



SME0822 Análise Multivariada e Aprendizado Não-Supervisionado

Aula 7a: **Análise de Componentes Principais**

Prof. Cibeles Russo

cibele@icmc.usp.br

<http://www.icmc.usp.br/~cibele>

Baseado em Johnson, R. A., & Wichern, D. W. (2007). Applied Multivariate Statistical Analysis. Prentice Hall.

Análise de componentes principais

A Análise de componentes principais (ACP ou PCA, de *principal component analysis*) é uma técnica que **transforma linearmente** um conjunto de p **variáveis correlacionadas** em um conjunto de k **variáveis não correlacionadas** (com $k < p$), que explicam uma parcela substancial das informações do conjunto original.

A ACP foi proposta por Karl Pearson em 1901 e estudada e nomeada por Harold Hotelling em 1933.

Se as variáveis originais tiverem distribuição normal multivariada, as componentes principais também terão distribuição normal multivariada e serão independentes.

Análise de componentes principais

A Análise de componentes principais (ACP ou PCA, de *principal component analysis*) é uma técnica que **transforma linearmente** um conjunto de p **variáveis correlacionadas** em um conjunto de k **variáveis não correlacionadas** (com $k < p$), que explicam uma parcela substancial das informações do conjunto original.

A ACP foi proposta por Karl Pearson em 1901 e estudada e nomeada por Harold Hotelling em 1933.

Se as variáveis originais tiverem distribuição normal multivariada, as componentes principais também terão distribuição normal multivariada e serão independentes.

Análise de componentes principais

A Análise de componentes principais (ACP ou PCA, de *principal component analysis*) é uma técnica que **transforma linearmente** um conjunto de p **variáveis correlacionadas** em um conjunto de k **variáveis não correlacionadas** (com $k < p$), que explicam uma parcela substancial das informações do conjunto original.

A ACP foi proposta por Karl Pearson em 1901 e estudada e nomeada por Harold Hotelling em 1933.

Se as variáveis originais tiverem distribuição normal multivariada, as componentes principais também terão distribuição normal multivariada e serão independentes.

Objetivos principais

- Reduzir a dimensionalidade dos dados.
- Obter combinações interpretáveis das variáveis originais.
- Descrever e compreender a estrutura de correlação das variáveis originais.

Começaremos com a obtenção das **componentes principais exatas**, extraídas da matriz de variâncias e covariâncias populacionais, Σ , e depois veremos as **componentes principais estimadas**, obtidas da matriz de variâncias e covariâncias amostrais, S , quando Σ é desconhecida.

Objetivos principais

- Reduzir a dimensionalidade dos dados.
- Obter combinações interpretáveis das variáveis originais.
- Descrever e compreender a estrutura de correlação das variáveis originais.

Começaremos com a obtenção das **componentes principais exatas**, extraídas da matriz de variâncias e covariâncias populacionais, Σ , e depois veremos as **componentes principais estimadas**, obtidas da matriz de variâncias e covariâncias amostrais, S , quando Σ é desconhecida.

Objetivos principais

- Reduzir a dimensionalidade dos dados.
- Obter combinações interpretáveis das variáveis originais.
- Descrever e compreender a estrutura de correlação das variáveis originais.

Começaremos com a obtenção das **componentes principais exatas**, extraídas da matriz de variâncias e covariâncias populacionais, Σ , e depois veremos as **componentes principais estimadas**, obtidas da matriz de variâncias e covariâncias amostrais, S , quando Σ é desconhecida.

Objetivos principais

- Reduzir a dimensionalidade dos dados.
- Obter combinações interpretáveis das variáveis originais.
- Descrever e compreender a estrutura de correlação das variáveis originais.

Começaremos com a obtenção das **componentes principais exatas**, extraídas da matriz de variâncias e covariâncias populacionais, Σ , e depois veremos as **componentes principais estimadas**, obtidas da matriz de variâncias e covariâncias amostrais, S , quando Σ é desconhecida.

Contexto

Seja \underline{X} um vetor aleatório de dimensão $p \times 1$ com vetor de médias (populacionais) $\underline{\mu}_{p \times 1}$ e matriz de variâncias e covariâncias (populacionais) de $\Sigma_{p \times p}$.

Estamos particularmente interessados no caso em que as variáveis X_1, \dots, X_p estão **correlacionadas**, isto é, algumas (ou muitas) das covariâncias $\text{Cov}(X_i, X_j), i, j = 1, \dots, p$ e $i \neq j$ são não-nulas.

Contexto

Quando isso ocorre, significa que existe **redundância entre dimensões**, e podemos ter interesse em **reduzir a dimensionalidade do problema**, construindo novas variáveis, **não correlacionadas** entre si, que sejam **combinações lineares das variáveis originais**.

Pode ser que poucas (k) novas variáveis expliquem grande parte da variabilidade existente nas (p) variáveis originais. Isso pode significar a **redução de custos** como tempo computacional e espaço para armazenamento de dados.

Análise de componentes principais

Seja $\underline{X} \sim (\underline{\mu}, \Sigma)$.

Sejam $\lambda_1 \geq \dots \geq \lambda_p$ os autovalores de Σ , com autovetores correspondentes $\underline{e}_1, \dots, \underline{e}_p$, tais que

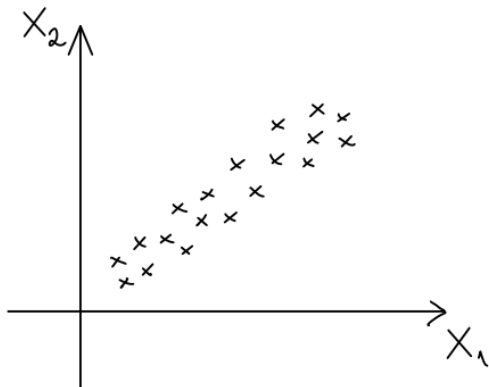
- 1 $\underline{e}_i^\top \underline{e}_j = 0$, para $i, j = 1, \dots, p$ e $i \neq j$,
- 2 $\underline{e}_i^\top \underline{e}_i = 1$, para $i = 1, \dots, p$,
- 3 $\Sigma \underline{e}_i = \lambda_i \underline{e}_i$, para $i = 1, \dots, p$.

Considere a matriz ortogonal $O_{p \times p} = (\underline{e}_1, \dots, \underline{e}_p)$.

Então $\underline{Y}_{p \times 1} = O^\top \underline{X}$ é o **vetor de componentes principais** de Σ .

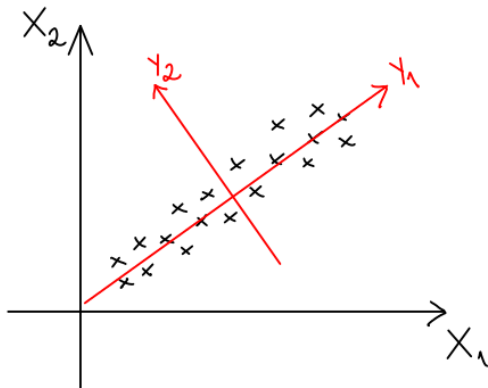
Análise de componentes principais

Interpretação geométrica ($p = 2$)



Análise de componentes principais

Interpretação geométrica ($p = 2$)



Análise de componentes principais

Propriedades:

- 1 A j -ésima componente principal de Σ é dada por

$$Y_j = \underline{e}_j^\top \underline{X}.$$

- 2 $E(Y_j) = \underline{e}_j^\top \underline{\mu}$.
- 3 $\text{Var}(Y_j) = \underline{e}_j^\top \Sigma \underline{e}_j = \lambda_j$.
- 4 $\text{Cov}(Y_i, Y_j) = \text{Cor}(Y_i, Y_j) = 0$ para $i, j = 1, \dots, p$ e $i \neq j$.
- 5 A proporção da variância total de \underline{X} que é explicada pela j -ésima componente principal é

$$\frac{\lambda_j}{\sum_{i=1}^p \lambda_i}.$$

Estimação das componentes principais

Como em geral a matriz Σ é desconhecida, utiliza-se a matriz S , de variâncias e covariâncias amostrais, para estimar as componentes principais.

Considere $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$ os autovalores de S , com autovetores correspondentes padronizados $\hat{\underline{e}}_1, \dots, \hat{\underline{e}}_p$.

A j -ésima componente principal amostral é dada por

$$\hat{Y}_j = \hat{\underline{e}}_j^T \underline{X}.$$

Estimação das componentes principais

Propriedades:

- 1 $\widehat{\text{Var}}(\hat{Y}_j) = \hat{\lambda}_j.$
- 2 $\text{Cov}(\hat{Y}_i, \hat{Y}_j) = \text{Cor}(\hat{Y}_i, \hat{Y}_j) = 0$ para $i, j = 1, \dots, p$ e $i \neq j.$
- 3 A proporção da variância total explicada pela j -ésima componente principal amostral é

$$\frac{\hat{\lambda}_j}{\sum_{i=1}^p \hat{\lambda}_i}$$

- 4 $\widehat{\text{Cor}}(\hat{Y}_j, X_i) = \frac{\hat{e}_{ji} \sqrt{\hat{\lambda}_j}}{\sqrt{s_{jj}}}$ (Exercício)

Estimação das componentes principais

- 5 Pelo teorema da decomposição espectral,

$$S_{p \times p} = \sum_{j=1}^p \hat{\lambda}_j \hat{e}_j \hat{e}_j^T$$

pode ser aproximada por

$$S_{p \times p} \approx \sum_{j=1}^k \hat{\lambda}_j \hat{e}_j \hat{e}_j^T$$

Escores das componentes

É comum utilizar os **escores das componentes** para análises estatísticas ou ordenação (*ranking*) dos elementos amostrais. Esses escores são obtidos ordenando os valores obtidos de

$$\hat{Y}_j = \hat{e}_j^\top \underline{X}.$$

Análise de componentes principais via matriz de correlação

Seja $\underline{X}_{p \times 1} \sim (\underline{\mu}, \Sigma)$, $\underline{X} = (X_1, \dots, X_p)^\top$.

Seja $\underline{Z} = (Z_1, \dots, Z_p)^\top$ tal que

$$Z_i = \frac{X_i - \mu_i}{\sigma_i},$$

em que $\mu_i = E(X_i)$ e $\sigma_i^2 = \text{Var}(X_i)$.

Temos que $\text{Cov}(\underline{Z}) = P = \text{Cor}(\underline{X})$.

Análise de componentes principais via matriz de correlação

Sejam $\lambda_1 \geq \dots \geq \lambda_p$ os autovalores de P , com autovetores correspondentes $\underline{e}_1, \dots, \underline{e}_p$, tais que

- 1 $\underline{e}_i^\top \underline{e}_j = 0$, para $i, j = 1, \dots, p$ e $i \neq j$,
- 2 $\underline{e}_i^\top \underline{e}_i = 1$, para $i = 1, \dots, p$,
- 3 $\Sigma \underline{e}_i = \lambda_i \underline{e}_i$, para $i = 1, \dots, p$.

Considere a matriz ortogonal $O_{p \times p} = (\underline{e}_1, \dots, \underline{e}_p)$.

Então $\underline{Y}_{p \times 1} = O^\top \underline{Z}$ é o **vetor de componentes principais** de P .

Análise de componentes principais via matriz de correlação

Propriedades:

- 1 A j -ésima componente principal de P é dada por

$$Y_j = e_j^\top \underline{Z}.$$

- 2 $E(Y_j) = 0$.
- 3 $\text{Var}(Y_j) = e_j^\top P e_j = \lambda_j$.
- 4 $\text{Cov}(Y_i, Y_j) = \text{Cor}(Y_i, Y_j) = 0$ para $i, j = 1, \dots, p$ e $i \neq j$.
- 5 A proporção da variância total de \underline{Z} que é explicada pela j -ésima componente principal é λ_j/p

Estimação das componentes principais

Como em geral a matriz P é desconhecida, utiliza-se a matriz R , de correlações amostrais, para estimar as componentes principais.

Considere $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$ os autovalores de R , com autovetores correspondentes padronizados $\hat{e}_1, \dots, \hat{e}_p$.

A j -ésima componente principal amostral é dada por

$$\hat{Y}_j = \hat{e}_j^\top \underline{Z}.$$

Estimação das componentes principais

Propriedades:

- 1 $\widehat{\text{Var}}(\hat{Y}_j) = \hat{\lambda}_j.$
- 2 $\text{Cov}(\hat{Y}_i, \hat{Y}_j) = \text{Cor}(\hat{Y}_i, \hat{Y}_j) = 0$ para $i, j = 1, \dots, p$ e $i \neq j.$
- 3 A proporção da variância total explicada pela j -ésima componente principal amostral é

$$\frac{\hat{\lambda}_j}{p}$$

- 4 $\widehat{\text{Cor}}(\hat{Y}_j, Z_i) = \hat{e}_{ji} \sqrt{\hat{\lambda}_j}$ (Exercício)

Estimação das componentes principais

- 5 Pelo teorema da decomposição espectral,

$$R_{p \times p} = \sum_{j=1}^p \hat{\lambda}_j \hat{e}_j \hat{e}_j^{\top}$$

pode ser aproximada por

$$R_{p \times p} \approx \sum_{j=1}^k \hat{\lambda}_j \hat{e}_j \hat{e}_j^{\top}$$

Escores das componentes

É comum utilizar os **escores das componentes** para análises estatísticas ou ordenação (*ranking*) dos elementos amostrais. Esses escores são obtidos ordenando os valores obtidos de

$$\hat{Y}_j = \hat{e}_j^\top \underline{Z}.$$

Determinação do número de componentes principais

Alguns métodos:

- 1 Proporção da variância total explicada.
- 2 Análise gráfica da variância explicada (scree-plot).
- 3 Aproximação da matriz S (ou R).
- 4 Análise prática das componentes principais.

Inferência assintótica sobre as componentes principais

Sejam $\underline{X} \sim (\underline{\mu}, \Sigma)$, $\lambda_1 \geq \dots \geq \lambda_p > 0$ os autovalores de Σ , com autovetores correspondentes $\underline{e}_1, \dots, \underline{e}_p$, $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p > 0$ os autovalores de S , com autovetores correspondentes $\hat{\underline{e}}_1, \dots, \hat{\underline{e}}_p$.

Sejam $\underline{\lambda} = (\lambda_1, \dots, \lambda_p)^\top$ e $\hat{\underline{\lambda}} = (\hat{\lambda}_1, \dots, \hat{\lambda}_p)^\top$.

Inferência assintótica sobre as componentes principais

Resultados:

- ① Se $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$, então, para n suficientemente grande, temos

$$\sqrt{n}(\hat{\lambda} - \lambda) \approx N_p(\mathbf{0}, 2\Lambda^2).$$

- ② Se $E_i = \lambda_i \sum_{k=1, k \neq i}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} \underline{e}_k \underline{e}_k^\top$, então

$$\sqrt{n}(\hat{e}_i - e_i) \approx N_p(\mathbf{0}, E_i).$$

- ③ A distribuição de cada $\hat{\lambda}_i$ não depende dos elementos de \hat{e}_i correspondentes.

Inferência assintótica sobre as componentes principais

Segue desses resultados que

$$\hat{\lambda}_i \stackrel{ind}{\sim} N\left(\lambda_i, 2\frac{\lambda_i^2}{n}\right).$$

Assim,

$$P\left(|\hat{\lambda}_i - \lambda_i| \leq z_{\alpha/2} \lambda_i \sqrt{\frac{2}{n}}\right) = 1 - \alpha.$$

Inferência assintótica sobre as componentes principais

Segue desses resultados que

$$\hat{\lambda}_i \stackrel{ind}{\sim} N\left(\lambda_i, 2\frac{\lambda_i^2}{n}\right).$$

Assim,

$$P\left(|\hat{\lambda}_i - \lambda_i| \leq z_{\alpha/2} \lambda_i \sqrt{\frac{2}{n}}\right) = 1 - \alpha.$$

Inferência assintótica sobre as componentes principais

Segue desses resultados que

$$\hat{\lambda}_i \stackrel{ind}{\sim} N\left(\lambda_i, 2\frac{\lambda_i^2}{n}\right).$$

Assim,

$$P\left(|\hat{\lambda}_i - \lambda_i| \leq z_{\alpha/2} \lambda_i \sqrt{\frac{2}{n}}\right) = 1 - \alpha.$$

Inferência assintótica sobre as componentes principais

E podemos então construir um intervalo com $100(1 - \alpha)\%$ de confiança para λ_i usando propriedades da distribuição normal:

$$IC_{100(1-\alpha)\%}(\lambda_i) = \left(\frac{\hat{\lambda}_i}{1 + z_{\frac{\alpha}{2}} \sqrt{\frac{2}{n}}}, \frac{\hat{\lambda}_i}{1 - z_{\frac{\alpha}{2}} \sqrt{\frac{2}{n}}} \right),$$

que também pode ser usado para fazer inferência sobre λ_i (por exemplo, avaliar se $H_0 : \lambda_i = \lambda_0$ contra $H_1 : \lambda_i \neq \lambda_0$).

Exemplo

Estudos mostram que grande parte de adultos e adolescentes norte-americanos usam regularmente substâncias psicoativas. Em um destes estudos (Huba et al. 1981, J. of Personality and Social Psychology), dados foram coletados de 1634 estudantes na área metropolitana de Los Angeles. Cada participante completou um questionário informando o número de vezes que cada item foi usado.

Os itens são os seguintes: cigarro, cerveja, vinho, licor, cocaína, tranquilizantes, medicamentos, heroína, maconha, haxixe, inalantes, alucinógenos e anfetaminas.

Exemplo

As respostas foram registradas em uma escala de cinco pontos: 1. nunca experimentei, 2. apenas uma vez, 3. poucas vezes, 4. muitas vezes e 5. regularmente.

A matriz de correlações das respostas encontra-se no arquivo Substancias.txt.

Foram obtidas as componentes principais a matriz de correlações dos dados. Como propor um “índice de utilização de substâncias psicoativas”?