

Data Science and Big Data

Taconeli, C.A.

23 de novembro, 2018

Análise de componentes principais

Análise de componentes principais

- A análise de componentes principais (ACP) é uma técnica multivariada, aplicada a um conjunto de p variáveis aleatórias, que tem como principais objetivos:
 - Reduzir a dimensão dos dados, projetando-os em uma dimensão $k < p$;
 - Explorar a estrutura de variância das variáveis;
 - Determinar índices e produzir escores com base nos resultados avaliados para as p variáveis.
- A ACP consiste na obtenção de novas variáveis, determinadas a partir das variáveis originais, de tal forma que um pequeno número de novas variáveis (componentes principais) seja capaz de explicar a maior parte possível da variação presente nos dados.

Análise de componentes principais

- As novas variáveis obtidas (denominadas componentes) são tais que o primeiro (principal) componente é aquele capaz de explicar a maior parte possível da variação dos dados, o segundo é aquele que explica a maior parte possível não explicada pelo primeiro...
- Algebricamente, os componentes (Z_1, Z_2, \dots, Z_p) correspondem a combinações lineares das variáveis originais Y_1, Y_2, \dots, Y_p .
- Geometricamente, as combinações lineares representam um novo sistema de coordenadas, obtido por meio da rotação do sistema original.

Análise de componentes principais

- Motivação de uma ACP, em que y_1 e y_2 são as variáveis originais e z_1 e z_2 os componentes

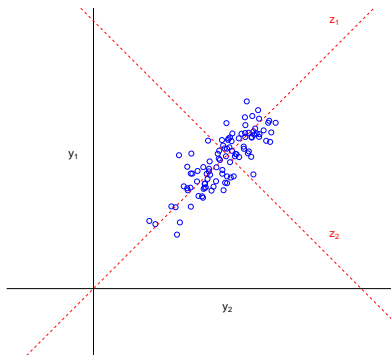


Figura 1: Ilustração - componentes principais.

Análise de componentes principais

- Em algumas aplicações, os componentes da ACP configuram o objetivo final do estudo.
- Em outras situações, os componentes obtidos servem como passo intermediário para realização de outras análises, como regressão, classificação, agrupamento. . .

Análise de componentes principais

- Considere o vetor aleatório $\mathbf{Y}' = (Y_1, Y_2, \dots, Y_p)$ com matriz de covariâncias $\mathbf{\Sigma}$ de autovalores $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.
- Considere as combinações lineares:

$$Z_1 = \mathbf{a}'_1 \mathbf{Y} = a_{11} Y_1 + a_{12} Y_2 + \dots + a_{1p} Y_p$$

$$Z_2 = \mathbf{a}'_2 \mathbf{Y} = a_{21} Y_1 + a_{22} Y_2 + \dots + a_{2p} Y_p$$

$$\vdots$$

$$Z_p = \mathbf{a}'_p \mathbf{Y} = a_{p1} Y_1 + a_{p2} Y_2 + \dots + a_{pp} Y_p$$

Análise de componentes principais

- Sabemos que:

$$\text{Var}(Z_i) = \mathbf{a}_i' \boldsymbol{\Sigma} \mathbf{a}_i, \quad i = 1, 2, \dots, p$$

$$\text{Cov}(Z_i, Z_k) = \mathbf{a}_i' \boldsymbol{\Sigma} \mathbf{a}_k, \quad i, k = 1, 2, \dots, p$$

- Os componentes principais são as combinações lineares de Y_1, Y_2, \dots, Y_p **não correlacionadas** e que tenham **maior variância possível**.
- No entanto, $\text{Var}(Z_i) = \mathbf{a}_i' \boldsymbol{\Sigma} \mathbf{a}_i$ pode crescer indefinidamente multiplicando \mathbf{a}_i por alguma constante. Assim, aplicamos a restrição $\mathbf{a}_i' \mathbf{a}_i = 1$.

Análise de componentes principais

- Resumo da análise de componentes principais:

Passo 1 - Primeiro componente principal = combinação linear $\mathbf{a}'_1 \mathbf{Y}$ que maximiza $Var(\mathbf{a}'_1 \mathbf{Y})$ sujeita a $\mathbf{a}'_1 \mathbf{a}_1 = 1$;

Passo 2 - Segundo componente principal = combinação linear $\mathbf{a}'_2 \mathbf{Y}$ que maximiza $Var(\mathbf{a}'_2 \mathbf{Y})$ sujeita a $\mathbf{a}'_2 \mathbf{a}_2 = 1$ e $Cov(\mathbf{a}'_1 \mathbf{Y}, \mathbf{a}'_2 \mathbf{Y}) = 0$;

⋮

Passo i - i-ésimo componente principal = combinação linear $\mathbf{a}'_i \mathbf{Y}$ que maximiza $Var(\mathbf{a}'_i \mathbf{Y})$ sujeita a $\mathbf{a}'_i \mathbf{a}_i = 1$ e $Cov(\mathbf{a}'_i \mathbf{Y}, \mathbf{a}'_k \mathbf{Y}) = 0$, para todo $k < i$.

Análise de componentes principais

- Para determinação dos componentes principais, com base no que foi exposto, usaremos o seguinte teorema:
- **Maximização de formas quadráticas** Seja B uma matriz positiva definida com autovalores $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ e autovetores associados normalizados $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$. Então:

$$\max_{\mathbf{x} \neq 0} \frac{\mathbf{x}' B \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \lambda_1, \text{ obtido quando } \mathbf{x} = \mathbf{e}_1;$$

$$\min_{\mathbf{x} \neq 0} \frac{\mathbf{x}' B \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \lambda_p, \text{ obtido quando } \mathbf{x} = \mathbf{e}_p.$$

Adicionalmente,

$$\max_{\mathbf{x} \perp \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_k} \frac{\mathbf{x}' B \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \lambda_{k+1}, \text{ obtido quando } \mathbf{x} = \mathbf{e}_{k+1}.$$

Análise de componentes principais

- Assim, no contexto de componentes principais, seja $\mathbf{Y}' = (Y_1, Y_2, \dots, Y_p)$ um vetor aleatório. Seja $\mathbf{\Sigma}$ a matriz de variâncias e covariâncias e $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$ seus autovalores e autovetores, tal que $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Então:

$$\max_{\mathbf{a} \neq \mathbf{0}} \frac{(\mathbf{a}'\mathbf{\Sigma}\mathbf{a})}{\mathbf{a}'\mathbf{a}} = \max_{\mathbf{a} \neq \mathbf{0}} (\mathbf{a}'\mathbf{\Sigma}\mathbf{a}) = \lambda_1, \text{ obtido quando } \mathbf{a} = \mathbf{e}_1;$$

$$\min_{\mathbf{a} \neq \mathbf{0}} \frac{(\mathbf{a}'\mathbf{\Sigma}\mathbf{a})}{\mathbf{a}'\mathbf{a}} = \min_{\mathbf{a} \neq \mathbf{0}} (\mathbf{a}'\mathbf{\Sigma}\mathbf{a}) = \lambda_p, \text{ obtido quando } \mathbf{a} = \mathbf{e}_p.$$

Adicionalmente,

$$\max_{\mathbf{a} \perp \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_k} \frac{(\mathbf{a}'\mathbf{\Sigma}\mathbf{a})}{\mathbf{a}'\mathbf{a}} = \max_{\mathbf{a} \perp \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_k} (\mathbf{a}'\mathbf{\Sigma}\mathbf{a}) = \lambda_{k+1}, \text{ obtido quando } \mathbf{a} = \mathbf{e}_{k+1}$$

Análise de componentes principais

- O i -ésimo componente principal é definido por:

$$Z_i = \mathbf{e}_i' \mathbf{Y} = e_{i1} Y_1 + e_{i2} Y_2 + \dots + e_{ip} Y_p, \quad i = 1, 2, \dots, p.$$

- Como consequências:

$$\text{Var}(Z_i) = \mathbf{e}_i' \mathbf{\Sigma} \mathbf{e}_i = \lambda_i, \quad i = 1, 2, \dots, p.$$

$$\text{Cov}(Z_i, Z_k) = \mathbf{e}_i' \mathbf{\Sigma} \mathbf{e}_k = 0, \quad i \neq k.$$

Análise de componentes principais

- Adicionalmente, como resultado da decomposição espectral:

$$\sum_{i=1}^p \text{Var}(Y_i) = \sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \lambda_1 + \lambda_2 + \lambda_p = \sum_{i=1}^p \text{Var}(Z_i)$$

- Assim, a proporção da variação total dos dados explicada pelo i -ésimo componente pode ser calculada por:

$$\frac{\lambda_i}{\lambda_1 + \lambda_2 + \dots + \lambda_p}, \quad i = 1, 2, \dots, p.$$

Análise de componentes principais

- A correlação entre o i -ésimo componente (Z_i) e a k -ésima variável (Y_k) é dada por:

$$\rho_{Z_i, Y_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sigma_{kk}}.$$

Análise de componentes principais

- Em geral, observamos as proporções explicadas pelos componentes principais e retemos um pequeno número de componentes capazes de explicar boa parte da variabilidade dos dados;
- Os coeficientes (**cargas**) dos componentes (elementos dos autovetores), seus sinais e magnitudes, permitem interpretar os componentes e avaliar a importância das variáveis em sua constituição;
- As correlações entre componentes e variáveis podem ser usadas também pra avaliar importância das variáveis na constituição dos componentes e produção de gráficos;

Análise de componentes principais

- Chamamos **escores** os valores dos componentes (resultados das combinações lineares) calculados para um particular indivíduo;
- Os escores podem ser usados para efeito de ranqueamento dos indivíduos, agrupamento, classificação. . .

Análise de componentes principais

- Nas situações em que as variáveis originais têm escalas diferentes, é natural esperarmos que as variâncias, conseqüentemente, sejam diferentes.
- Para evitar distorções nos resultados, ocasionadas pela predominância de variáveis com maior variação, é recomendável considerar a análise de componentes principais com base nas variáveis padronizadas:

$$Y_1^* = \frac{Y_1 - \mu_1}{\sqrt{\sigma_{11}}}, \quad Y_2^* = \frac{Y_2 - \mu_2}{\sqrt{\sigma_{22}}}, \dots, \quad Y_p^* = \frac{Y_p - \mu_p}{\sqrt{\sigma_{pp}}}$$

Análise de componentes principais

- Os componentes principais para $Y_1^*, Y_2^*, \dots, Y_p^*$ são determinados pela matriz de variâncias e covariâncias das variáveis padronizadas, o que equivale à matriz de correlações de Y_1, Y_2, \dots, Y_p .
- Assim, teremos $Z_i^* = \mathbf{e}_i' \mathbf{Y}_i^*, i = 1, 2, \dots, p$, com:

$$\sum \text{Var} Z_i^* = \sum \text{Var} Y_i^* = p$$

e

$$\rho_{Y_i^*, Z_k^*} = e_{ik} \sqrt{\lambda_i}, \quad i, k = 1, 2, \dots, p,$$

sendo $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$ os autovalores e autovetores da matriz de correlações.

Análise de componentes principais

- Na prática, aplicaremos a ACP a uma amostra, de tal forma que a análise é conduzida com base em uma amostra $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$, com estatísticas amostrais $\bar{\mathbf{y}}$, \mathbf{S} e \mathbf{R} .
- Seja \mathbf{S} a matriz de covariâncias amostral, com $(\hat{\lambda}_1, \hat{\mathbf{e}}_1), (\hat{\lambda}_2, \hat{\mathbf{e}}_2), \dots, (\hat{\lambda}_p, \hat{\mathbf{e}}_p)$ os pares de autovalores e autovetores, tal que $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$.
- O i -ésimo componente principal é dado por:

$$\hat{z}_i = \hat{\mathbf{e}}_i' \mathbf{y} = \hat{e}_{i1}y_1 + \hat{e}_{i2}y_2 + \dots + \hat{e}_{ip}y_p, \quad i = 1, 2, \dots, p$$

Análise de componentes principais

- A variância amostral de \hat{z}_i é igual a $\hat{\lambda}_i$, e a covariância amostral entre \hat{z}_i e \hat{z}_j é igual a zero;
- A variância amostral total fica dada por:

$$\sum_{i=1}^n s_{ii} = \hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p$$

- A correlação entre o i-ésimo componente e a k-ésima variável é dada por:

$$r_{\hat{z}_i, y_k} = \frac{\hat{e}_{ik} \sqrt{\hat{\lambda}_i}}{\sqrt{s_{kk}}}, i, k = 1, 2, \dots, p.$$

Análise de componentes principais

- Um resultado fundamental da ACP é a produção de escores, que correspondem a valores calculados para os indivíduos nos componentes obtidos.
- Mesmo que a ACP seja realizada a partir das variáveis em sua escala original, é usual calcular os escores com base nas variáveis centradas nas respectivas médias. Por exemplo, para um indivíduo g , no i -ésimo componente:

$$\hat{z}_{gi} = \hat{\mathbf{e}}_i'(\mathbf{y} - \bar{\mathbf{y}}) = e_{i1}(y_{g1} - \bar{y}_1) + e_{i2}(y_{g2} - \bar{y}_2) + \dots + e_{ip}(y_{gp} - \bar{y}_p).$$

- Como resultado, os componentes principais obtidos terão suas médias amostrais iguais a 0.

Análise de componentes principais

- Já no caso em que a análise é conduzida com base nas variáveis padronizadas (matriz de correlações), os escores são calculados a partir delas:

$$\hat{z}_{gi} = e_{i1} \frac{y_{g1} - \bar{y}_1}{\sqrt{s_{11}}} + e_{i2} \frac{y_{g2} - \bar{y}_2}{\sqrt{s_{22}}} + \dots + e_{ip} \frac{y_{gp} - \bar{y}_p}{\sqrt{s_{pp}}}.$$

Análise de componentes principais

- Um dos principais atrativos da ACP é a possibilidade de visualização dos resultados em gráficos de duas ou três dimensões, como resultado da projeção. Algumas alternativas:
 - Plotar os escores de cada indivíduo em cada componente, a fim de identificar grupos de indivíduos com escores (e características) semelhantes, identificar indivíduos extremos, . . .
 - Plotar as correlações entre variáveis e componentes, a fim de visualizar as associações entre variáveis e componentes, as intercorrelações entre variáveis. . .
 - Plotar conjuntamente dados e variáveis num mesmo gráfico, permitindo visualizar globalmente relações entre variáveis e indivíduos (gráfico biplot, veremos adiante).

Análise de componentes principais

- Uma das questões mais recorrentes em ACP é o número de componentes a reter e explorar na análise.
- Um dos aspectos a se levar em consideração é a interpretação prática dos componentes, possível justificativa à luz da teoria da área. . .
- Obviamente, no entanto, aspectos referentes ao desempenho da técnica devem ser considerados.

Análise de componentes principais

- Alguns critérios que fundamentam a definição do número de componentes a serem “retidos”:
 - Reter o menor número k de componentes principais capaz de explicar uma porcentagem mínima desejável da variação dos dados (Ex: 70 ou 80%);
 - Reter os componentes cujos autovalores são maiores que a média dos autovalores, $(\sum_{i=1}^p \lambda_i / p)$. Para a matriz de correlações a média é igual a 1;
 - Decidir o número de componentes com base na apreciação do *scree plot*;
 - Testar a significância dos componentes principais.

Análise de componentes principais

- Um teste a ser considerado em ACP é o teste de completa independência entre as variáveis.
- Caso as variáveis sejam conjuntamente independentes, não há porque fazer uma ACP, uma vez que, nesse caso, cada variável individualmente dominará um componente.
- Nesse caso, consideramos a hipótese nula $H_0 : \mathbf{P} = \mathbf{I}$, sendo \mathbf{P} a matriz de correlações populacional. A estatística do teste é dada por:

$$u = - \left(n - 1 - \frac{2p + 5}{6} \right) * \ln |\mathbf{R}|,$$

que segue, sob H_0 , distribuição χ_f^2 , com $f = \frac{1}{2}p(p - 1)$. Rejeitaremos H_0 , a um nível α , se $u > \chi_f^2(\alpha)$ (Teste de Bartlett).