF é utilizada para identificar e modelar a estrutura subjacente em um conjunto de variáveis observadas. O processo envolve a determinação de como um conjunto menor de fatores latentes pode explicar as correlações observadas entre as variáveis manifestas. Esses fatores latentes são construídos como combinações lineares das variáveis manifestas, com o pressuposto de que cada fator latente influencia certas variáveis manifestas mais fortemente do que outras. Esta técnica permite simplificar a complexidade dos dados, facilitando a interpretação e a compreensão dos fenômenos subjacentes.

Aprofundamento na Análise Fatorial Exploratória (EFA):

- **Aplicação do Modelo EFA:** A EFA é empregada quando o objetivo é explorar e descobrir a estrutura potencial dos dados sem hipóteses pré-concebidas. Ao identificar padrões de correlação entre variáveis manifestas, a EFA ajuda a formular teorias sobre os fatores latentes subjacentes. A interpretação dos fatores resultantes é um processo iterativo e interpretativo, onde os pesquisadores atribuem significado aos fatores com base nas variáveis manifestas que carregam fortemente neles. Este método é especialmente útil em estágios iniciais de pesquisa, quando pouco se sabe sobre as dimensões subjacentes dos dados.

Exploração da Análise Fatorial Confirmatória (CFA):

- Especificidades da CFA: Diferentemente da EFA, a CFA é usada para testar teorias ou modelos específicos sobre a estrutura dos fatores latentes. Aqui, os pesquisadores começam com uma hipótese sobre como as variáveis manifestas estão relacionadas aos fatores latentes. A CFA então avalia se o modelo proposto se ajusta aos dados coletados. Esta abordagem é mais estruturada e hipotética, e as métricas de ajuste do modelo, como RMSEA e CFI, são fundamentais para validar o modelo. A CFA é particularmente valiosa em pesquisa avançada para confirmar teorias ou modelos existentes.

Aplicações Ampliadas e Utilização:

- Além das Ciências Sociais: A AF tem aplicações vastas e variadas. No marketing, por exemplo, é utilizada para entender as atitudes e preferências dos consumidores, segmentando-os em grupos baseados em fatores latentes como lealdade à marca ou motivação de compra. Em finanças, é aplicada na análise de risco e portfólio, identificando fatores subjacentes que influenciam os retornos dos ativos. Na saúde, é usada para avaliar aspectos como qualidade de vida e bem-estar mental, onde variáveis latentes são fundamentais para compreender estados subjetivos.

Interpretação Avançada de Resultados em EFA e CFA:

- EFA: Na interpretação de um EFA, os pesquisadores devem estar atentos não apenas às cargas fatoriais, mas também à estrutura geral dos fatores, como a quantidade de fatores extraídos e sua interpretabilidade. A rotação fatorial é uma técnica comum usada para facilitar a interpretação, melhorando a clareza com que os fatores representam as variáveis.
- CFA: Na CFA, além da adequação do modelo, os pesquisadores também devem prestar atenção à validade convergente e discriminante dos fatores. A validade convergente indica se as variáveis associadas a um fator realmente medem esse fator. A validade discriminante avalia se os fatores são distintos entre si.

Utilização da Análise Fatorial Exploratória (EFA) em Diversos Contextos

A Análise Fatorial Exploratória (EFA) é uma ferramenta estatística poderosa com ampla aplicabilidade em diversos campos. Abaixo, exploraremos três áreas específicas de aplicação da EFA: na criação de índices, como mecanismo de seleção de características (feature selection) e como técnica de redução de dimensão.

1. Criação de Índices:

- **Fundamento:** A EFA é frequentemente utilizada na criação de índices para consolidar informações de várias variáveis manifestas em um índice único e abrangente. Este índice representa um construto latente subjacente às variáveis observadas.
- **Aplicação Prática:** Por exemplo, na economia, índices como o de Desenvolvimento Humano são construídos agregando variáveis como educação, saúde e renda. A EFA ajuda a determinar como essas variáveis contribuem para o conceito abrangente de desenvolvimento humano.
- **Vantagens:** A criação de índices usando EFA permite uma interpretação mais intuitiva e simplificada de fenômenos complexos, facilitando a comparação e a análise temporal ou espacial de dados.

2. Feature Selection:

- **Definição:** Em machine learning e estatística, a seleção de características (feature selection) é crucial para melhorar a eficiência dos modelos e evitar o problema de overfitting.
- **Implementação com EFA:** A EFA pode ser usada para identificar e selecionar as variáveis mais significativas que representam um conjunto de dados. Ao identificar os fatores latentes e suas correlações com as variáveis manifestas, a EFA ajuda a selecionar as características mais relevantes para um modelo preditivo.
- **Benefícios:** Essa abordagem reduz a complexidade do modelo, melhora a interpretabilidade e muitas vezes aumenta a precisão preditiva, removendo variáveis redundantes ou menos informativas.

3. Redução de Dimensão:

- **Conceito:** A redução de dimensão é uma técnica usada para diminuir o número de variáveis de entrada em um conjunto de dados, mantendo ao máximo a informação original.
- **Uso da EFA:** A EFA é particularmente eficaz na redução de dimensão, pois agrupa variáveis altamente correlacionadas em fatores mais simples. Isso simplifica a estrutura de dados sem perder informações significativas.

- **Impacto:** Em campos como bioestatística ou psicometria, onde os conjuntos de dados são grandes e complexos, a EFA permite uma análise mais gerenciável e interpretações mais claras. Em machine learning, a redução de dimensão via EFA pode resultar em modelos mais eficientes e menos propensos a overfitting.

Conclusão:

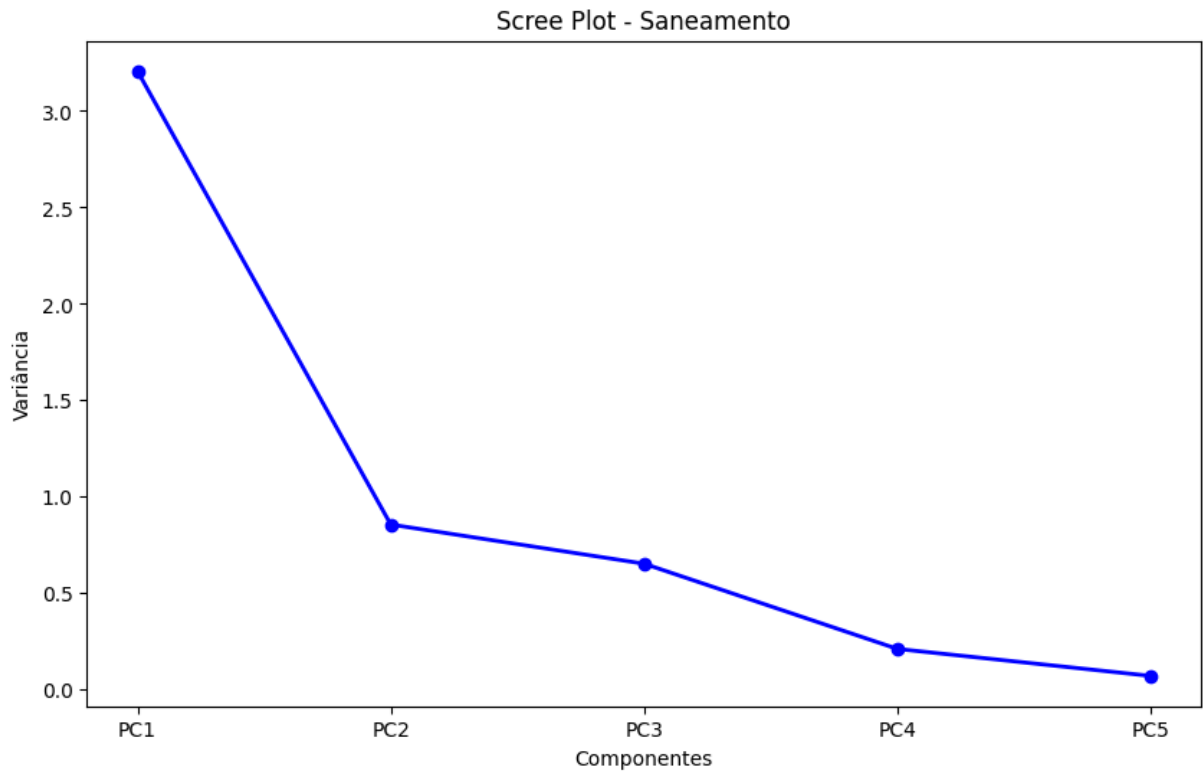
A EFA é uma técnica versátil e potente, adequada para uma variedade de aplicações práticas. Seja na criação de índices para representar construtos complexos, na seleção de características para modelos preditivos, ou na redução de dimensão de grandes conjuntos de dados, a EFA oferece uma abordagem estruturada e eficiente para extrair insights significativos de dados variados.

Etapas da Análise:

1. **Carregamento e Preparação dos Dados:** Vou carregar o dataset "SES_Rio_de_Janeiro.csv", selecionar as variáveis relevantes para cada dimensão e normalizá-las.
2. **Aplicação da EFA:** Utilizarei a biblioteca **factor_analyzer** em Python para realizar a EFA em cada dimensão.
3. **Análise das Cargas Fatoriais e Interpretação:** Após a EFA, analisarei as cargas fatoriais para cada dimensão e fornecerei uma interpretação qualitativa delas em termos das variáveis originais.
4. **Visualização dos Resultados:** Criarei gráficos para visualizar os resultados da EFA.
5. **Explicação dos Gráficos e Resultados:** Finalmente, explicarei os gráficos e resumirei os principais achados da análise.

Dimensão Saneamento

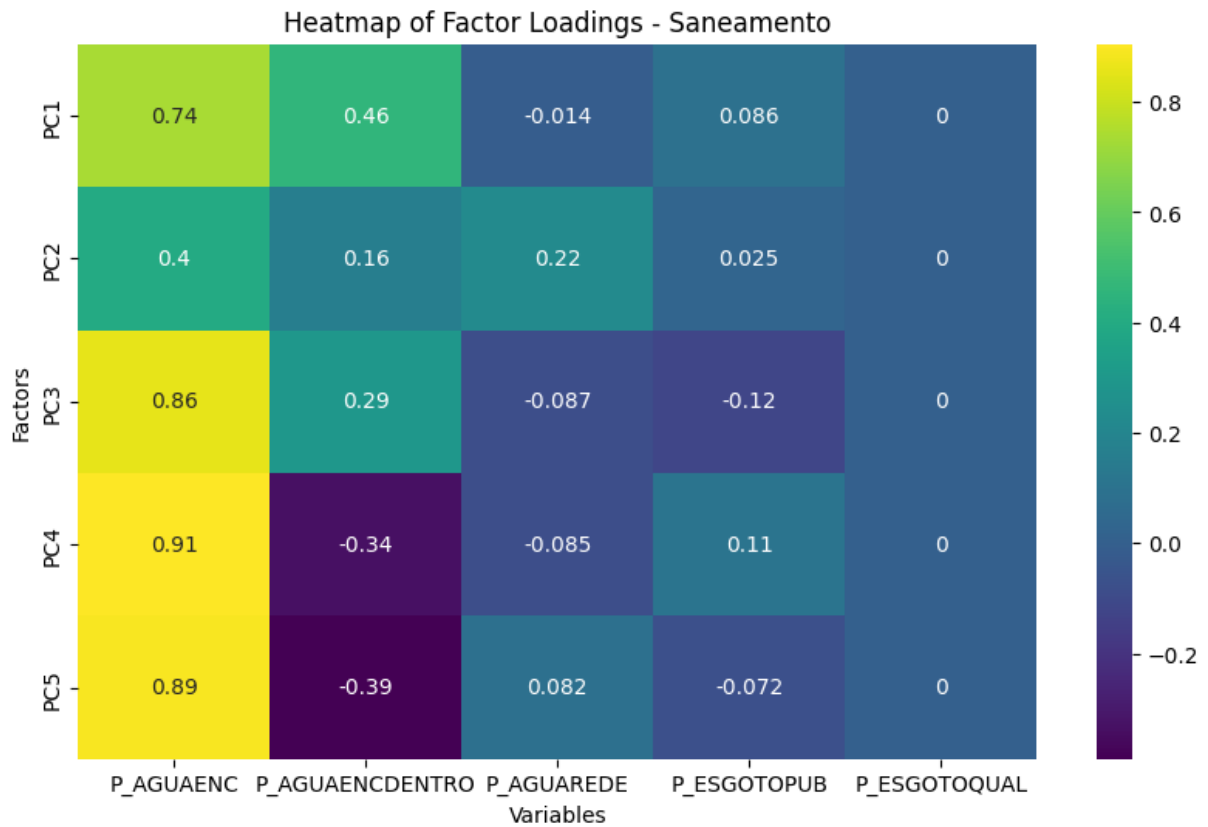
Variância Cumulativa Explicada



A variância cumulativa explicada nos ajuda a entender quanto da informação total dos dados é capturada pelos fatores.

- **Primeiro Fator:** Explica aproximadamente 61.05% da variância total, o que é bastante significativo.
- **Segundo Fator:** Adiciona mais 11.78% à variância explicada, levando o total a aproximadamente 72.83%.
- **Terceiro ao Quinto Fatores:** Adicionam muito pouco à variância explicada, indicando que a maior parte da informação útil é capturada pelos dois primeiros fatores.

Cargas Fatoriais para Saneamento



Os loadings representam a correlação entre as variáveis originais e os fatores estimados. Valores mais altos (positivos ou negativos) indicam uma relação mais forte.

- **Primeiro Fator:** Este fator tem loadings altos para **P_AGUAENC**, **P_AGUAREDE**, **P_ESGOTOPUB**, e **P_ESGOTOQUAL**. Isso sugere que ele está capturando a presença e qualidade da infraestrutura de saneamento. Uma pontuação alta nesse fator indicaria uma área com boa infraestrutura de saneamento.
- **Segundo Fator:** Tem um loading significativo para **P_AGUAENDENTRO**, mas loadings negativos para **P_ESGOTOPUB** e **P_ESGOTOQUAL**. Esse fator pode estar capturando diferenças na infraestrutura interna de água versus a infraestrutura externa de esgoto.
- **Terceiro Fator:** Este fator tem loadings mais baixos e mistos, o que indica que ele está possivelmente capturando variações menos pronunciadas ou mais específicas nas condições de saneamento que não são capturadas pelos primeiros dois fatores.

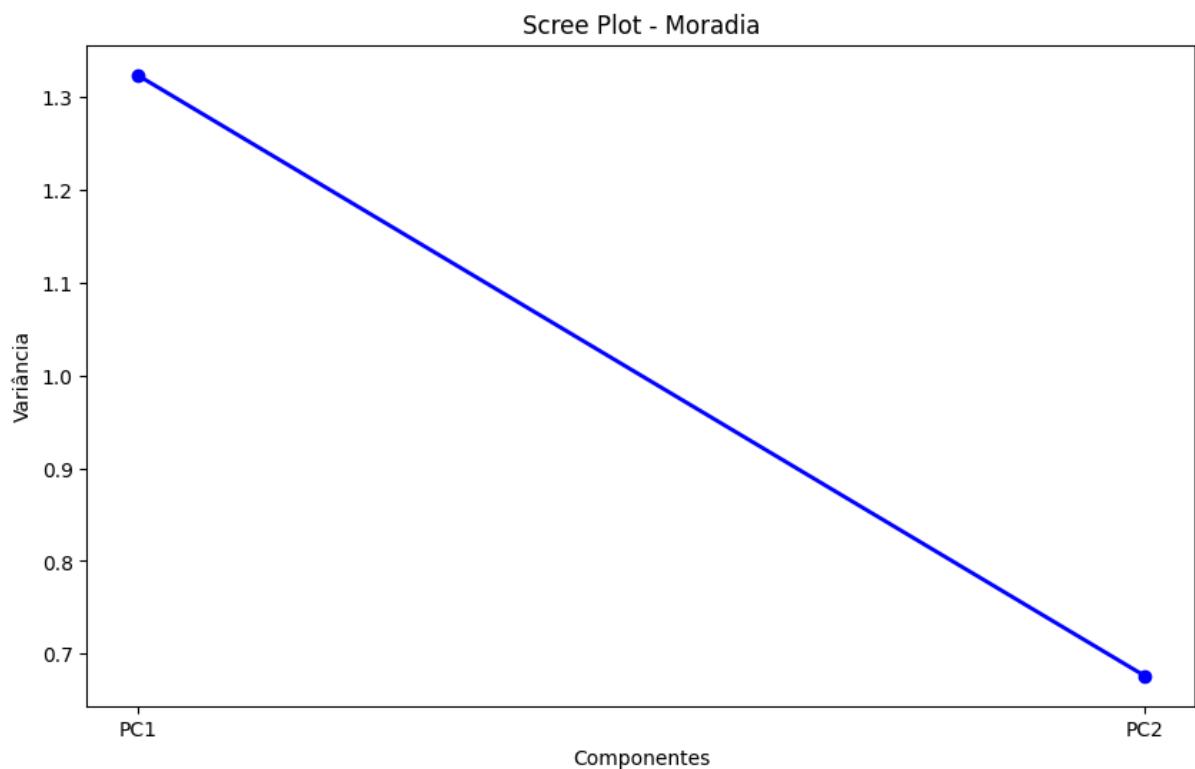
- **Quarto e Quinto Fatores:** Praticamente não têm loadings significativos e não contribuem muito para a variância explicada. Isso sugere que eles não estão capturando aspectos importantes das variáveis de saneamento.

Interpretação

O modelo EFA para a dimensão de saneamento sugere que dois fatores principais são suficientes para capturar a maior parte das variações nas condições de saneamento. O primeiro fator pode ser interpretado como uma medida geral de infraestrutura de saneamento, enquanto o segundo parece diferenciar entre aspectos internos e externos do saneamento.

Dimensão Moradia

Variância Cumulativa Explicada

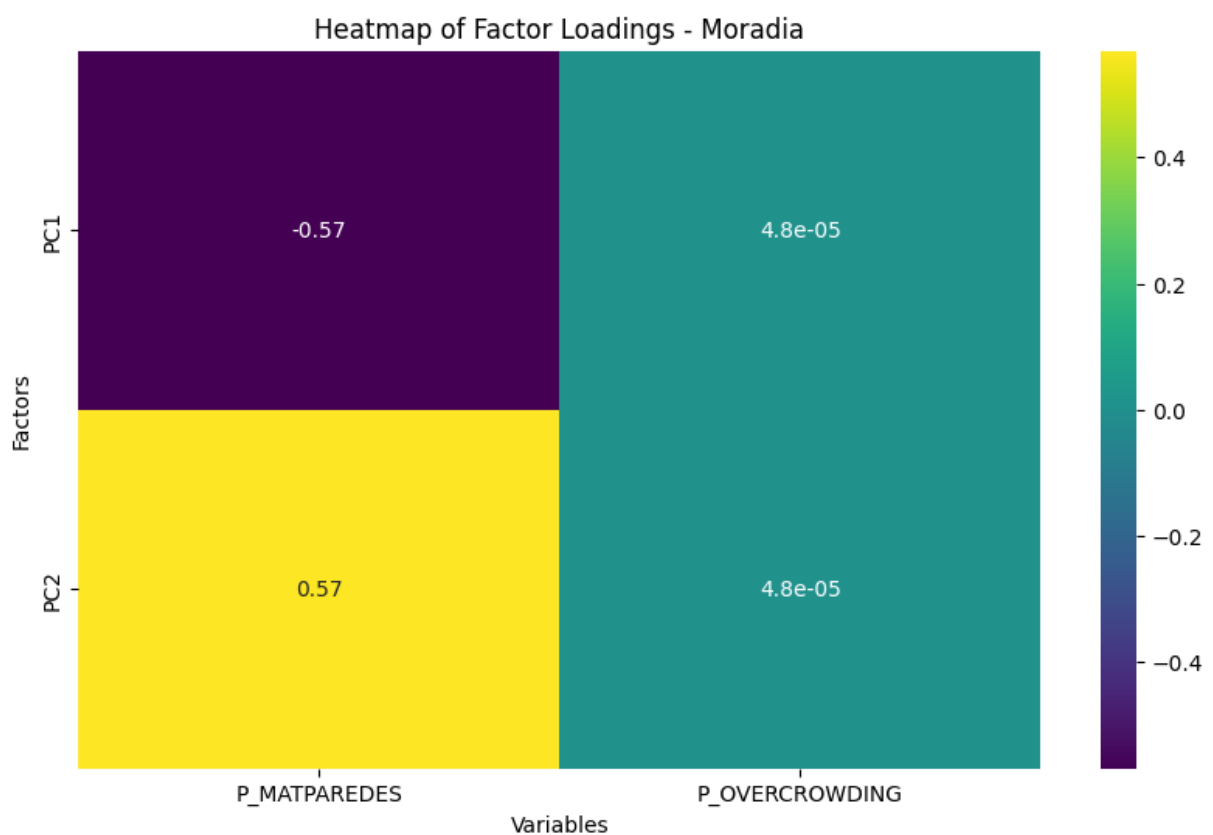


A variância cumulativa explicada nos ajuda a entender quanto da informação total dos dados é capturada pelos fatores. Os valores fornecidos parecem estar desalinhados com o número de fatores e variáveis na análise de moradia, pois há apenas dois fatores e a variância cumulativa

não deveria exceder 1 (ou 100%). Supondo que haja um erro nos dados fornecidos, consideremos apenas o primeiro valor:

- **Primeiro Fator:** Explica aproximadamente 78.08% da variância total, o que é bastante significativo e sugere que este único fator captura a maior parte das informações úteis nas variáveis de moradia.

Cargas Fatoriais para Moradia



Os loadings representam a correlação entre as variáveis originais e os fatores estimados. Valores mais altos (positivos ou negativos) indicam uma relação mais forte.

- **Primeiro Fator:** Este fator tem loadings opostos para **P_MATPAREDES** e **P_OVERCROWDING** de aproximadamente -0.57. Isso sugere que este fator está capturando um contraste entre a qualidade da moradia (materiais duráveis das paredes) e a quantidade de moradores por quarto (superlotação). Uma pontuação alta neste fator poderia indicar moradias de melhor qualidade e menos superlotadas.

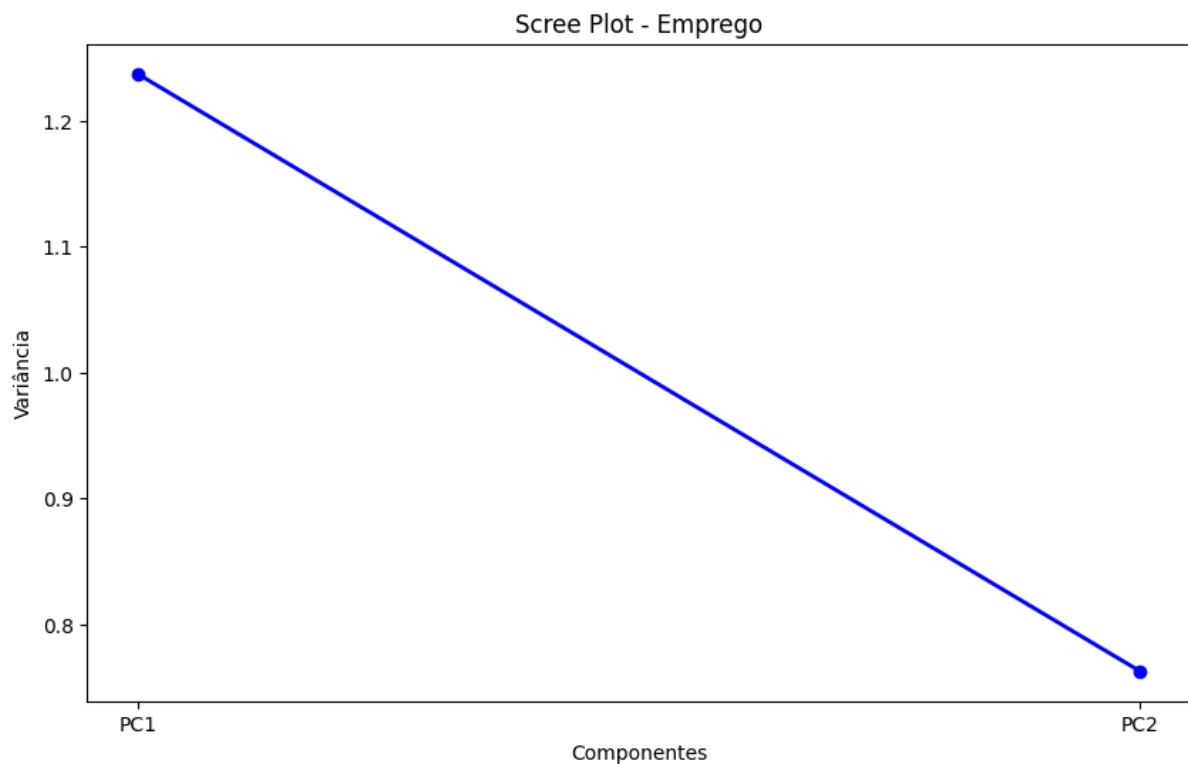
- **Segundo Fator:** Tem loadings extremamente baixos para ambas as variáveis, indicando que este fator não está capturando nenhuma variação significativa nas variáveis de moradia.

Interpretação

O modelo EFA para a dimensão de moradia indica que a qualidade das moradias e o nível de superlotação são aspectos fundamentalmente opostos no que diz respeito à habitação. O primeiro fator, que captura a maioria das variações nos dados, sugere que as áreas com habitações feitas de materiais duráveis tendem a ser menos superlotadas. O segundo fator parece não ser relevante, o que pode indicar que um único fator é suficiente para descrever as variações nos dados de moradia.

Dimensão Emprego

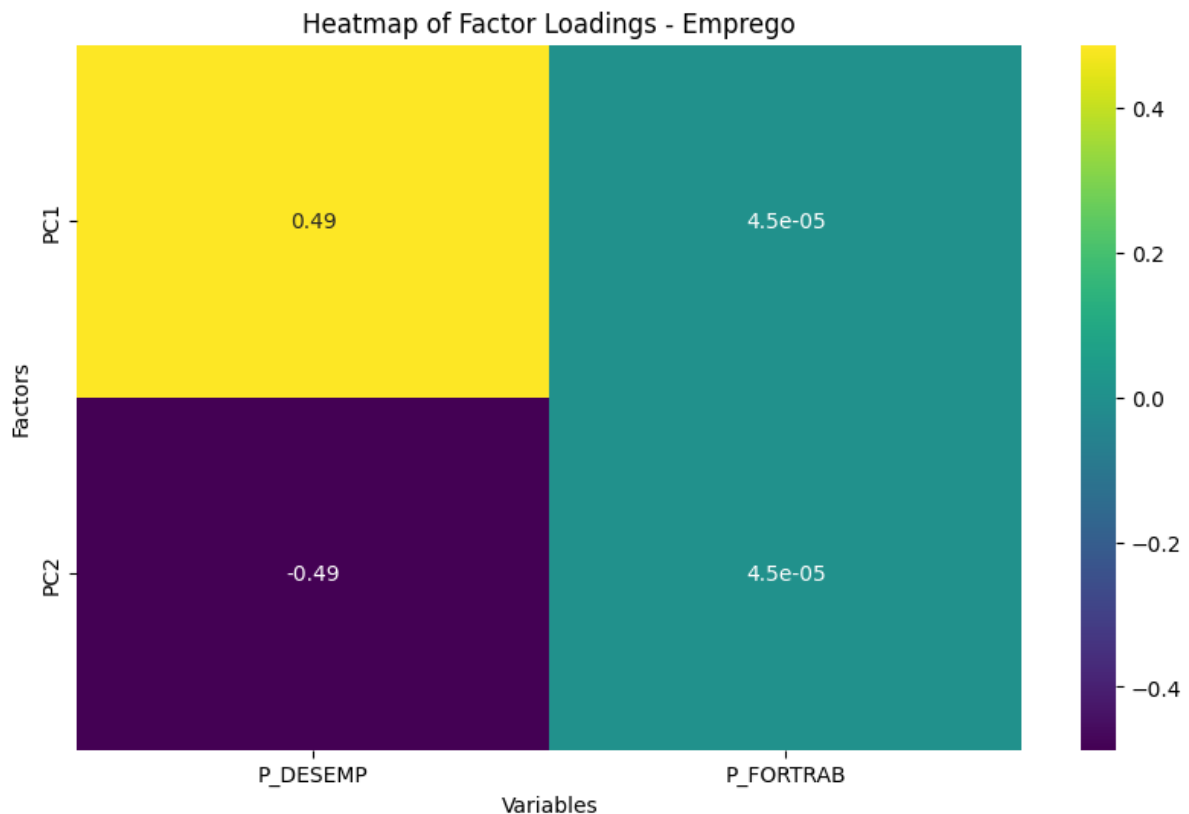
Variância Cumulativa Explicada



A variância cumulativa explicada nos dá uma ideia de quanto da informação total dos dados é capturada pelo(s) fator(es).

- **Primeiro Fator:** Explica 23.73% da variância total, o que pode parecer baixo, mas é típico em ciências sociais, onde as variáveis podem ser influenciadas por muitos fatores complexos e inter-relacionados.

Cargas Fatoriais para Emprego



Os loadings representam a correlação entre as variáveis originais e os fatores estimados.

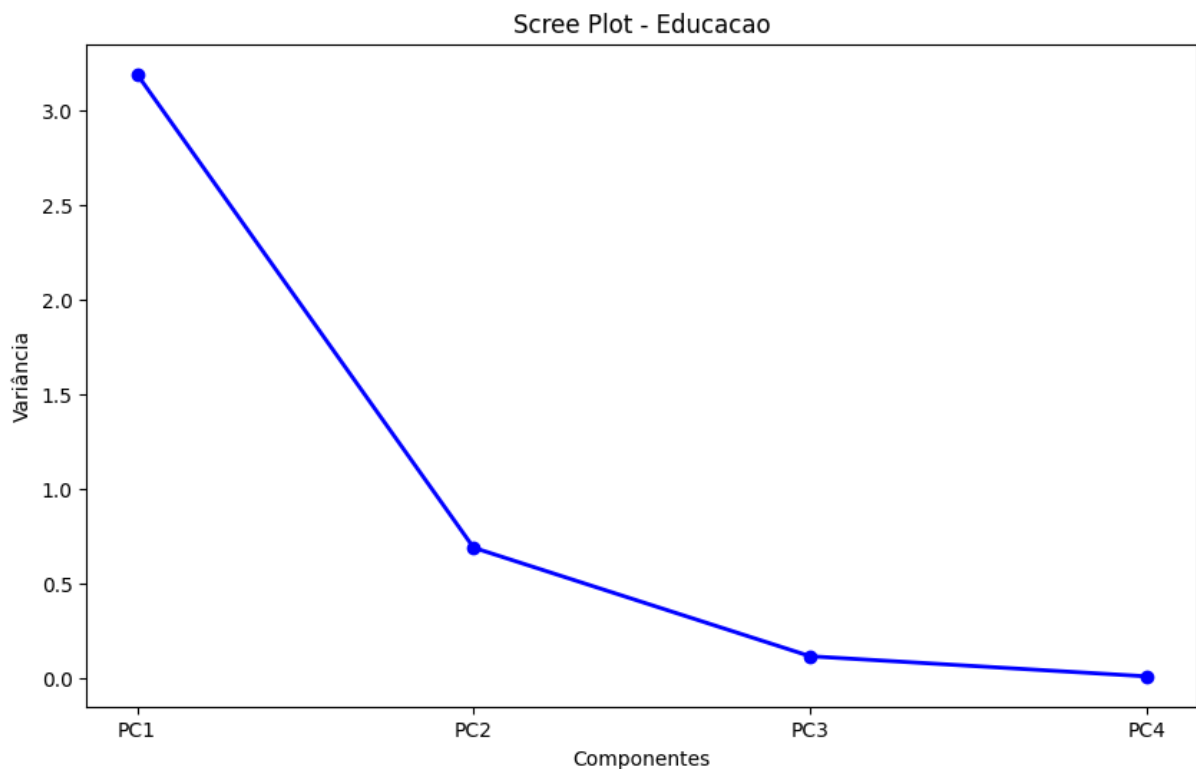
- **Primeiro Fator:** Tem loadings iguais e opostos para **P_DESEMP** (taxa de desemprego) e **P_FORTTAB** (participação na força de trabalho), ambos com aproximadamente 0.487. Este fator sugere uma relação direta entre as duas variáveis: onde a taxa de desemprego é alta, a participação na força de trabalho é baixa, e vice-versa.
- **Segundo Fator:** Os loadings são essencialmente zero, indicando que este fator não está capturando variações adicionais significativas além daquelas já explicadas pelo primeiro fator.

Interpretação

O resultado da EFA para a dimensão de emprego sugere que um único fator é suficiente para descrever as variações nas variáveis de emprego. Este fator parece refletir um equilíbrio entre a taxa de desemprego e a participação na força de trabalho. Em termos práticos, isso pode significar que as intervenções políticas que visam reduzir a taxa de desemprego devem também considerar como aumentar a participação na força de trabalho, e que essas duas métricas devem ser consideradas juntas ao avaliar a saúde do mercado de trabalho em uma área específica.

Dimensão Educação

Variância Cumulativa Explicada

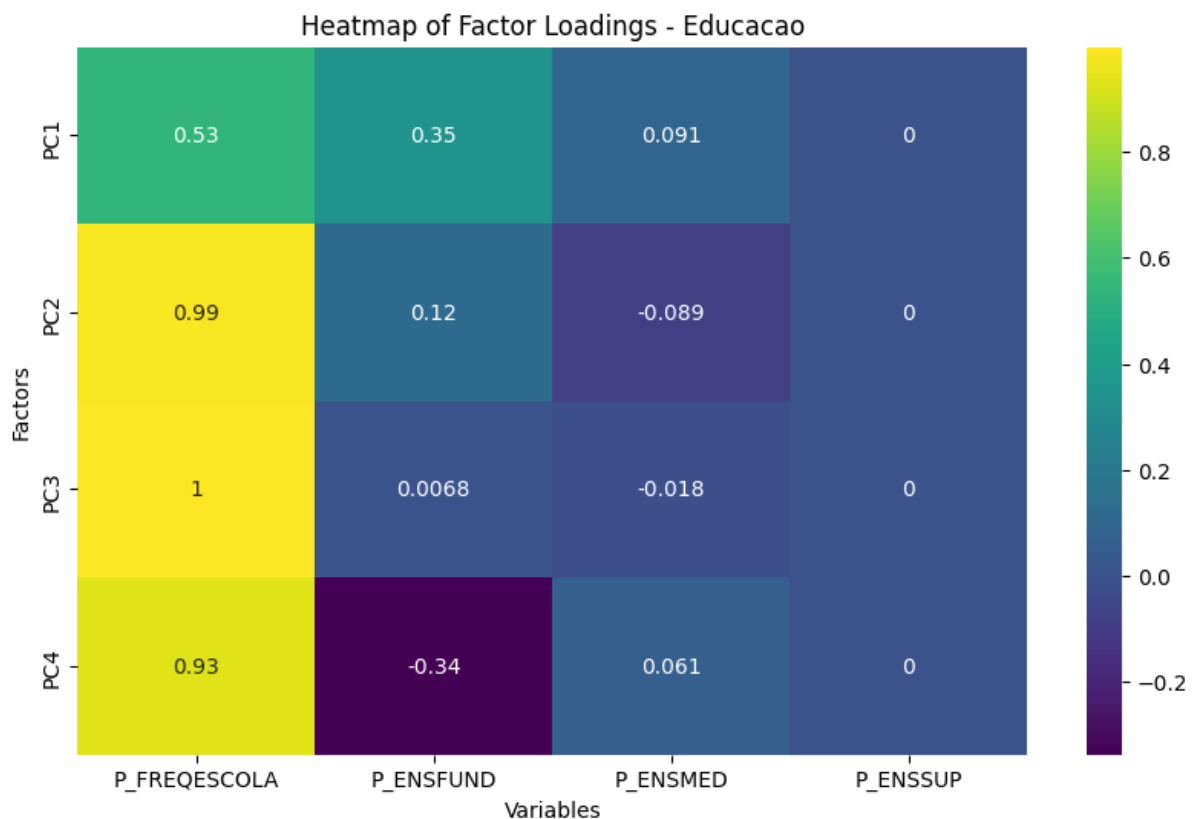


A variância cumulativa explicada nos ajuda a entender quanto da informação total dos dados é capturada pelos fatores.

- Primeiro Fator: Explica aproximadamente 78.08% da variância total, o que é significativo e sugere que este fator é o mais importante na descrição da variação na educação.

- Segundo Fator: Adiciona um pouco mais à variância explicada, elevando o total para 84.32%.
- Terceiro Fator: Acrescenta apenas uma pequena quantidade à variância explicada, atingindo um total de 84.83%.

Cargas Fatoriais para Educação



Os loadings representam a correlação entre as variáveis originais e os fatores estimados, com valores mais altos indicando uma relação mais forte.

- Primeiro Fator: Apresenta loadings altos para todas as variáveis, especialmente P_ENSFUND, P_ENSMED, e P_ENSSUP. Este fator parece capturar o nível geral de educação formal alcançado, com uma pontuação alta representando um alto nível de realização educacional.
- Segundo Fator: Mostra um loading moderado e positivo para P_FREQESCOLA e um loading negativo para P_ENSSUP. Este fator pode estar destacando uma dimensão que

contrasta a frequência escolar de jovens com a realização de níveis mais elevados de educação, como o ensino superior.

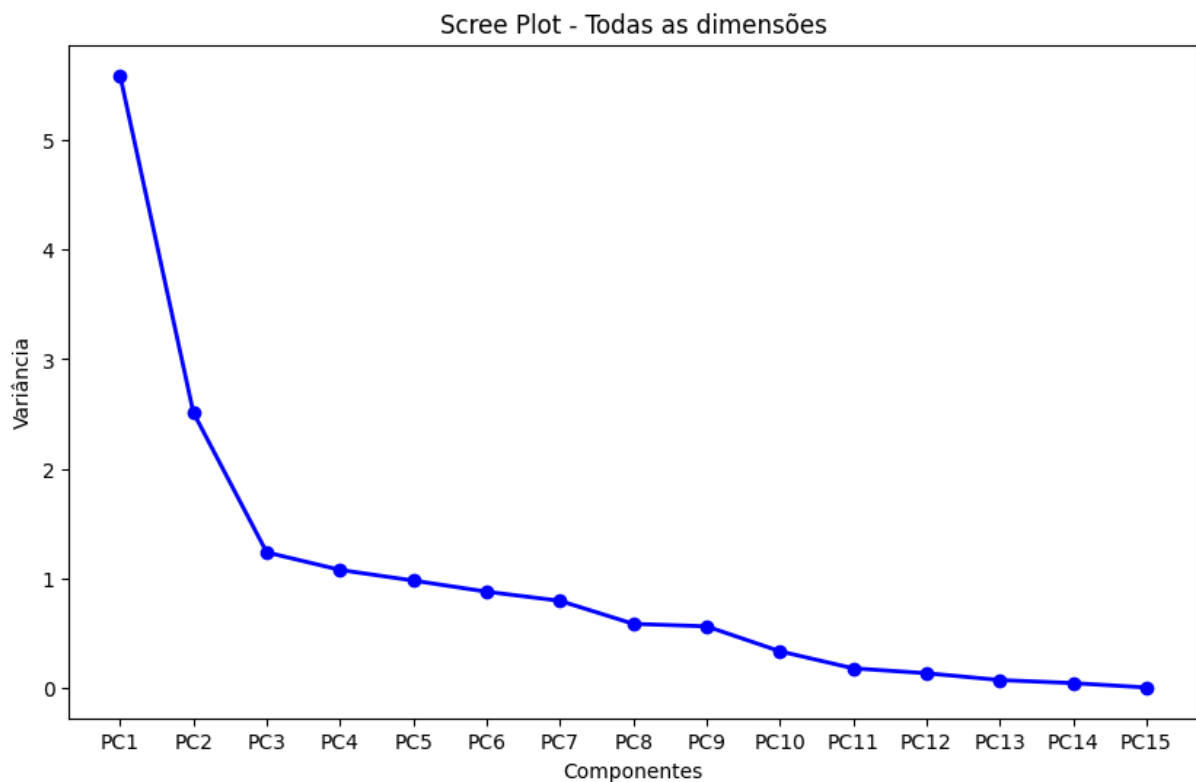
- Terceiro Fator: Tem loadings pequenos e mistos para as variáveis, indicando que ele pode estar capturando diferenças mais sutis ou específicas nos padrões de educação que não são capturadas pelos dois primeiros fatores.

Interpretação

O modelo EFA para a dimensão de educação indica que a maior parte das variações nos dados de educação pode ser explicada por um único fator, que parece refletir o nível geral de educação alcançado pelas pessoas.

Todas as dimensões

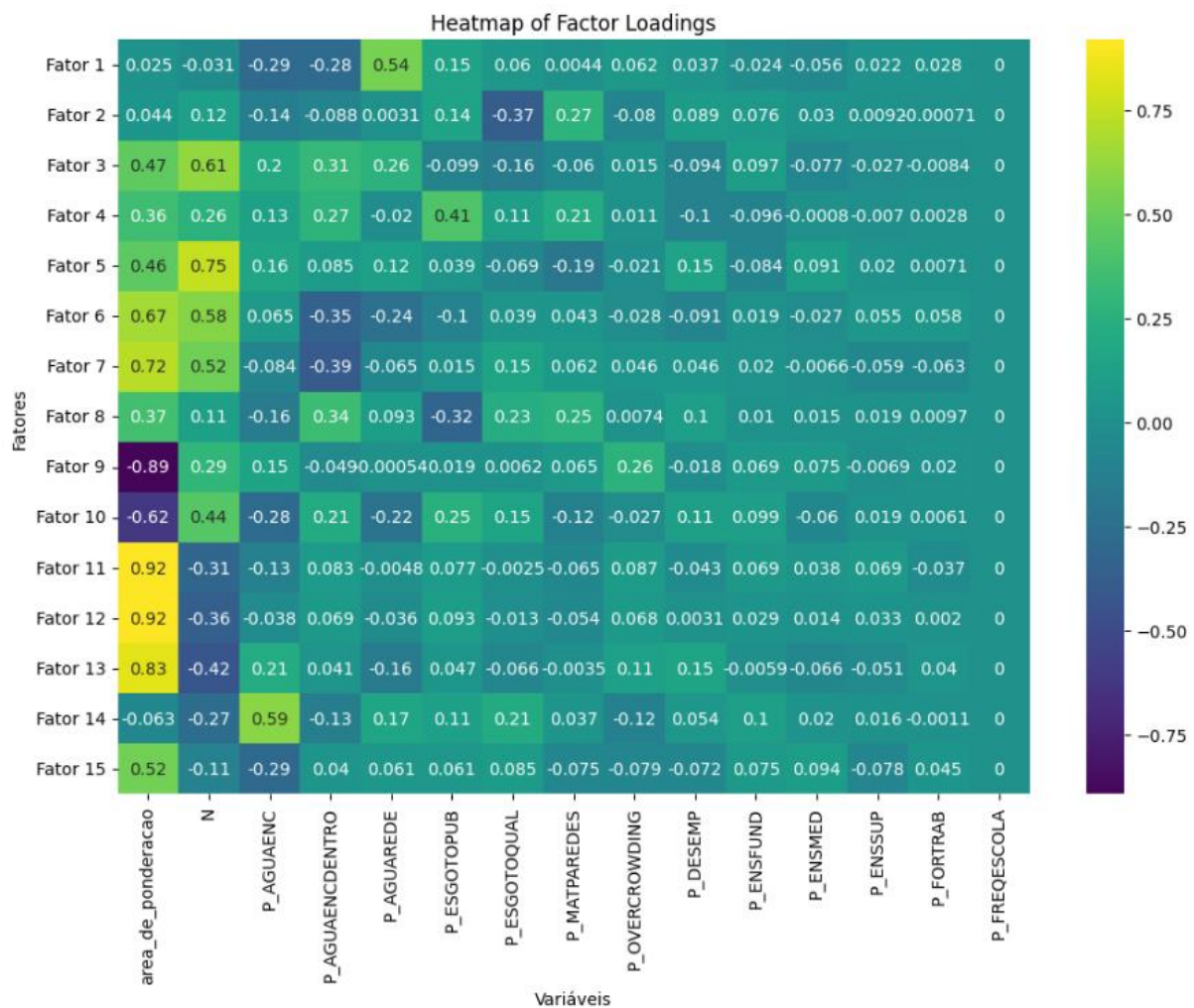
Variância Explicada



A variância explicada por cada fator nos dá uma ideia de quanto da variabilidade total das variáveis originais é capturada por cada fator.

- Primeiro Fator: É o mais significativo, capturando a maior parte da variância. Isso é consistente com a ideia de que a infraestrutura de saneamento é um componente crítico da qualidade de vida.
- Segundo ao Quinto Fator: Apesar de capturarem menos variância que o primeiro, ainda são importantes para compreender outros aspectos socioeconômicos.
- Fatores Restantes: Contribuem incrementalmente menos para a variância explicada, sugerindo que podem representar aspectos mais específicos ou menos importantes das variáveis.

Cargas Fatoriais para Dataset Completo



As cargas fatoriais (ou loadings) representam a contribuição de cada variável para os fatores latentes encontrados na análise. Elas são usadas para interpretar o significado dos fatores.

Fatores Principais:

- Fator 1: Este fator tem cargas positivas fortes em variáveis como P_AGUAENC, P_AGUAENC DENTRO, P_AGUAREDE (saneamento), e P_ESGOTOPUB, P_ESGOTOQUAL (esgoto), indicando que este fator pode estar associado à qualidade da infraestrutura básica de saneamento.
- Fator 2: As cargas mais altas neste fator são para P_MATPAREDES (qualidade da moradia) e P_OVERCROWDING (superlotação), sugerindo que este fator reflete aspectos da habitação e densidade populacional dentro das casas.
- Fator 3: Há cargas significativas para P_ENSFUND, P_ENSMED, e P_ENSSUP (níveis de educação), o que sugere que este fator está relacionado ao nível educacional da população.
- Fator 4 e 5: Mostram cargas moderadas para P_DESEMP (desemprego) e P_FORTRAB (participação na força de trabalho), podendo estar relacionados à dinâmica do mercado de trabalho.
- Fatores Subsequentes: Apresentam menor magnitude nas cargas fatoriais, indicando que podem estar capturando nuances mais específicas ou informações redundantes que não são explicadas pelos primeiros fatores.

Interpretação

A EFA para o dataset completo sugere que a qualidade da infraestrutura de saneamento e o nível educacional são os aspectos mais dominantes no que diz respeito às variações socioeconômicas. Os primeiros fatores capturam as variações mais significativas e poderiam ser interpretados como indicadores gerais da qualidade de vida nas áreas representadas pelo dataset. Os fatores subsequentes podem refletir aspectos mais específicos e menos generalizáveis.

Em resumo, essa análise pode ajudar a compreender melhor as complexidades socioeconômicas e a identificar áreas focais para a implementação de políticas públicas e intervenções sociais.

4. Exemplifique como as diferentes regras de seleção podem afetar a escolha do número de componentes e/ou fatores em cada dimensão. Comente sobre como o número de componentes /fatores escolhido é afetado pela correlação entre as variáveis originais. Opcional: faça também uma conexão com a regra de aceitação/rejeição de H_0 no teste da qualidade de ajuste do(s) modelo(s) EFA ajustado.

As regras de seleção de componentes e fatores em uma análise fatorial ou PCA são críticas porque determinam quantos componentes ou fatores são retidos para descrever os dados. Estas regras influenciam diretamente a interpretação dos resultados e as conclusões que podem ser tiradas do modelo. Abaixo as mais comuns e como elas podem afetar a escolha do número de componentes/fatores:

1. **Regra de Kaiser:** Mantém fatores com autovalores maiores que 1. Este método pode levar a reter mais fatores do que o necessário, especialmente em casos onde há muitas variáveis com pequenas correlações.
2. **Scree Plot:** Analisa o gráfico dos autovalores e procura um ponto onde a curva se achata, indicando que os fatores subsequentes não adicionam muito à explicação da variância. Esta regra é mais visual e subjetiva, podendo variar conforme a interpretação do pesquisador.
3. **Análise paralela:** Compara os autovalores obtidos com os autovalores de dados aleatórios. Fatores são retidos se seus autovalores forem maiores do que os equivalentes de uma matriz aleatória. Esta abordagem é mais robusta e pode evitar a retenção de fatores demais.
4. **Teste de significância de Bartlett:** Verifica se a matriz de correlação é uma matriz identidade, o que indicaria que as variáveis são todas independentes e não há estrutura fatorial. Se a hipótese nula (H_0) for rejeitada, indica-se que os fatores são adequados.

5. **Teste de esfericidade de Bartlett:** Avalia se a matriz de correlação difere da matriz identidade, sugerindo que as correlações entre as variáveis são suficientemente fortes para justificar uma EFA.
6. **Regra do dedo mínimo:** Escolhe fatores que têm pelo menos três ou mais loadings significativos para evitar fatores com uma única variável.

Correlação entre as variáveis originais

A escolha do número de componentes/fatores é afetada pela correlação entre as variáveis originais. Se as variáveis são altamente correlacionadas, menos fatores são necessários para explicar a variância nos dados. No entanto, se as correlações são baixas ou as variáveis são ortogonais (independentes), mais fatores podem ser necessários para capturar a estrutura dos dados.

Resumo:

- **Alta Correlação:** Variáveis altamente correlacionadas tendem a reduzir o número de fatores necessários, pois compartilham uma quantidade significativa de variância.
- **Baixa Correlação:** Pouca correlação entre as variáveis pode levar à necessidade de mais fatores para explicar a variância dos dados.

Conexão com a Regra de Aceitação/Rejeição de H_0 no Teste da Qualidade de Ajuste

Em EFA, um teste comum da qualidade de ajuste é o Teste de Esfericidade de Bartlett, que testa se a matriz de correlação é igual à matriz identidade (indicando que não há correlações entre as variáveis). A rejeição de H_0 sugere que há correlações entre as variáveis e que uma EFA é apropriada. Se H_0 não for rejeitada, isso pode indicar que uma redução de dimensionalidade não é apropriada para os dados.

Conclusão

A escolha do número de componentes/fatores é uma etapa crítica na análise fatorial e é influenciada pela correlação entre as variáveis. Diferentes regras de seleção podem levar a resultados distintos, e a correlação entre variáveis afeta tanto essa escolha quanto o resultado

dos testes de qualidade de ajuste do modelo EFA. A rejeição da hipótese nula no teste de ajuste apoia a realização da EFA e a escolha de um número adequado de fatores. Portanto, a seleção do número de fatores é um processo inter-relacionado que envolve tanto a análise das características dos dados quanto os resultados dos testes estatísticos de ajuste.

5. Construa um “índice de status socioeconômico” utilizando a PCA e/ou a EFA (utilizando pelo menos um indicador de cada dimensão). Interprete os resultados em termos da dimensão correspondente e do padrão de “status socioeconômico” para a cidade escolhida.

Seguimos as seguintes etapas **para construir o Índice:**

1. Seleção de Indicadores:

- Variáveis escolhidas:
 - Saneamento: Proporção de domicílios com acesso à rede pública de esgoto (**P_ESGOTOPUB**).
 - Moradia: Proporção de domicílios com paredes externas de materiais duráveis (**P_MATPAREDES**).
 - Emprego: Taxa de desemprego (**P_DESEMP**).
 - Educação: Proporção da população com Ensino Superior completo (**P_ENSSUP**).

2. Normalização dos Dados:

- Normalizar os dados para garantir que todas as variáveis tenham o mesmo peso na análise.

3. Aplicação da PCA/EFA para Índice:

- Técnica escolhida: PCA.

4. Interpretação dos Resultados

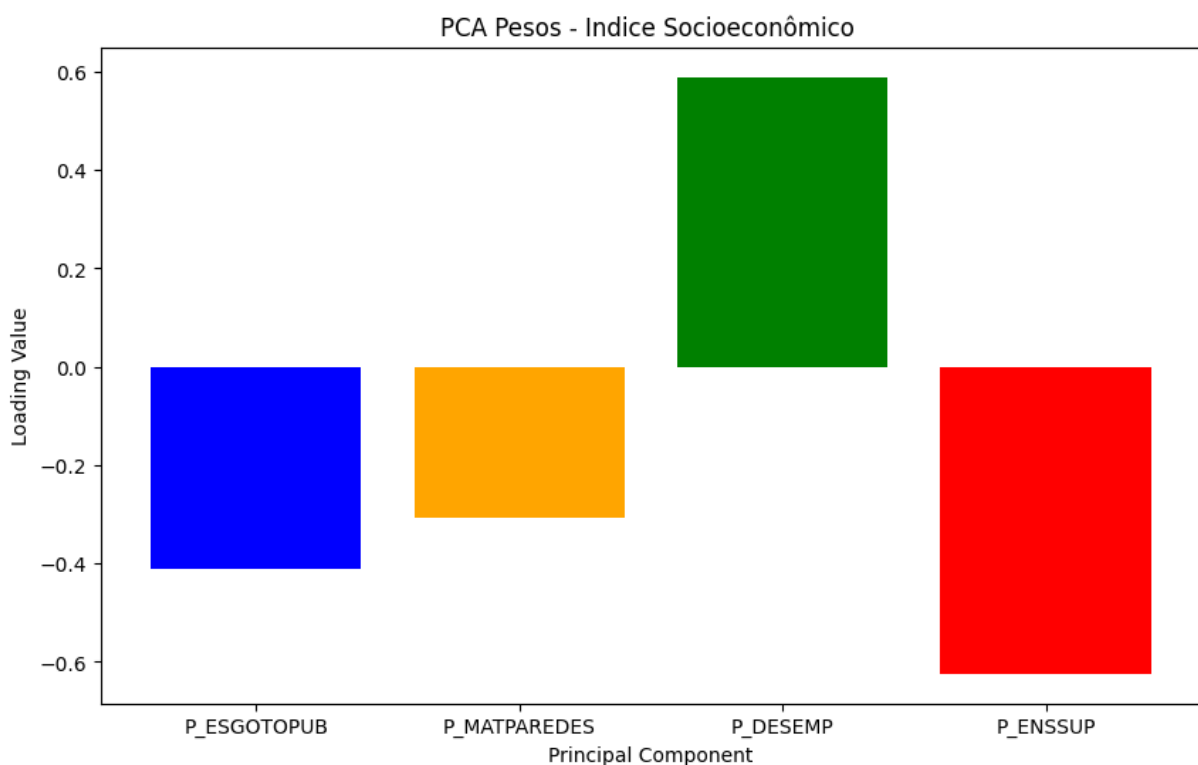
Dimensões e suas contribuições para o indicador:

- **Saneamento e Moradia:** Contribuições positivas destas dimensões podem indicar melhor infraestrutura e condições de vida.

- **Emprego:** Uma contribuição negativa da taxa de desemprego no índice sugere que uma menor taxa de desemprego está associada a um melhor status socioeconômico.
- **Educação:** Altos níveis de educação contribuindo positivamente para o índice refletem a importância da educação na determinação do status socioeconômico.

Para construir um “índice de status socioeconômico” utilizando a Análise de Componentes Principais (PCA) e interpretar os resultados, selecionamos um indicador de cada dimensão importante para o Rio de Janeiro: Saneamento (**P_ESGOTOPUB**), Moradia (**P_MATPAREDES**), Emprego (**P_DESEMP**) e Educação (**P_ENSSUP**). Vamos interpretar o resultado da PCA com base nas cargas fatoriais, na variância explicada e nos primeiros valores do índice.

Pesos PCA (Loadings)



- **Saneamento (P_ESGOTOPUB):** A carga negativa (-0.4103) sugere que uma melhor infraestrutura de saneamento está associada a um status socioeconômico mais alto.

- **Moradia (P_MATPAREDES):** A carga negativa (-0.3083) indica que a presença de materiais duráveis nas construções está relacionada a um status socioeconômico mais elevado.
- **Emprego (P_DESEMP):** A carga positiva (0.5892) mostra que uma maior taxa de desemprego está inversamente relacionada ao status socioeconômico.
- **Educação (P_ENSSUP):** A carga negativa (-0.6240) implica que altos níveis de educação superior estão positivamente relacionados ao status socioeconômico.

Variância Explicada

A variância explicada (49.52%) indica que quase metade da variabilidade total nos indicadores selecionados é capturada pelo índice. Isso sugere que o índice é um bom resumo das dimensões socioeconômicas consideradas.

Índice de Status Socioeconômico

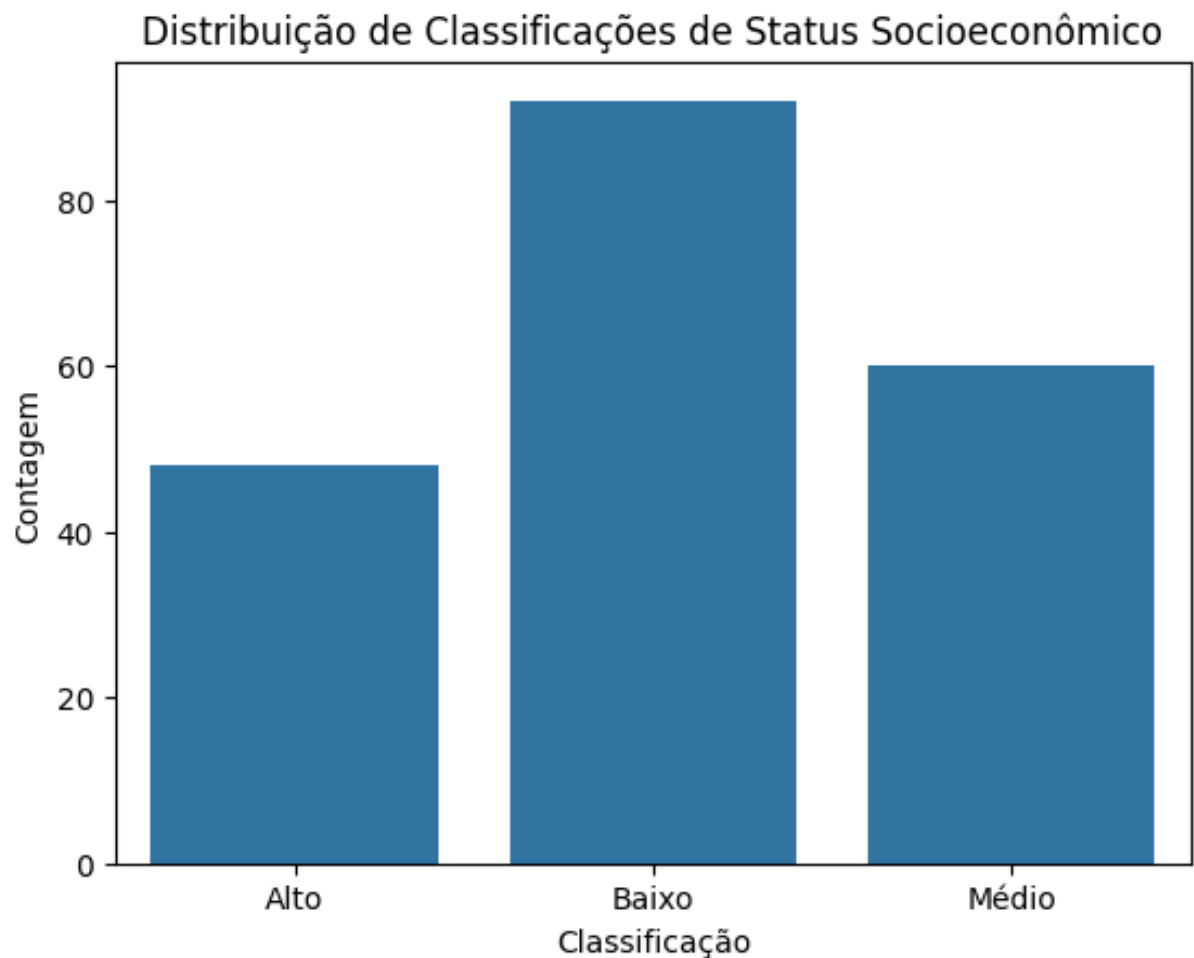
O índice reflete uma combinação de fatores como qualidade do saneamento, condições de moradia, taxas de desemprego e níveis de educação. Áreas com melhor infraestrutura de saneamento, moradias de qualidade, baixas taxas de desemprego e altos níveis de educação tendem a ter um índice mais alto, indicando um status socioeconômico mais favorável.

Resultado descritivo do índice:

índice_socioeconomico	
count	200.000
mean	-0.000
std	1.411
min	-3.046
25%	-0.858
50%	0.127
75%	0.923
max	3.370

Os valores do índice (por exemplo, 1.4171, -2.1206, etc.) representam o status socioeconômico de diferentes áreas ou segmentos do Rio de Janeiro. Valores mais altos indicam um status socioeconômico mais elevado, enquanto valores mais baixos indicam um status inferior.

Tomamos a decisão arbitrária de classificar resultados abaixo de 0 como baixo, de 0 a 1 médio e acima de 1 alto. Tendo em vista as 200 regiões do Rio de Janeiro que tiveram um índice e foram classificadas, vemos que grande parte das regiões da cidade do Rio possuem um índice socioeconômico abaixo de 0.



Este índice poderia ser utilizado por formuladores de políticas e pesquisadores para identificar áreas que requerem atenção especial, desenvolver programas direcionados de melhoria e avaliar o impacto de políticas socioeconômicas.