

# Imputação múltipla para matriz de interação genótipo e ambiente

Welinton Y. Hirai & Carlos T. S. Dias

Departamento de Ciências Exatas  
Estatística e Experimentação Agronômica

---



Piracicaba  
2018

# Estrutura da apresentação

---

1. Introdução
2. Problemática
3. Metodologia
4. Aplicação
5. Estudos Futuros

# Estrutura da apresentação

---

1. Introdução
2. Problemática
3. Metodologia
4. Aplicação
5. Estudos Futuros

# Abordagem estatística na genética

---

“As metodologias estatísticas são ferramentas cruciais para análises de dados genômicos, no qual atualmente está se tornando cada vez mais útil para avaliações para variedades de espécies, e ainda oferecer com precisão descrições sobre a variação genética que ocorre entre espécies, populações, famílias e indivíduos.”

*Wu, R., Ma, C., & Casella, G. (2007). Statistical genetics of quantitative traits: linkage, maps and QTL. Springer Science & Business Media.*

FLORES, F.; MORENO, M. .; CUBERO, J. . . A comparison of univariate and multivariate methods to analyze GxEx interaction. *Field Crops Research*, v. 56, n. 3, p. 271–286, 1 abr. 1998.

M é t o d o s  e s t a t í s t i c o s	univariados paramétricos	Roemer (1917) Becker e Leon (1998)	Mensura os desvios para a média genética. Variância mínima em diferentes ambientes representa estabilidade genética.
		Eberhart e Russel (1966) Tai (1971)	Estuda a estabilidade de um experimental ótimo por meio de regressão.
		Shukla (1972)	Estima os componentes de variância da IGA atribuído por cada genótipo.
	univariados não- paramétricos	Hühn (1979)	Propôs métricas utilizando estatística não paramétrica, no qual valores baixo representa alta estabilidade.
		Ketata et al. (1989)	Representa graficamente a média dos “ <i>rankeamentos</i> ” em relação aos ambiente, em função dos desvios de <i>rank</i> para todos os genótipos .
		Flores (1993)	Gera polígonos para cada genótipos, no qual cada desenho representa os valores médios em cada ambiente.
	multivariados	Lin (1982)	Calcula valores de similaridades para cada par de genótipos, e mensura utilizando análise de agrupamentos.
		Westcott (1987)	Utiliza escalonamento multidimensional de cada genótipo em cada ambiente.

# Interação genótipo x ambiente

“Para o pesquisador, este processo que diferencia a adaptabilidade da estabilidade é necessária, dado que, um caso implica a não perda de produção em diferentes ecossistemas, enquanto o outro caso pode providenciar uma vantagem significativa em ambientes específicos.”

Wade, L., C. McLaren, L. Quintana, D. Harnpichitvitaya, S. Rajatasereekul, A. Sarawgi, A. Kumar, H. Ahmed, A. Singh, R. Rodriguez, J. Siopongco, and S. Sarkarung, 1999  
Genotype by environment interactions across diverse rainfed lowland rice environments. *F. Crop. Res.* **64**: 35–50.

“Há casos em que as respostas dos efeitos genotípicos e ambientais não são capazes de serem descritas separadamente, levando assim a necessidade de utilizar um termo que representa conjuntamente a contribuição destes efeitos.”

Gregorius, H.-R. and G. Namkoong, 1986 Joint analysis of genotypic and environmental effects. *Theor Appl Genet* **72**: 413–422.



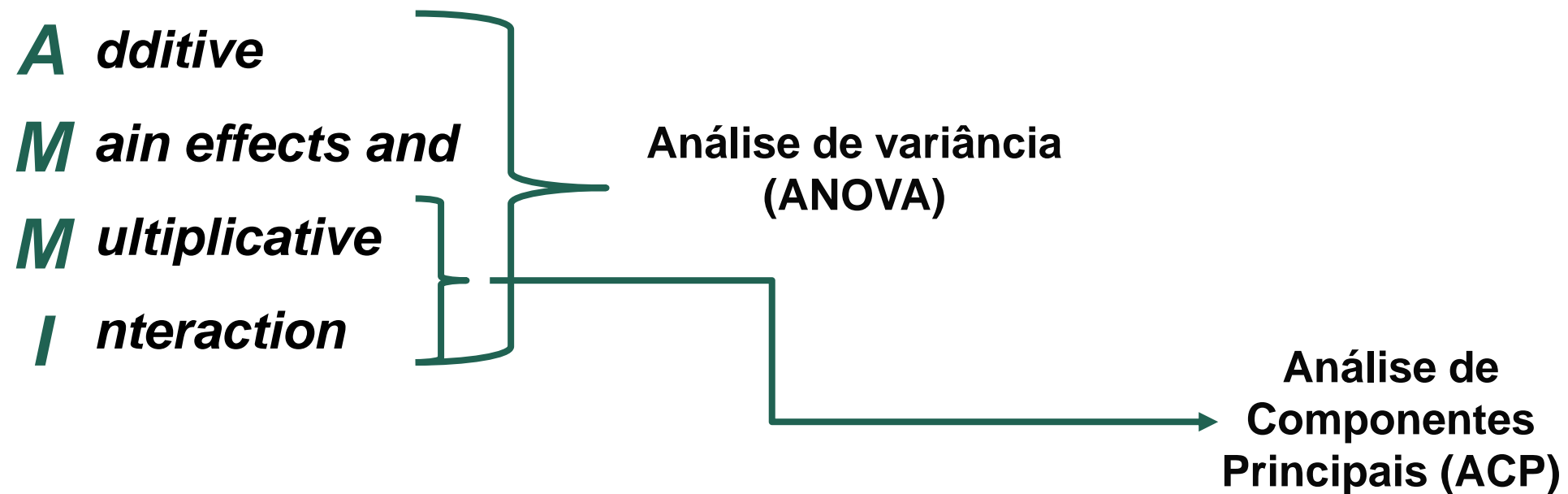
# Modelo AMMI

MANDEL, J. A New Analysis of Variance Model for Non-Additive Data. **Technometrics**, v. 13, n. 1, p. 1–18, 1971.

GAUCH, H. G.; ZOBEI, R. W. Predictive and postdictive success of statistical analyses of yield trials\*. **Theor Appl Genet**, v. 76, p. 1–10, 1988.

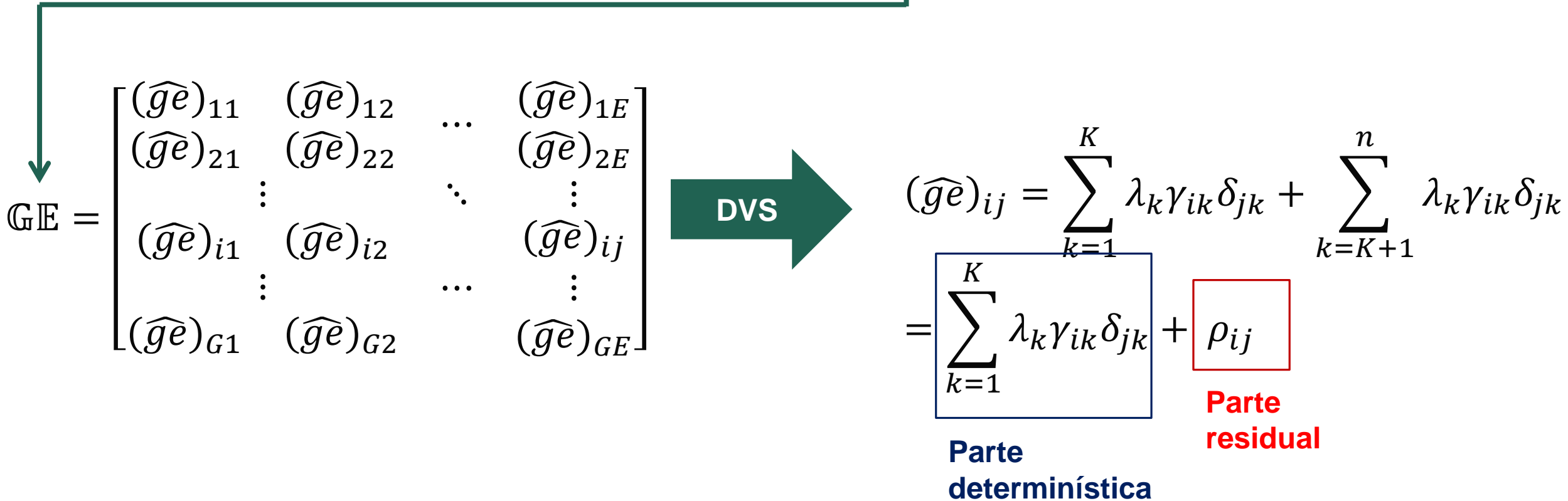
---

## Model



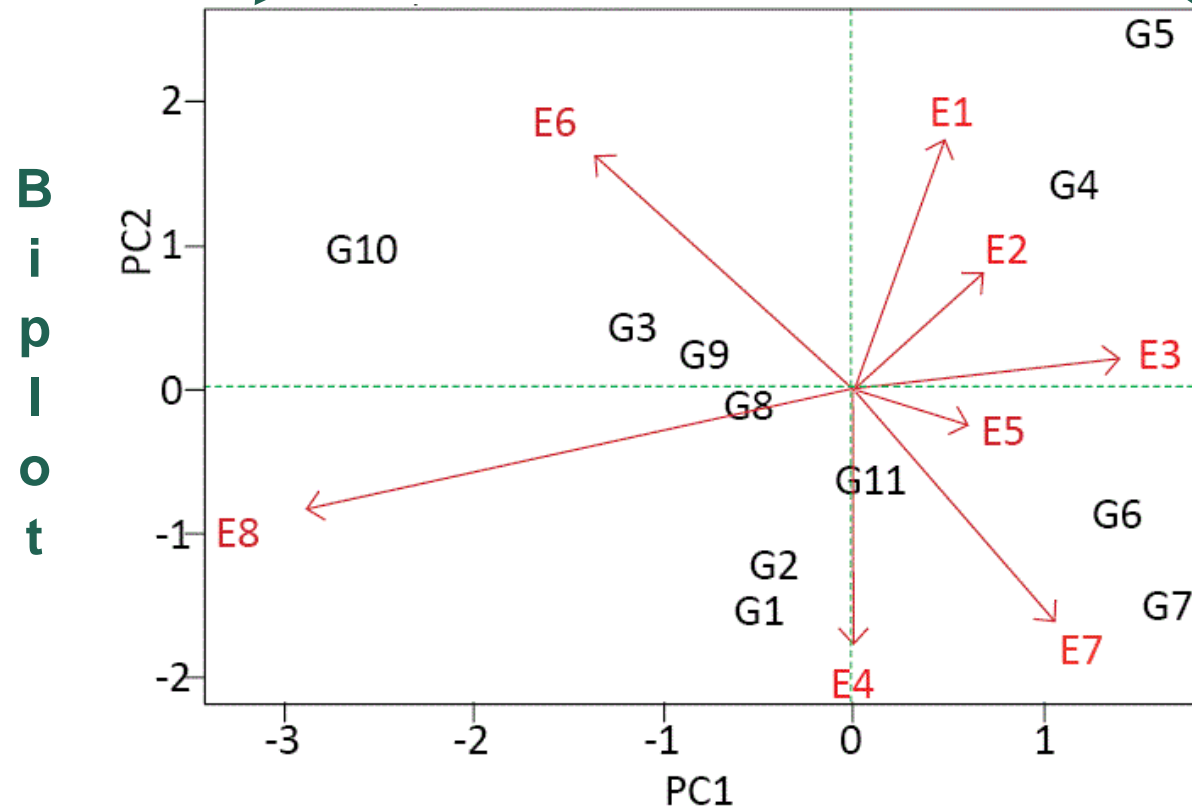
## Modelo estatístico para ANOVA conjunta

$$Y_{ijk} = \mu + g_i + e_j + \boxed{(ge)_{ij}} + \epsilon_{ijk}$$



## Modelo AMMI

$$Y_{ijk} = \mu + g_i + e_j + \sum_{k=1}^K \lambda_k \gamma_{ik} \delta_{jk} + \rho_{ij} + \epsilon_{ijk}$$



# Estrutura da apresentação

---

1. Introdução
- 2. Problemática**
3. Metodologia
4. Aplicação
5. Estudos Futuros

## Ausência no conjunto de dados

---

$$\mathbb{G}E = \begin{bmatrix} (\widehat{ge})_{11} & (\widehat{ge})_{12} & \dots & (\widehat{ge})_{1E} \\ (\widehat{ge})_{21} & (\widehat{ge})_{22} & \dots & (\widehat{ge})_{2E} \\ NA & (\widehat{ge})_{32} & & (\widehat{ge})_{3E} \\ (\widehat{ge})_{41} & (\widehat{ge})_{42} & \dots & (\widehat{ge})_{4E} \\ & \vdots & & \vdots \\ (\widehat{ge})_{i1} & (\widehat{ge})_{i2} & & NA \\ & \vdots & \dots & \vdots \\ (\widehat{ge})_{G1} & (\widehat{ge})_{G2} & & (\widehat{ge})_{GE} \end{bmatrix}$$

# Estrutura da apresentação

---

1. Introdução
2. Problemática
- 3. Metodologia**
4. Aplicação
5. Estudos Futuros

# Imputação Múltipla Livre de Distribuição

Missing value imputation in multivariate data using the singular value decomposition of a matrix	<i>Krzanowski, Wojtek Janusz</i>	1988
Imputação múltipla livre de distribuição utilizando a decomposição por valor singular em matriz de interação	<i>Bergamo, Genevile Carife</i>	2007
Imputação de dados em experimentos com interação genótipo por ambiente: uma aplicação a dados de algodão	<i>Arciniegas-Alarcón, Sergio</i>	2008
An alternative methodology for imputing missing data in trials with genotype-by-environment interaction: some new aspects	<i>Arciniegas-Alarcón, Sergio García-Peña, Marisol Krzanowski, Wojtek Janusz Dias, Carlos Tadeu dos Santos</i>	2014
Imputação de dados em experimentos multiambientais: novos algoritmos utilizando a decomposição por valores singulares	<i>Arciniegas-Alarcón, Sergio</i>	2015
Missing value imputation in multi-environment trials: Reconsidering the Krzanowski method	<i>Arciniegas-Alarcón, Sergio García-Peña, Marisol Krzanowski, Wojtek Janusz</i>	2016



## Teorema da Decomposição em Valores Singulares (DVS)

---

$$\mathbb{X}_{(n,p)} = \mathbb{U}_{(n,p)} \Phi_{(p,p)} \mathbb{V}_{(p,p)}^T$$

Suponha-se 3 valores ausentes, numa matriz  $\mathbb{X}_{(n,p)}$

$$\left[ x_{miss}^{[1]}, x_{miss}^{[2]}, x_{miss}^{[3]} \right]$$

$$\mathbb{X}_{miss} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \dots & x_{1p} \\ x_{miss}^{[1]} & x_{22} & x_{23} & x_{24} & \dots & x_{2p} \\ x_{31} & x_{32} & x_{miss}^{[2]} & x_{34} & \dots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \dots & x_{miss}^{[3]} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \dots & x_{np} \end{bmatrix}$$

# 1º processo iterativo da imputação múltipla ( $impu_1$ )

$$\mathbb{X}_{miss} = \begin{array}{cccccc} \bar{x}_1 S_1 & \bar{x}_2 S_2 & \bar{x}_3 S_3 & \bar{x}_4 S_4 & \dots & \bar{x}_p S_p \\ \begin{array}{c} x_{11} \\ x_{miss}^{[1]} \\ x_{31} \\ x_{41} \\ \vdots \\ x_{n1} \end{array} & \begin{array}{c} x_{12} \\ x_{22} \\ x_{32} \\ x_{42} \\ \vdots \\ x_{n2} \end{array} & \begin{array}{c} x_{13} \\ x_{23} \\ x_{miss}^{[2]} \\ x_{43} \\ \vdots \\ x_{n3} \end{array} & \begin{array}{c} x_{14} \\ x_{24} \\ x_{34} \\ x_{44} \\ \vdots \\ x_{n4} \end{array} & \dots & \begin{array}{c} x_{1p} \\ x_{2p} \\ x_{3p} \\ x_{miss}^{[3]} \\ \vdots \\ x_{np} \end{array} \end{array}$$

Substitui-se os valores ausente pelas média referente à coluna

$$\mathbb{X}_{miss} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \cdots & x_{1p} \\ \bar{x}_1 & x_{22} & x_{23} & x_{24} & \cdots & x_{2p} \\ x_{31} & x_{32} & \bar{x}_3 & x_{34} & \cdots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \cdots & \bar{x}_p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \cdots & x_{np} \end{bmatrix}$$

Para a **primeira observação** que será imputada, tem-se os passos:

$$\mathbb{X}_{miss} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \dots & x_{1p} \\ x_{[1]}^{miss} & x_{22} & x_{23} & x_{24} & \dots & x_{2p} \\ x_{31} & x_{32} & x_{[2]}^{miss} & x_{34} & \dots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \dots & x_{[3]}^{miss} \\ \vdots & & & \vdots & & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \dots & x_{np} \end{bmatrix}$$

# 1ª Passo: Gerar duas matriz a partir da matriz com os valores médios

1.1 Matriz sem a linha que está o primeiro valor ausente

1.2 Matriz sem a coluna que está o primeiro valor ausente

$$\mathbb{X}_{-i_{miss}} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \cdots & x_{1p} \\ \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ x_{31} & x_{32} & \bar{x}_3 & x_{34} & \cdots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \cdots & \bar{x}_p \\ \vdots & & & \vdots & & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \cdots & x_{np} \end{bmatrix}$$

$$\mathbb{X}_{-j_{miss}} = \begin{bmatrix} \text{---} & x_{12} & x_{13} & x_{14} & \cdots & x_{1p} \\ \text{---} & x_{22} & x_{23} & x_{24} & \cdots & x_{2p} \\ \text{---} & x_{32} & \bar{x}_3 & x_{34} & \cdots & x_{3p} \\ \text{---} & x_{42} & x_{43} & x_{44} & \cdots & \bar{x}_p \\ \vdots & \vdots & & \vdots & & \vdots \\ \text{---} & x_{n2} & x_{n3} & x_{n4} & \cdots & x_{np} \end{bmatrix}$$

**2ª Passo:** Aplicar a padronização nas duas matrizes

$$\mathbb{X}_{-i_{miss}}^* = \begin{bmatrix} x_{11}^* & x_{12}^* & x_{13}^* & x_{14}^* & \cdots & x_{1p}^* \\ x_{31}^* & x_{32}^* & \bar{x}_3 & x_{34}^* & \cdots & x_{3p}^* \\ x_{41}^* & x_{42}^* & x_{43}^* & x_{44}^* & \cdots & \bar{x}_p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1}^* & x_{n2}^* & x_{n3}^* & x_{n4}^* & \cdots & x_{np}^* \end{bmatrix}$$

$$\mathbb{X}_{-i_{miss}}^* = \begin{bmatrix} & x_{12}^* & x_{13}^* & x_{14}^* & \cdots & x_{1p}^* \\ & x_{22}^* & x_{23}^* & x_{24}^* & \cdots & x_{2p}^* \\ & x_{32}^* & \bar{x}_3 & x_{34}^* & \cdots & x_{3p}^* \\ & x_{42}^* & x_{43}^* & x_{44}^* & \cdots & \bar{x}_p \\ & \vdots & \vdots & \vdots & \vdots & \vdots \\ & x_{n2}^* & x_{n3}^* & x_{n4}^* & \cdots & x_{np}^* \end{bmatrix}$$

**3ª Passo:** Calcular a DVS para as duas matrizes padronizadas

$$\mathbb{X}_{-i_{miss}}^* = \begin{bmatrix} x_{11}^* & x_{12}^* & x_{13}^* & x_{14}^* & \dots & x_{1p}^* \\ x_{31}^* & x_{32}^* & \bar{x}_3 & x_{34}^* & \dots & x_{3p}^* \\ x_{41}^* & x_{42}^* & x_{43}^* & x_{44}^* & \dots & \bar{x}_p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1}^* & x_{n2}^* & x_{n3}^* & x_{n4}^* & \dots & x_{np}^* \end{bmatrix}$$

$$\mathbb{X}_{(n-1,p)}^* = \begin{matrix} \mathbb{U}_{(n-1,p)} \Phi_{(p,p)} \mathbb{V}_{(p,p)}^T \\ \mathbb{U}_{-i_{miss}} \Phi_{-i_{miss}} \mathbb{V}_{-i_{miss}}^T \end{matrix}$$

$$\mathbb{X}_{-i_{miss}}^* = \begin{bmatrix} x_{12}^* & x_{13}^* & x_{14}^* & \dots & x_{1p}^* \\ x_{22}^* & x_{23}^* & x_{24}^* & \dots & x_{2p}^* \\ x_{32}^* & \bar{x}_3 & x_{34}^* & \dots & x_{3p}^* \\ x_{42}^* & x_{43}^* & x_{44}^* & \dots & \bar{x}_p \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n2}^* & x_{n3}^* & x_{n4}^* & \dots & x_{np}^* \end{bmatrix}$$

$$\mathbb{X}_{(n,p-1)}^* = \begin{matrix} \mathbb{U}_{(n,p-1)} \Phi_{(p-1,p-1)} \mathbb{V}_{(p-1,p-1)}^T \\ \mathbb{U}_{-j_{miss}} \Phi_{-j_{miss}} \mathbb{V}_{-j_{miss}}^T \end{matrix}$$

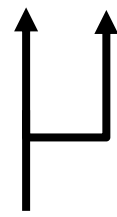


**4ª Passo:** Calcular a matriz com o valor imputado, utilizando as decomposições das duas matrizes

$$Y_{impu1}^* = U_{-j_{miss}} \Phi_{-j_{miss}} \left( \frac{a}{b} \right) \Phi_{-i_{miss}} \left( \frac{b-a}{b} \right) V_{-i_{miss}}^T$$

Segundo Bergamo (2007) são utilizados cinco processos iterativos na imputação, com valores para  $a = 8, 9, 10, 11, 12$  e  $b = 20$

$$Y_{(n,p)}^* = U_{(n,p-1)} \Phi_{(p-1,p-1)} \left( \frac{a}{b} \right) \Phi_{(p,p)} \left( \frac{b-a}{b} \right) V_{(p,p)}^T$$



Desta forma,  
é descartado  
o último  
autovalor



Em por  
consequência, é  
retirada a última  
coluna

**5ª Passo:** “Despadronizar” a matriz calculada utilizando os valores de média e desvio-padrão encontradas no início

$$Y_{impu_1} = \begin{bmatrix} y_{11} & y_{12} & y_{13} & y_{14} & \cdots & y_{1p} \\ y_{impu_1}^{[1]} & y_{22} & y_{23} & y_{24} & \cdots & y_{2p} \\ y_{31} & y_{32} & y_{33} & y_{34} & \cdots & y_{3p} \\ y_{41} & y_{42} & y_{43} & y_{44} & \cdots & y_{4p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ y_{n1} & y_{n2} & y_{n3} & y_{n4} & \cdots & y_{np} \end{bmatrix}$$

## Para a todas as observações que serão imputadas, tem-se os passos:

**1ª Passo:** Gerar duas matriz a partir da matriz com os valores médios

1.1 Matriz sem a linha que está o primeiro valor ausente

1.2 Matriz sem a coluna que está o primeiro valor ausente

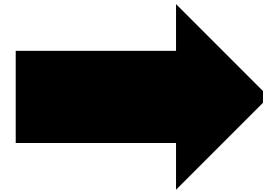
**2ª Passo:** Aplicar a padronização nas duas matrizes

**3ª Passo:** Calcular a DVS para as duas matrizes padronizadas

**4ª Passo:** Calcular a matriz com o valor imputado, utilizando as decomposições das duas matrizes

**5ª Passo:** “Despadronizar” a matriz calculada utilizando os valores de média e desvio-padrão encontradas no início

$$\left[ x_{miss}^{[1]}, x_{miss}^{[2]}, x_{miss}^{[3]} \right]$$



$$\left[ y_{impu_1}^{[1]}, y_{impu_1}^{[2]}, y_{impu_1}^{[3]} \right]$$

## No 1º processo iterativo da imputação múltipla ( $impu_1$ )

$$Y_{impu_1} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \dots & x_{1p} \\ y_{impu_1}^{[1]} & x_{22} & x_{23} & x_{24} & \dots & x_{2p} \\ x_{31} & x_{32} & y_{impu_1}^{[2]} & x_{34} & \dots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \dots & y_{impu_1}^{[3]} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \dots & x_{np} \end{bmatrix}$$

## No 2º processo iterativo da imputação múltipla ( $impu_2$ )

$$\mathbb{Y}_{impu_1} = \begin{array}{c} \bar{x}_1 S_1 \quad \bar{x}_2 S_2 \quad \bar{x}_3 S_3 \quad \bar{x}_4 S_4 \quad \dots \quad \bar{x}_p S_p \\ \left[ \begin{array}{ccccccc} x_{11} & x_{12} & x_{13} & x_{14} & \dots & x_{1p} \\ y_{impu_1}^{[1]} & x_{22} & x_{23} & x_{24} & \dots & x_{2p} \\ x_{31} & x_{32} & y_{impu_1}^{[2]} & x_{34} & \dots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \dots & y_{impu_1}^{[3]} \\ \vdots & & & & & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \dots & x_{np} \end{array} \right] \end{array}$$

## Para o 2º processo iterativo, segue-se os passos:

**1ª Passo:** Gerar duas matrizes a partir da matriz com os valores médios

1.1 Matriz sem a linha que está o primeiro valor ausente

1.2 Matriz sem a coluna que está o primeiro valor ausente

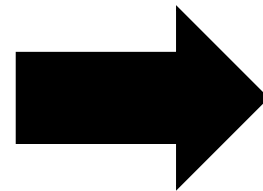
**2ª Passo:** Aplicar a padronização nas duas matrizes

**3ª Passo:** Calcular a DVS para as duas matrizes padronizadas

**4ª Passo:** Calcular a matriz com o valor imputado, utilizando as decomposições das duas matrizes

**5ª Passo:** “Despadronizar” a matriz calculada utilizando os valores de média e desvio-padrão encontradas no início

$$\left[ x_{miss}^{[1]}, x_{miss}^{[2]}, x_{miss}^{[3]} \right]$$



$$\left[ y_{impu_2}^{[1]}, y_{impu_2}^{[2]}, y_{impu_2}^{[3]} \right]$$

**No final do processo iterativo da imputação múltipla, tem-se que os valores imputados serão:**

$$\hat{y}_{impu}^{[1]} = \sum_{l=1}^M \frac{y_{impu_k}^{[1]}}{M} \quad \hat{y}_{impu}^{[2]} = \sum_{l=1}^M \frac{y_{impu_k}^{[2]}}{M} \quad \hat{y}_{impu}^{[3]} = \sum_{l=1}^M \frac{y_{impu_k}^{[3]}}{M}$$

$$\mathbb{X}_{miss} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \dots & x_{1p} \\ \hat{y}_{impu}^{[1]} & x_{22} & x_{23} & x_{24} & \dots & x_{2p} \\ x_{31} & x_{32} & \hat{y}_{impu}^{[2]} & x_{34} & \dots & x_{3p} \\ x_{41} & x_{42} & x_{43} & x_{44} & \dots & \hat{y}_{impu}^{[3]} \\ \vdots & & & \vdots & & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} & \dots & x_{np} \end{bmatrix}$$

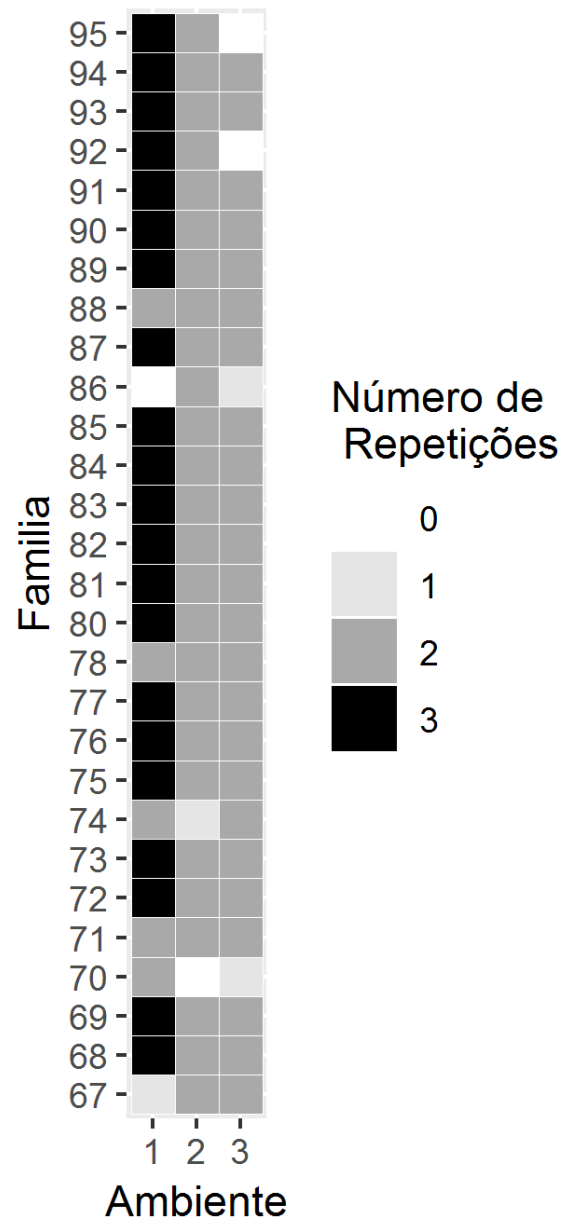
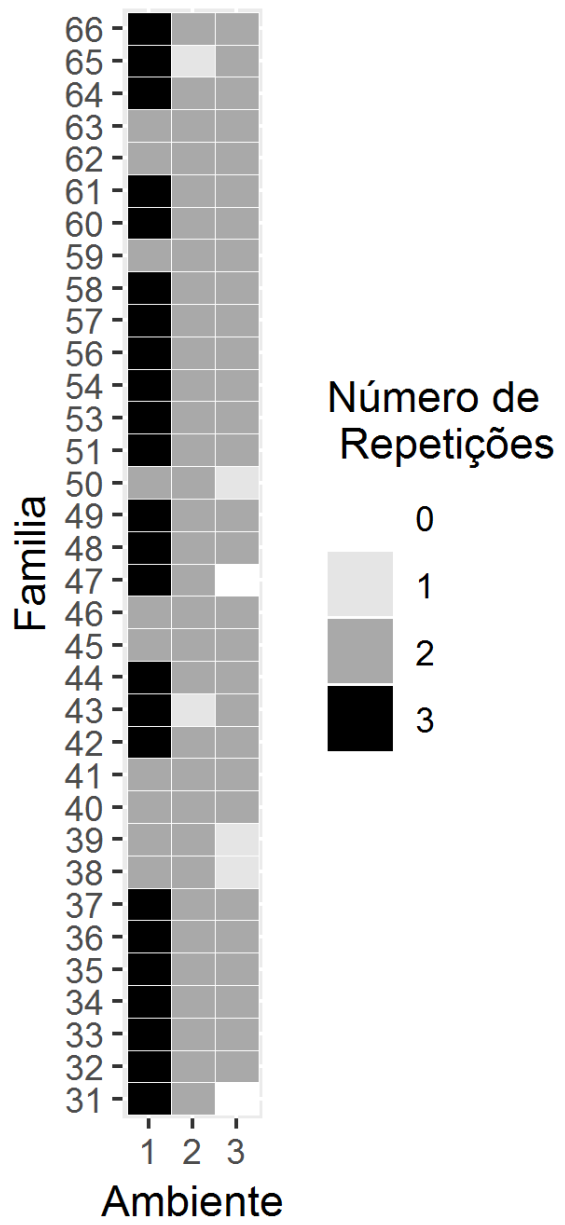
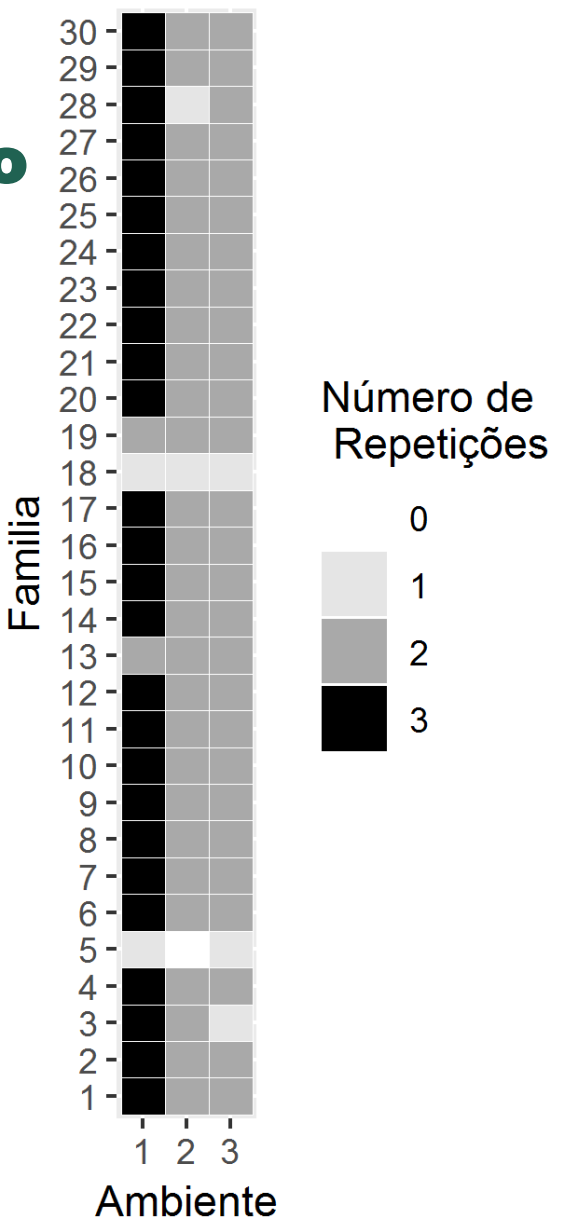
# Estrutura da apresentação

---

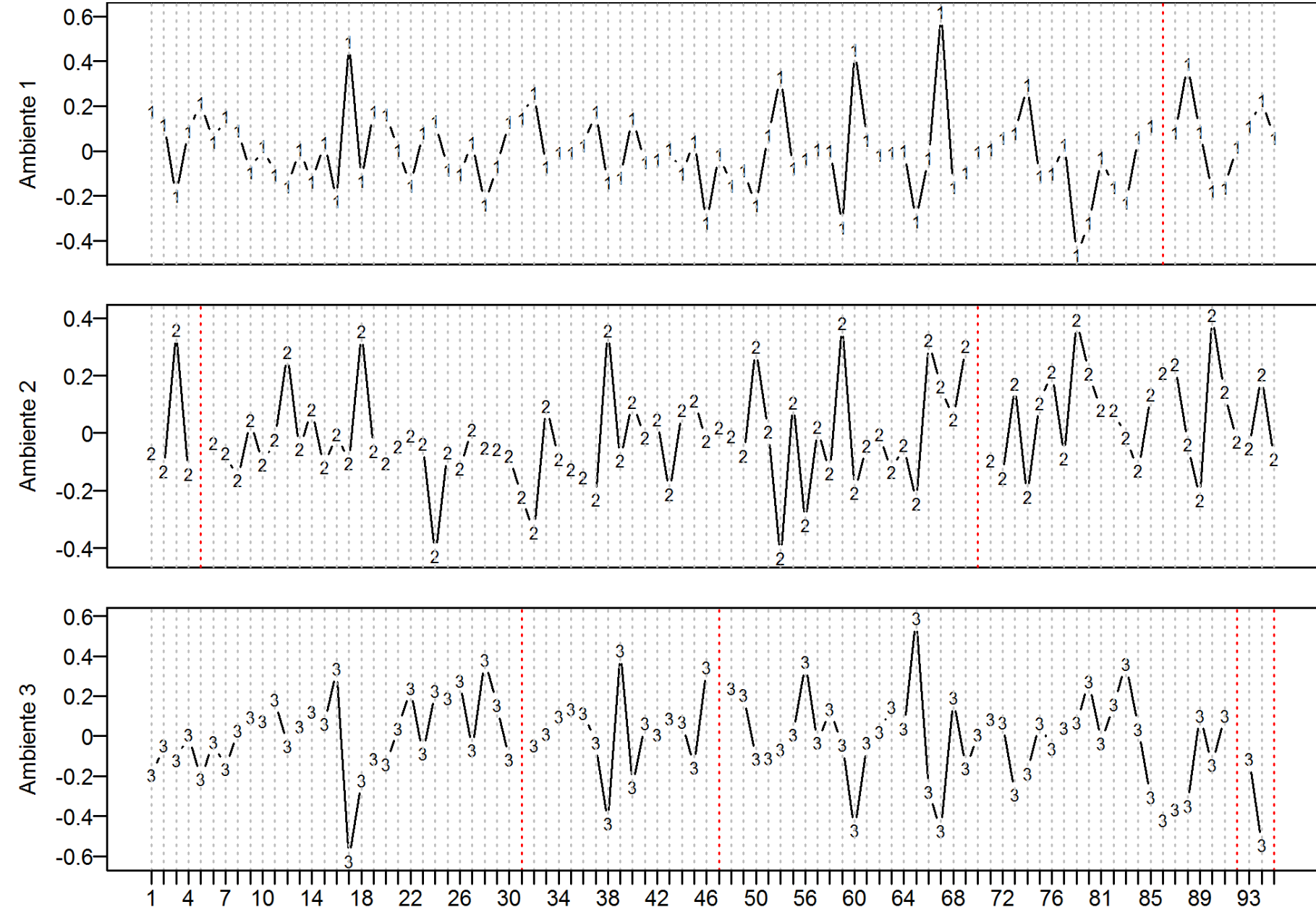
1. Introdução
2. Problemática
3. Metodologia
- 4. Aplicação**
5. Estudos Futuros



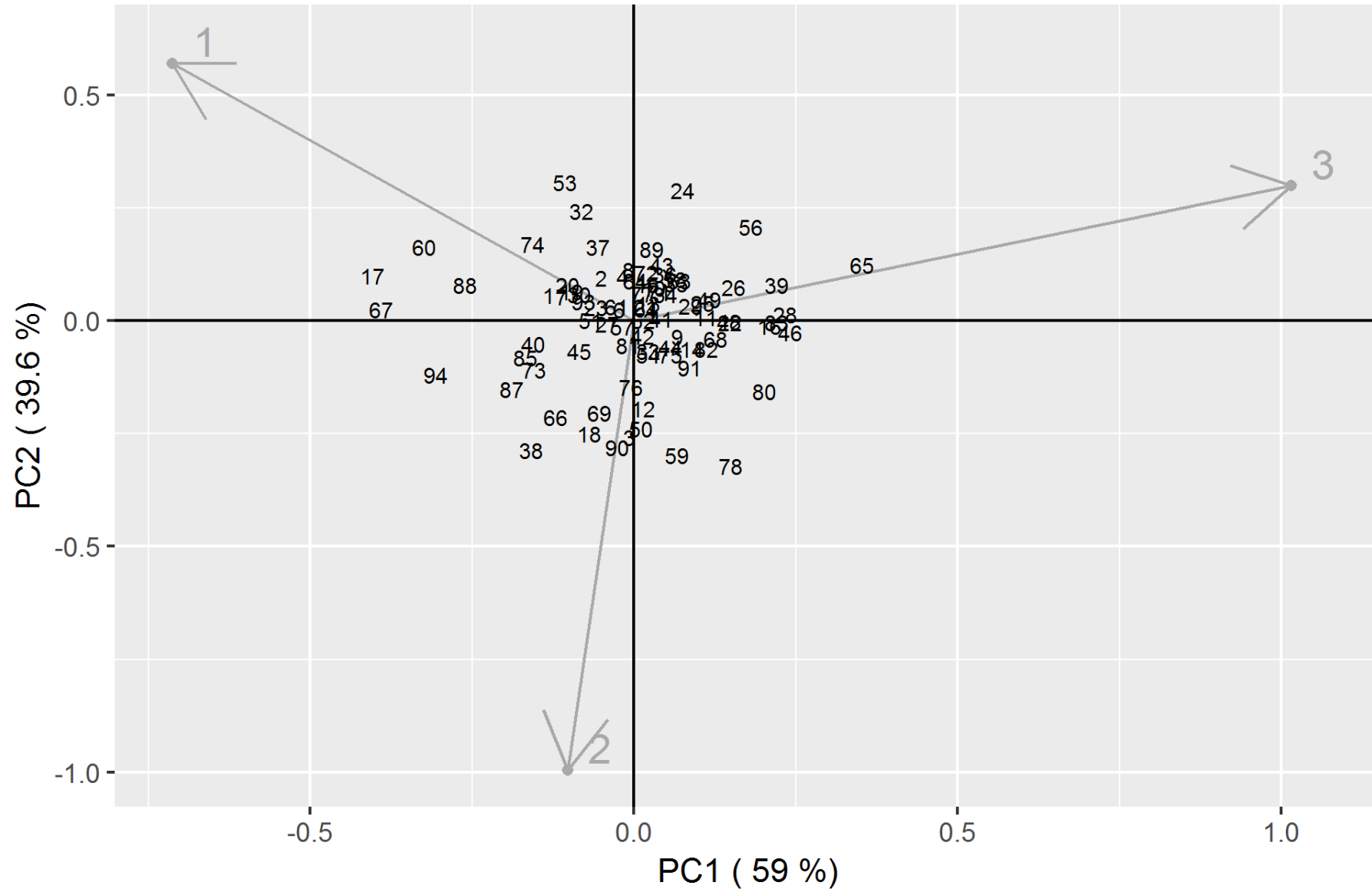
# Descrição do experimento



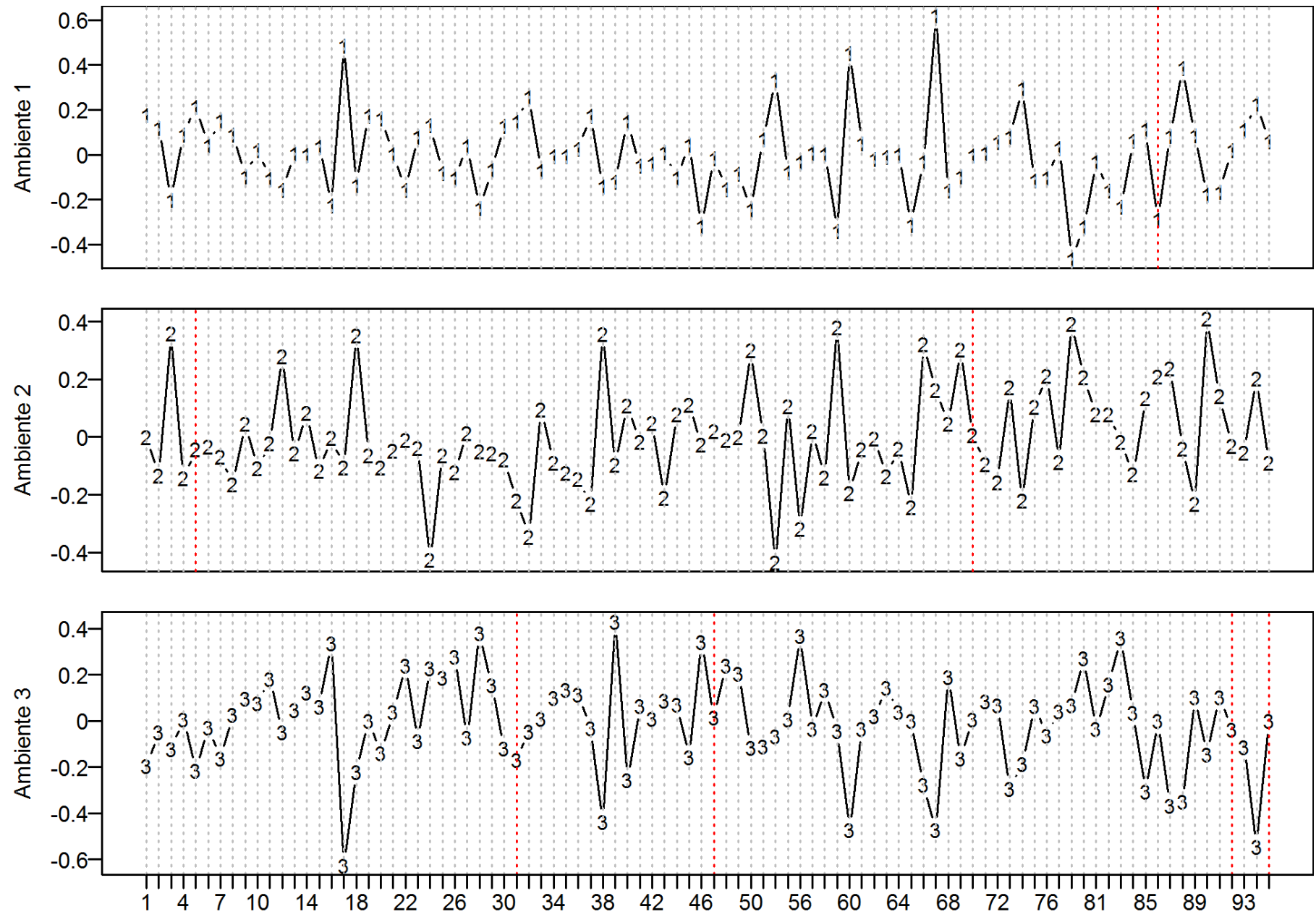
# Matriz da interação genótipo e ambiente com valores ausentes



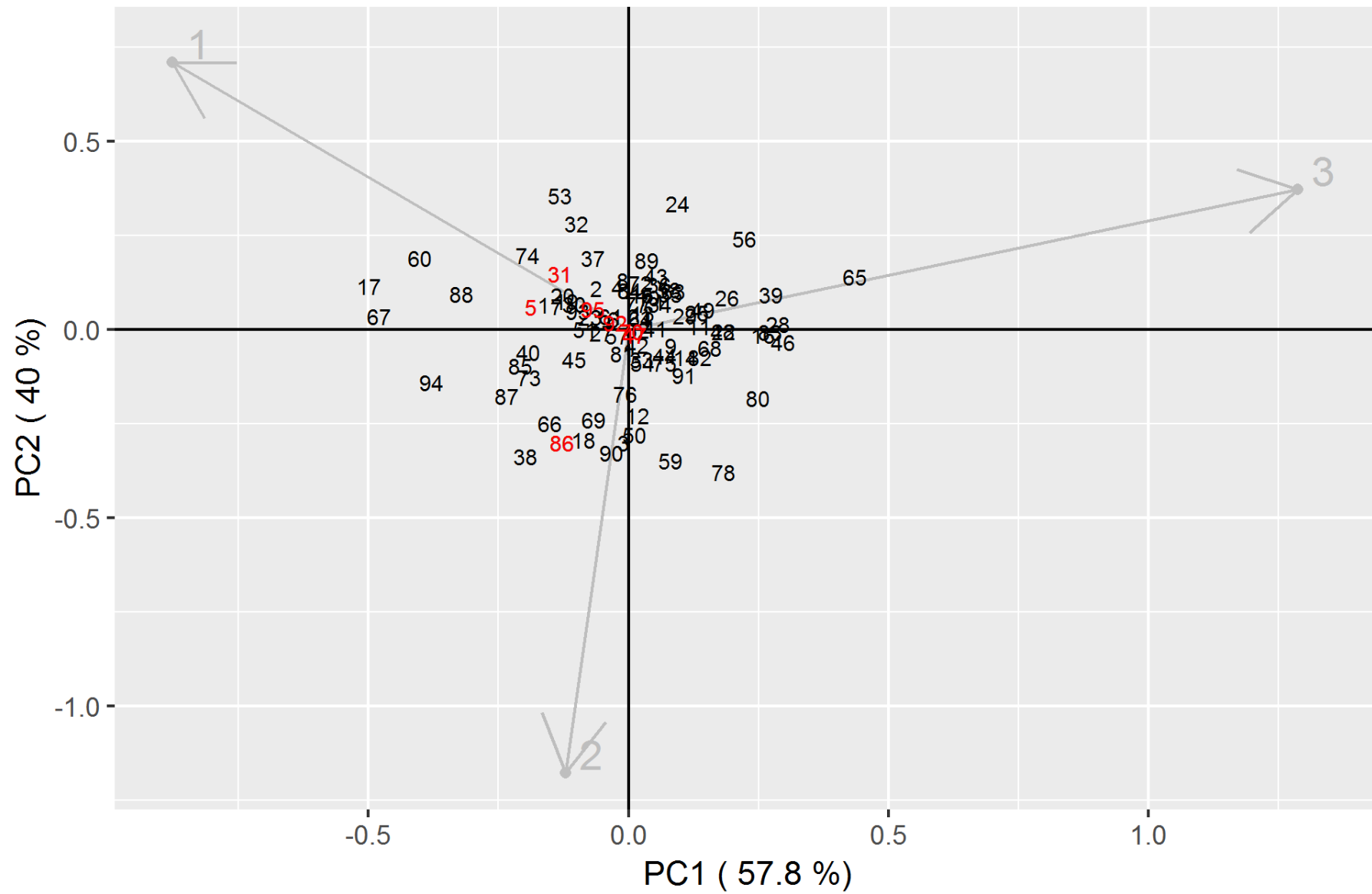
# AMMI sem imputação



# Matriz da interação genótipo e ambiente com valores imputados



# AMMI com imputação



# Estrutura da apresentação

---

1. Introdução
2. Problemática
3. Metodologia
4. Aplicação
5. Estudos Futuros

- 
1. Fazer testes em relação ao número de iterações, e sobre os coeficientes;
  2. Estudar sobre a teoria de imputação e encontrar um critério de teste para os valores imputados;
  3. Outros métodos de imputação;

**Obrigado pela atenção ^^**