

**Universidade de Brasília  
Instituto de Ciências Exatas  
Departamento de Estatística**

**Dissertação de Mestrado**

**Tópicos em análise de experimentos longitudinais  
para aplicações em estudos de sinais biopotenciais**

por

**Thaysa Guimarães Souza**

**Orientador: Prof. George Freitas von Borries, PhD**

**Julho de 2013**

Thaysa Guimarães Souza

# **Tópicos em análise de experimentos longitudinais para aplicações em estudos de sinais biopotenciais**

Dissertação apresentada ao Departamento de Estatística do Instituto de Ciências Exatas da Universidade de Brasília como requisito parcial à obtenção do título de Mestre em Estatística.

**Universidade de Brasília**

**Brasília, Julho de 2013**

TERMO DE APROVAÇÃO

Thaysa Guimarães Souza

**Tópicos em análise de experimentos longitudinais  
para aplicações em estudos de sinais biopotenciais**

Dissertação apresentada ao Departamento de Estatística do Instituto de Ciências Exatas da Universidade de Brasília como requisito parcial à obtenção do título de Mestre em Estatística.

Data da defesa: 11 de Julho de 2013

Orientador:

---

Prof. George Freitas von Borries, PhD  
Departamento de Estatística, UnB

Comissão Examinadora:

---

Prof. Dr. Afrânio Márcio Corrêa Vieira  
Departamento de Estatística, UnB

---

Prof<sup>a</sup>. Dr<sup>a</sup>. Roseli Aparecida Leandro  
Departamento de Ciências Exatas, ESALQ

**Brasília, Julho de 2013**

## Ficha Catalográfica

**SOUZA, THAYSA GUIMARÃES**

Tópicos em análise de experimentos longitudinais para aplicações em estudos de sinais biopotenciais, (UnB - IE, Mestre em Estatística, 2013).

Dissertação de Mestrado - Universidade de Brasília. Departamento de Estatística - Instituto de Ciências Exatas.

1. dados longitudinais 2. estrutura assintótica 3. teste de ausência de efeito simples 4. poder do teste 5. PPCLUSTEL 6. microarranjo 7. eletroencefalografia 8. eletromiografia.

É concedida à Universidade de Brasília a permissão para reproduzir cópias desta dissertação de mestrado e para emprestar ou vender tais cópias somente para propósitos acadêmicos e científicos. O autor reserva outros direitos de publicação e nenhuma parte desta dissertação de mestrado pode ser reproduzida sem a autorização por escrito do autor.

Thaysa Guimarães Souza

*A meus pais*

*A meu irmão*

# Agradecimentos

Agradeço primeiramente a Deus, que me permitiu trilhar toda essa estrada.

Aos meus pais, Maria Aparecida e Valtercides, pelo incansável apoio e por toda a educação e caráter que me proporcionaram.

Ao meu irmão, Pedro, por toda a ajuda e amizade.

Ao professor George von Borries, pelos conhecimentos compartilhados, pela orientação e pelo incentivo durante todo o transcorrer do curso de mestrado. Pela confiança e amizade em mim depositadas, sendo um exemplo de ética e seriedade no trabalho.

Aos funcionários do Departamento de Estatística da UnB, pelo excelente atendimento ao longo das minhas constantes idas ao departamento para verificar o andamento das simulações.

A todos os meus amigos que apoiaram, aturaram e participaram desta conquista.

E, por fim, agradeço à Coordenação de Aperfeiçoamento de Pessoal e Nível Superior (CAPES) pelo apoio financeiro ao longo de alguns semestres.

# Sumário

<b>Lista de Figuras</b>	<b>3</b>
<b>Lista de Tabelas</b>	<b>5</b>
<b>Resumo</b>	<b>6</b>
<b>Abstract</b>	<b>7</b>
<b>1 Introdução</b>	<b>8</b>
<b>2 Revisão de Literatura</b>	<b>12</b>
2.1 Análise Univariada . . . . .	13
2.2 Análise Multivariada . . . . .	17
2.3 Outras abordagens . . . . .	18
2.4 Limitações . . . . .	19
2.5 ADF e ADL . . . . .	20
<b>3 Casos Assintóticos</b>	<b>22</b>
3.1 Testes (Brunner e Puri, 2001) . . . . .	22
3.2 Testes (Wang, 2004) . . . . .	24
3.2.1 Teste baseados nas observações originais . . . . .	26
3.2.2 Testes baseados em postos . . . . .	27
3.3 Testes (von Borries, 2008) . . . . .	29
3.3.1 Testes baseados nas observações originais . . . . .	29
3.3.2 Testes baseados em postos . . . . .	31
3.4 Testes (Zhang, 2008) . . . . .	31
3.5 Algoritmo de agrupamento por partição . . . . .	34

<i>SUMÁRIO</i>	2
<b>4 Simulação e Aplicação de Algoritmos</b>	<b>37</b>
4.1 Resultados em dados simulados . . . . .	41
4.1.1 Comparação entre os algoritmos . . . . .	45
4.2 Resultados em dados reais . . . . .	47
4.2.1 Eletroencefalografia . . . . .	47
4.2.2 Microarranjo . . . . .	54
4.2.3 Eletromiografia . . . . .	57
<b>5 Conclusões</b>	<b>62</b>
5.1 Estudos futuros . . . . .	64
<b>A Poder do teste</b>	<b>70</b>
<b>B Heatmaps</b>	<b>83</b>
<b>C Programações em SAS</b>	<b>88</b>



# Lista de Figuras

4.1	Poder do teste. 50, 100 e 200 fatores com 5% de contaminação e 20 repetições por tratamento. . . . .	42
4.2	Poder do teste. 400 e 800 fatores com 5% de contaminação e 20 repetições por tratamento. . . . .	43
4.3	Poder do teste. 800 fatores e diferentes repetições por tratamento. . . . .	44
4.4	Poder do teste. 100 fatores e diferentes níveis de contaminação. . . . .	44
4.5	Poder do teste segundo diferentes metodologias. . . . .	45
4.6	Poder do teste segundo diferentes metodologias. . . . .	46
4.7	Coleta de dados de EEG no MSPL com o uso de estímulos visuais. . . . .	48
4.8	Estímulos visuais adotados na coleta dos sinais de EEG no MSPL. . . . .	49
4.9	PPCLUSTEL-RG com limiar $\alpha = 10^{-4}$ . . . . .	55
4.10	PPCLUSTEL-RW com limiar $\alpha = 10^{-4}$ . . . . .	56
4.11	PPCLUSTEL-RZ com limiar $\alpha = 10^{-8}$ . . . . .	56
4.12	Músculos envolvidos na análise de fadiga muscular. (a) <i>Sternocleidomastoid</i> - Ação: flexionar e rotacionar lateralmente a espinha cervical. (b) <i>Splenius capitis</i> - Ação: estender e rotacionar a espinha cervical. (c) <i>Trapezius</i> - Ação: estabilizar, elevar, retrain e rotacionar a escápula. (Imagens disponíveis em: < <a href="http://en.wikipedia.org/wiki/">http://en.wikipedia.org/wiki/</a> > Acesso em: 06/04/2012). . . . .	58
B.1	Dados não ordenados . . . . .	84
B.2	PPCLUSTEL-RG com limiar $\alpha = 10^{-4}$ . . . . .	85
B.3	PPCLUSTEL-RW com limiar $\alpha = 10^{-4}$ . . . . .	86
B.4	PPCLUSTEL-RZ com limiar $\alpha = 10^{-8}$ . . . . .	87

# Lista de Tabelas

1.1	Dados longitudinais - Estruturas. . . . .	9
2.1	Estrutura para Dados Longitudinais. . . . .	13
2.2	Análise de Variância para o Modelo de Parcelas Divididas com a Estrutura de Simetria Composta. . . . .	15
3.1	Postos e Distribuições das observações $Y_{ijk}$ . . . . .	23
4.1	Critério de decisão para o teste de efeito simples com base na metodologia de Zhang (2008). . . . .	38
4.2	Estrutura da base de dados de EEG. . . . .	50
4.3	Resultado do agrupamento com o algoritmo PPCLUSTEL-RW. Foi considerado o limiar $\alpha=0,1$ . . . . .	52
4.4	Resultado do agrupamento com o algoritmo PPCLUSTEL-RW. Foi considerado o limiar $\alpha=0,1$ . . . . .	53
4.5	Estrutura da base de dados de EMG. . . . .	59
4.6	P-valores dos testes de ausência de efeitos principais e de interação. . . . .	60
A.1	Poder do teste de ausência de efeito simples de Wang (2004) para 50, 100 e 200 fatores com 5% de contaminação. . . . .	71
A.2	Poder do teste de ausência de efeito simples de Wang (2004) para 400 e 800 fatores com 5% de contaminação. . . . .	72
A.3	Poder do teste de ausência de efeito simples de von Borries (2008) para 50, 100 e 200 fatores com 5% de contaminação. . . . .	73
A.4	Poder do teste de ausência de efeito simples de von Borries (2008) para 400 e 800 fatores com 5% de contaminação. . . . .	74

A.5 Poder do teste de ausência de efeito simples de Zhang (2008) para 50, 100 e 200 fatores com 5% de contaminação. . . . .	75
A.6 Poder do teste de ausência de efeito simples de Zhang (2008) para 400 e 800 fatores com 5% de contaminação. . . . .	76
A.7 Poder do teste de ausência de efeito simples de Wang (2004) para 50 e 100 fatores com 10% de contaminação. . . . .	77
A.8 Poder do teste de ausência de efeito simples de Wang (2004) para 200 e 400 fatores com 10% de contaminação. . . . .	78
A.9 Poder do teste de ausência de efeito simples de von Borries (2008) para 50 e 100 fatores com 10% de contaminação. . . . .	79
A.10 Poder do teste de ausência de efeito simples de von Borries (2008) para 200 e 400 fatores com 10% de contaminação. . . . .	80
A.11 Poder do teste de ausência de efeito simples de Zhang (2008) para 50 e 100 fatores com 10% de contaminação. . . . .	81
A.12 Poder do teste de ausência de efeito simples de Zhang (2008) para 200 e 400 fatores com 10% de contaminação. . . . .	82

# Resumo

Este trabalho busca revisar e comparar numericamente diferentes metodologias de análise de dados longitudinais com estrutura assintótica. Especificamente, são estudados testes de ausência de efeito simples com base nos trabalhos de Wang (2004), von Borries (2008) e Zhang (2008). As curvas de poder desses testes são construídas para diferentes cenários de simulação e a partir disso, constata-se que o teste de Zhang (2008) apresenta resultados superiores, mesmo nos casos em que os testes de von Borries (2008) e Wang (2004) eram tidos como adequados. Como consequência, esses testes são adaptados ao algoritmo de agrupamento PPCLUSTEL e utilizados na análise de dados de microarranjo, eletroencefalografia e eletromiografia. Os softwares SAS e Gnuplot são adotados na obtenção dos resultados.

**Palavras Chave:** *dados longitudinais com estrutura assintótica, teste de ausência de efeito simples, poder do teste, PPCLUSTEL, microarranjo, eletroencefalografia, eletromiografia.*

# Abstract

This work looks at different methodologies for analyzing longitudinal data with asymptotic structure. The studies focus specifically on tests of no simple effect based on the works of Wang (2004), von Borries (2008) and Zhang (2008). The power curves of these tests are then built for different simulation scenarios and from these curves it can be seen that Zhang's (2008) tests presents superior results, even in the cases where von Borries's (2008) and Wang's (2004) tests were considered adequate. Therefore, these tests are adapted to the clustering algorithms PPCLUSTEL and used in the analysis of microarray, electroencephalography and electromyography data. SAS and Gnuplot software were adopted for obtaining the results.

**Key words:** *longitudinal data with asymptotic distribution, tests of no simple effect, power curves, PPCLUSTEL, microarray, electroencephalography, electromyography.*

# Capítulo 1

## Introdução

Em áreas como medicina e biologia, não raramente, depara-se com a necessidade do tratamento de dados longitudinais. Basicamente, eles correspondem a medidas feitas em diferentes pontos no tempo para uma mesma unidade experimental e são caracterizados pela correlação usualmente existente entre os pares dessas medidas.

A literatura estatística contém material extenso sobre a análise desses dados. As abordagens incluem modelos mistos lineares e não-lineares com modelagens paramétrica (Pinheiro e Bates, 1995), semiparamétrica (Davidian e Gallant, 1993), não-paramétrica (Mallet, 1986) e bayesiana (Smith e Roberts, 1993). Além disso, modelos lineares generalizados são utilizados, por exemplo, no estudo de situações nas quais a variável resposta é ordinal (LIANG; ZEGUER, 1986).

Apesar das diferentes abordagens, a maior parte dos exemplos encontrados envolvem um número fixo de medidas repetidas no tempo, tratamentos e repetições por tratamento. As diferentes estruturas podem ser visualizadas na Tabela 1.1.

Um exemplo é análise de problemas que envolvem muitos pontos no tempo e tratamentos para um número fixo de repetições por tratamento. A análise desse caso é dificultada porque, segundo Bathke et al. (2011), não é possível obter estimadores consistentes das grandes matrizes de covariância. Especialmente, se essas matrizes precisam ser estimadas quando existem evidências empíricas da heterocedasticidade dos dados.

Como a estimação não pode ser realizada, testes de hipóteses de ausência de efeito de tratamento, tempo, interação entre ambos e efeito simples são afetados. De modo que, a fidedignidade dos resultados é comprometida.

Diante das dificuldades expostas, alguns pesquisadores têm estudado testes não-paramétricos alternativos na busca de soluções.

Tabela 1.1: Dados longitudinais - Estruturas.

Casos	Tratamentos	Repetição	Pontos no tempo	$a/n_i$
C1:	$a$ fixo	$n_i$ fixo	$b$ fixo	fixo
C2:	$a$ fixo	$n_i \rightarrow \infty$	$b$ fixo	$\rightarrow 0$
C3:	$a$ fixo	$n_i$ pequeno	$b \rightarrow \infty$	$\rightarrow \infty$
C4:	$a$ fixo	$n_i \rightarrow \infty$	$b \rightarrow \infty$	$\rightarrow 0$
C5:	$a \rightarrow \infty$	$n_i$ pequeno	$b$ fixo	$\rightarrow \infty$

Brunner e Puri (2001) trabalham com a avaliação de situações em que o número de tratamentos e pontos no tempo é fixo e o número de repetições por tratamento é fixo ou tende a infinito. Especificamente, os casos C1 e C2 descritos na Tabela 1.1.

Já Wang (2004) realiza inferências sobre os efeitos dos níveis de fator e pontos no tempo nos casos em que  $b \rightarrow \infty$ ,  $n_i$  pode ser pequeno ou grande e  $a$  é fixo. Inicialmente, os testes são aplicados sobre as observações originais e, posteriormente, sobre os postos dessas observações.

Outro caso é encontrado em von Borries (2008) e Zhang (2008), onde a estrutura de dados é caracterizada por  $a$  relativamente grande ou  $a \rightarrow \infty$ ,  $b$  fixo e  $n_i$  pequeno.

Sendo assim, o estudo dessas abordagens mostra-se de extrema relevância devido à inviabilidade de utilização da literatura tradicional. Um exemplo são os problemas encontrados na análise do caso C3, uma vez que a maior parte dos procedimentos tradicionais requer que o tamanho de amostra disponível seja superior ao número de pontos no tempo (AHMAD, 2008).

Desse modo, este trabalho apresenta como principais objetivos, a revisão bibliográfica da literatura recente de análise de dados longitudinais, a descrição de testes F em dados com estrutura não convencional e a investigação da aplicabilidade desses testes em diferentes situações.

Devido à subjetividade dos termos pequeno, fixo e infinito apresentados na Tabela 1.1, revela-se de fundamental importância verificar o desempenho dos testes de ausência de efeito simples apresentados por Wang (2004), von Borries (2008) e Zhang

(2008). Tal verificação é conduzida sob diferentes circunstâncias com o intuito de melhor compreender a dimensão dos termos utilizados pelos autores.

Tem-se como objetivo averiguar se um dos testes descritos como adequado para um determinado caso pode ou não ser utilizado em outras situações. E caso isso seja possível, observar se ele apresenta desempenho superior ou não ao que de fato foi desenvolvido para o caso de interesse.

A tabela abaixo ilustra a proposta, na qual as setas indicam a possibilidade de uma abordagem passar de uma situação para outra.

Tratamentos	Réplicas	Pontos no Tempo				
		pequeno	...	fixo	...	$\infty$
fixo	pequeno					C3
⋮	⋮	⋮	⋮	⋮	↗↘	⋮
fixo	fixo			C1		
⋮	⋮	⋮	⋮	↑↓	⋮	⋮
fixo	$\infty$			C2	↔	C4
⋮	⋮	⋮	⋮	↑↓	⋮	⋮
$\infty$	pequeno			C5		
⋮	⋮	⋮	⋮	↑↓	⋮	⋮
$\infty$	fixo			C5	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
$\infty$	$\infty$				⋮	⋮

Para avaliar a adequação dos testes de efeito simples de von Borries (2008), Wang (2004) e Zhang (2008), as curvas de poder desses testes foram construídas segundo diferentes cenários de simulação.

O critério de estimação da matriz de covariância foi realizado com base no trabalho de Silva (2012), no qual o autor analisa uma série de algoritmos de agrupamento de dados superdimensionados, longitudinais e com amostras pequenas. Os testes não-paramétricos presentes nos algoritmos estudados são simulados para três formas distintas de estimação da matriz de covariância. O processo de estimação  $\Sigma$  foi adotado neste trabalho.



É importante ressaltar que os testes que aqui serão estudados são válidos apenas para situações que envolvem efeitos fixos. E que apesar de poderem ser adotados no estudo de delineamentos não balanceados, serão estudadas apenas situações com estrutura balanceada ao longo do tempo e das repetições.

Destaca-se, também, que os testes de efeito simples serão utilizados na construção dos algoritmos de agrupamento PPCLUSTEL - RG, PPCLUSTEL - RW e PPCLUSTEL - RZ, baseados nas metodologias de von Borries (2008), Wang (2004) e Zhang (2008) com a estimação da matriz de covariância por meio da utilização de todo o conjunto de dados.

Como última observação, a estrutura deste trabalho está dividida da seguinte forma: no Capítulo 2, é realizada uma breve revisão da literatura de análise de dados longitudinais. São apresentadas as principais limitações de técnicas convencionais e as diferenças entre os termos funcional e longitudinal. No capítulo 3, são descritas diferentes abordagens de análise de dados longitudinais com estrutura assintótica. E por fim, no capítulo 4 são discutidos os resultados das simulações realizadas e a aplicação dos testes na análise de dados de eletromiografia, eletroencefalografia e microarranjo.

# Capítulo 2

## Revisão de Literatura

Estudos longitudinais envolvem observações de um conjunto de unidades experimentais classificadas segundo diferentes tratamentos e pontos no tempo. Basicamente, esses pontos nos quais as medidas são tomadas podem ser representados por  $t_1, \dots, t_b$ , o número total de tratamentos pode variar de  $i = 1, \dots, a$  e o número de repetições por tratamento de  $k = 1, \dots, n_i$ . A estrutura é apresentada na tabela 2.1.

De uma forma geral, a análise de dados longitudinais (ADL) permite investigar mudanças de comportamento de uma mesma unidade experimental ou variações entre unidades ao longo do tempo. Tal que, é possível descrever um perfil individual de respostas associado a cada unidade e um perfil médio observado em cada tratamento.

É importante notar que o uso de observações repetidas no tempo costuma ser mais eficiente que a adoção de uma única unidade experimental para cada ponto  $t_1, \dots, t_b$ . Além de exigir uma menor quantidade de unidades, promove uma redução nos custos e um aumento na precisão da estimação da tendência da variável resposta ao longo do tempo (KUEHL, 2001).

Apesar das vantagens mencionadas, a análise de dados longitudinais é um pouco mais complexa que a de dados transversais. Isso porque é preciso levar em consideração os efeitos do tempo sobre as respostas da unidade experimental e a correlação existente entre as respostas de uma mesma unidade.

Milliken e Johnson (2009) mostram que existem diferentes formas de analisar tais dados. Na análise univariada, o tempo é um dos fatores do modelo ANOVA e existem restrições rigorosas quanto à estrutura de covariância dos dados. Já no método multivariado, adota-se uma matriz de covariância não estruturada, ou seja,

o modelo MANOVA não exige que a variância das medidas repetidas e a covariância dos pares dessas medidas permaneçam constantes ao longo do tempo.

Destaca-se ainda que, em geral, os métodos apresentados são adotados na avaliação de situações experimentais tradicionais, ou seja, com um número fixo de tratamentos, repetições e pontos no tempo. Algumas referências comumente adotadas na análise tradicional são Crowder e Hand (1990), Milliken e Johnson (2009) e Hedecker e Gibbons (2006). Modelos mais complexos são encontrados nos trabalhos de Davis (2002), Diggle et al. (2002), Vonesh e Chinchilli (1996) e Davidian e Giltinan (1995).

Tabela 2.1: Estrutura para Dados Longitudinais.

Tratamentos	Indivíduos	Pontos no Tempo			
		$t_1$	$t_2$	$\dots$	$t_b$
1	1	$Y_{111}$	$Y_{121}$	$\dots$	$Y_{1b1}$
	2	$Y_{112}$	$Y_{122}$	$\dots$	$Y_{1b2}$
	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
	$n_1$	$Y_{11n_1}$	$Y_{12n_1}$	$\dots$	$Y_{2bn_1}$
2	1	$Y_{211}$	$Y_{221}$	$\dots$	$Y_{2b1}$
	2	$Y_{212}$	$Y_{222}$	$\dots$	$Y_{2b2}$
	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
	$n_2$	$Y_{21n_2}$	$Y_{22n_2}$	$\dots$	$Y_{2bn_2}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$a$	1	$Y_{a11}$	$Y_{a21}$	$\dots$	$Y_{ab1}$
	2	$Y_{a12}$	$Y_{a22}$	$\dots$	$Y_{ab2}$
	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
	$n_a$	$Y_{a1n_a}$	$Y_{a2n_a}$	$\dots$	$Y_{abn_a}$

## 2.1 Análise Univariada

Adotar técnicas de parcelas divididas (*split-plot designs*) na análise de dados longitudinais é uma prática comum e que exige cuidado. Para que isso possa ser feito, é necessário considerar que os termos de erro apresentam distribuição normal e que são

independentes e identicamente distribuídos. O modelo abaixo pode ser adotado ao longo do processo de análise.

$$Y_{ijk} = \mu + \alpha_i + d_{ik} + \tau_j + (\alpha\tau)_{ij} + \epsilon_{ijk} \quad (2.1)$$

$$i=1,\dots,a \quad j=1,\dots,b \quad k=1,\dots,n_i$$

$\mu$  : média geral;

$\alpha_i$ : efeito do  $i$ -ésimo tratamento;

$d_{ik}$ : erro experimental dentro do tratamento (erro de parcela), tal que  $d_{ik} \stackrel{i.i.d.}{\sim} N(0, \sigma_d^2)$ ;

$\tau_j$ : efeito do  $j$ -ésimo ponto no tempo;

$(\alpha\tau)_{ij}$ : efeito da interação entre tempo e tratamento;

$\epsilon_{ijk}$ : erro aleatório (erro de subparcela), tal que  $\epsilon_{ijk} \stackrel{i.i.d.}{\sim} N(0, \sigma_e^2)$ ;

$\epsilon_{ijk}$  e  $d_{ik}$  independentes.

Além das suposições mencionadas, assume-se que os dados são balanceados ao longo do tempo e que podem não ser em relação às repetições por tratamento.

De acordo com Xavier (2000), em um delineamento de parcelas divididas com medidas repetidas no tempo, o teste F com relação à parcela (tratamento) tem distribuição F exata, mas em relação à subparcela (tempo e interação entre tempo e tratamento), só terá essa distribuição se a matriz de covariância satisfizer certa condição. Dessa forma, uma condição suficiente para garantir a validade do modelo e teste mencionados é que a matriz de covariância dos dados tenha forma de simetria composta, ou seja, apresente correlações constantes para qualquer par de diferentes instantes de tempo.

A tabela 2.2 apresenta a análise de variância (ANOVA) para a estrutura considerada. Nessa tabela,  $\theta_A$ ,  $\theta_T$  e  $\theta_{TA}$  são parâmetros de não centralidade que medem os efeitos de tratamento, tempo e interação entre esses fatores, respectivamente. De modo que, tais parâmetros são iguais a zero se os correspondentes efeitos são nulos. Além disso,  $\bar{y}$  e o ponto subscrito indicam a média e as unidades das quais ela é tomada.

O primeiro teste de interesse deve ser o da interação. Assim, há evidências para rejeitar  $H_0(\gamma)$ :  $\theta_{TA}=0$  se

$$F = \frac{MSTA}{MSE} > F_{\alpha, (a-1)(b-1), a(b-1)(n-1)} \quad (2.2)$$

Tabela 2.2: Análise de Variância para o Modelo de Parcelas Divididas com a Estrutura de Simetria Composta.

Fontes de Variação	Graus de Liberdade (GL)	Somas de Quadrados (SQ)	Quadrados Médios (QM)	QMs Esperados
Tratamento (A)	a-1	$SSA = b \sum_{i=1}^a n_i (\bar{y}_{i..} - \bar{y}_{...})^2$	$MSA = SSA / a - 1$	$\sigma_e^2 + b\sigma_d^2 + \theta_A$
Erro inter-individual	N-a	$SSD = b \sum_{i=1}^a \sum_{k=1}^{n_i} (\bar{y}_{i.k} - \bar{y}_{i..})^2 - SSA$	$MSD = SSD / N - a$	$\sigma_e^2 + b\sigma_d^2$
Tempo (T)	b-1	$SST = N \sum_{j=1}^b (\bar{y}_{.j.} - \bar{y}_{...})^2$ , onde $N = \sum_{i=1}^a n_i$	$MST = SST / b - 1$	$\sigma_e^2 + \theta_T$
Tempo x Tratamento (TA)	(a-1)(b-1)	$SSTA = \sum_{i=1}^a \sum_{j=1}^b n_i (\bar{y}_{i.j.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2$	$MSTA = SSTA / (a-1)(b-1)$	$\sigma_e^2 + \theta_{TA}$
Erro intra-individual	(N-a)(b-1)	$SSE = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} (\bar{y}_{i.j.k} - \bar{y}_{i.k} - \bar{y}_{i.j.} + \bar{y}_{i..})^2$	$MSE = SSE / (N-a)(b-1)$	$\sigma_e^2$
Total	Nb-1	$SSTo = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} (\bar{y}_{i.j.k} - \bar{y}_{...})^2$		

Quando a hipótese nula de ausência do efeito de interação é rejeitada, as diferenças entre os tratamentos não são as mesmas ao longo do tempo, os efeitos de tempo e tratamento são confundidos com a interação e não podem ser separadamente testados.

Para o caso em que hipótese em questão não é rejeitada, os testes para os efeitos principais, tempo e tratamento, são válidos. Assim, há evidências para rejeitar  $H_0(\alpha)$ :  $\theta_A=0$  quando

$$F = \frac{MSA}{MSD} > F_{\alpha, a-1, a(n-1)} \quad (2.3)$$

E por fim, o mesmo ocorre com  $H_0(\beta)$ :  $\theta_T=0$  se

$$F = \frac{MST}{MSE} > F_{\alpha, b-1, a(b-1)(n-1)} \quad (2.4)$$

Os estimadores das componentes de variância do modelo considerado são apresentados nas equações abaixo:

$$\hat{\sigma}_e^2 = MSE$$

$$\hat{\sigma}_d^2 = \frac{MSD - \hat{\sigma}_e^2}{b}$$

Segundo Hedecker e Gibbons (2006), a simetria composta é altamente restritiva e freqüentemente irrealista (especialmente à medida que  $b$  aumenta). Essa estrutura de covariância é um caso especial de uma situação mais geral denominada esfericidade. Quando a suposição de esfericidade não é rejeitada, os testes F mencionados são válidos. Caso contrário, ocorre uma inflação na probabilidade de erro do tipo I associada ao teste F de efeito de subparcela.

O teste de esfericidade de Mauchly (1940) é uma alternativa existente para verificar a validade da esfericidade. No entanto, é altamente sensível à violação de normalidade dos dados e apresenta um baixo poder nas situações que envolvem um número pequeno de observações.

Quando a suposição de esfericidade é rejeitada, o pesquisador pode trabalhar com a análise multivariada de medidas repetidas (MANOVA). Outra alternativa estatística

clássica inclui o ajuste do número de graus de liberdade do teste F proposto por Box (1945) e aperfeiçoado por Greenhouse e Geisser (1959) e por Huynh e Feldt (1976).

## 2.2 Análise Multivariada

O modelo adotado na Análise de Variância Multivariada (MANOVA) pode ser representado conforme indicado na equação 2.5. Especificamente,  $\mathbf{Y}_{(N \times p)}$  denota o conjunto de dados, no qual cada linha caracteriza uma unidade experimental e cada coluna corresponde a uma das medidas repetidas.  $\mathbf{X}_{(N \times r)}$  representa a matriz de delineamento com posto  $t$ .  $\boldsymbol{\beta}_{(r \times 1)}$  denota a matriz de parâmetros. E, por fim,  $\boldsymbol{\epsilon}_{(N \times p)}$  corresponde à matriz de erros aleatórios.

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (2.5)$$

É importante notar que embora as linhas de  $\boldsymbol{\epsilon}$  sejam independentes com distribuição  $N_p(\mathbf{0}, \boldsymbol{\Sigma})$ , os elementos de uma mesma linha podem ser correlacionados e apresentar variâncias heterogêneas. Além disso, uma matriz de covariância não estruturada é adotada ao longo do processo. Para que métodos MANOVA possam ser aplicados, o conjunto de dados deve ser balanceado no tempo e o número de unidades experimentais alocadas a cada tratamento não precisa ser o mesmo.

As hipóteses de perfis coincidentes, constantes e paralelos, abordadas em Johnson e Wichern (2007), podem ser representadas na forma da hipótese linear geral apresentada abaixo:

$$H_0 : \mathbf{CBM} = 0 \text{ contra } H_a : \mathbf{CBM} \neq 0 \quad (2.6)$$

onde  $\mathbf{C}_{(g \times r)}$  e  $\mathbf{M}_{(p \times q)}$  são matrizes conhecidas com postos  $g$  e  $q$ , respectivamente.

Para testar a hipótese de interesse é necessário obter os estimadores de mínimos quadrados dos parâmetros da matriz  $\boldsymbol{\beta}_{(r \times 1)}$ . De modo que, “ $\hat{\boldsymbol{\epsilon}}$  corresponde à matriz de somas de quadrados e produtos cruzados devido ao erro” (MILLIKEN; JOHNSON,

2009) .

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \text{ e } \hat{\epsilon} = \mathbf{Y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{Y} \quad (2.7)$$

O critério da razão de verossimilhança de Wilks ( $\Lambda$ ) (equação 2.8) adotado como estatística do teste em questão apresenta distribuição amostral complexa.

$$\Lambda = \frac{|\mathbf{R}|}{|\mathbf{H} + \mathbf{R}|} \quad (2.8)$$

Uma alternativa é trabalhar com um teste aproximado que rejeita  $H_0$  quando a condição estabelecida na equação 2.9 é satisfeita.

$$\begin{aligned} \mathbf{R} = \mathbf{M}'\hat{\epsilon}\mathbf{M} \quad \mathbf{H} = \mathbf{M}'\hat{\beta}\mathbf{C}'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}\mathbf{C}\hat{\beta}\mathbf{M} \\ - \left( N - t - \frac{|q - g| + 1}{2} \right) \log_e(\Lambda) > \chi_{\alpha, qg}^2 \end{aligned} \quad (2.9)$$

Destaca-se ainda que existem melhores aproximações para  $\Lambda$  quando  $g$  e  $q$  são maiores que 2. Conforme Rencher (2002), testes F exatos ocorrem quando  $g$  ou  $q$  apresentam valores 1 ou 2. Ressalta-se, “que outros testes multivariados podem ser adotados e encontrados em pacotes estatísticos” (JOHNSON; WICHERN, 2007).

## 2.3 Outras abordagens

Grande parte dos esforços adotados na ADL está relacionada com a modelagem da estrutura de covariância decorrente da medição de uma mesma variável, na mesma unidade experimental, em diferentes instantes de tempo.

Com esse objetivo, Laird e Ware (1982) propõem a utilização de modelos lineares mistos. Já Liang e Zeger (1986) apresentam uma extensão de modelos lineares generalizados para a análise dos dados em questão. O trabalho apresentado em Fitzmaurice et al. (2009) expõe com uma maior riqueza de detalhes as particularidades dessas técnicas.

Outro modo de analisar dados longitudinais é através da análise de curvas de crescimento, por meio de modelos mistos lineares ou não-lineares, que possibilitam o



uso de diferentes estruturas de covariância, de modo a descrever o comportamento dos perfis médios através de curvas.

Portanto, existem inúmeros métodos de análise e eles variam de acordo com as características do experimento em estudo. Por exemplo, deve ser considerado se a variável resposta é quantitativa ou não, se as suposições de normalidade e homocedasticidade dos dados são atendidas e se a estrutura dos dados de fato se enquadra na realização de análises tradicionais. Caso contrário, é preciso buscar técnicas na literatura que forneçam maior suporte.

## 2.4 Limitações

Apesar das inúmeras aplicações, as técnicas mencionadas não se adequam ao estudo de casos assintóticos (tabela 1.1). Como exemplo, tem-se o problema encontrado na especificação da estatística do teste de efeito do fator tempo ( $F = MST/MSE$ ) quando  $b \rightarrow \infty$ , já que o número de graus de liberdade das fontes tempo e erro também passa a apresentar o mesmo comportamento. Situação semelhante ocorre para o teste de ausência de efeito de tratamento quando  $a \rightarrow \infty$ .

Outras limitações surgem quando os dados são superdimensionados, ou seja, com um número de pontos no tempo bem maior que o tamanho de amostra disponível. Segundo Ahmad (2008), no caso univariado, a suposição de esfericidade dificilmente é satisfeita. E, no multivariado, a estimação da matriz de covariância é comprometida, já que ela deixa de ser não-singular.

Assim, quando o conjunto de dados não satisfaz as suposições usuais impostas pelas técnicas tradicionais, outras abordagens devem ser consideradas. Trabalhos recentes têm buscado soluções para analisar tais casos.

Conforme mencionado anteriormente, Wang (2004) considera um modelo não-paramétrico para realizar inferências sobre os efeitos dos níveis de fator e pontos no tempo nas situações em que  $b \rightarrow \infty$ ,  $n_i$  é pequeno ou grande e  $a$  é fixo. Os testes são aplicados sobre as observações originais e sobre os postos dessas observações.

Outro caso é encontrado em von Borries (2008) e Zhang (2008), onde a estrutura de dados é caracterizada por  $a$  relativamente grande ou  $a \rightarrow \infty$ ,  $b$  fixo e  $n_i$  pequeno. De modo que, na primeira abordagem, a hipótese nula e a estatística do teste de

efeito simples são definidas para dados originais. E na segunda, o autor trabalha com testes de ausência de efeito de tratamento e interação entre tempo e tratamento.

Estudos com situações assintóticas podem ainda ser vistos na literatura de análise de dados funcionais (ADF). Um exemplo é a situação estudada por Fan e Lin (1998), na qual os autores propõem novos testes para comparar um conjunto de curvas de dados comerciais. Essas curvas são observadas ao longo do tempo e analisadas por meio de técnicas envolvendo *wavelets* e o teste adaptado de Neyman.

## 2.5 ADF e ADL

No estudo de situações que envolvem um número relativamente elevado de pontos no tempo, comumente se encontra o termo funcional como sinônimo de longitudinal. Contudo, ainda que ambos sejam aplicados a estruturas com medidas repetidas em diferentes instantes, alguns autores enxergam divergências entre as duas abordagens.

De acordo com Souza (2008), ao serem registrados densamente ao longo do tempo, na maior parte das vezes por máquinas, os dados são denominados funcionais. Ao contrário dos longitudinais, que são captados de forma mais esparsa.

Fan e Lin (1998) apresentam idéias semelhantes. Para os autores, a diferença entre tais dados reside no fato de que vetores de dados funcionais apresentam dimensionalidade consideravelmente superior e que, portanto, necessitam de técnicas de redução de dimensionalidade.

Assim, a análise de dados funcionais (ADF) se difere por utilizar como dados individuais funções das medidas repetidas. Nesse caso, o interesse não é restrito apenas ao estudo das curvas formadas, mas também das derivadas e integrais dessas curvas. Considere o seguinte exemplo: no estudo de crescimento de indivíduos pode-se estar interessado em, além de estimar a curva de crescimento, simultaneamente estimar a velocidade de crescimento e a aceleração como função do tempo para cada indivíduo.

Desse modo, assume-se a existência de uma função  $x(t)$  baseada nos dados observados, a qual implica, a princípio, dados funcionais continuamente definidos. Entretanto, caso apenas valores discretos estejam disponíveis, é possível avaliar  $x(t)$  por meio do uso de técnicas de suavização.

Ademais, as funções podem ser independentes umas das outras, de modo que a mesma suposição não é válida para valores de uma mesma função. Maiores detalhes sobre ADF são encontrados em (Ramsay e Silverman, 2002) e (Rice, 2004).

Nesse contexto, Faraway (1997) sugeriu a redução da dimensionalidade de dados funcionais para que posteriormente pudessem ser aplicadas técnicas tradicionais de ADL. Entretanto, o objetivo deste trabalho não está em reduzir a dimensionalidade e sim buscar técnicas de ADL que possam ser diretamente aplicadas a dados funcionais.

Apesar das diferenças conceituais apresentadas, as análises também apresentam pontos em comum. Segundo Rice (2004), destacam-se os seguintes: o estudo de fontes importantes de padrão e variação entre as observações, a análise da variação de resultados com base nas informações da variável independente e a comparação de dois ou mais fatores com respeito a certos tipos de variação.

# Capítulo 3

## Casos Assintóticos

Quando um conjunto de dados não satisfaz as suposições usuais impostas pelas técnicas tradicionais, métodos não-paramétricos surgem como uma boa opção para a análise de dados.

Recentemente, diferentes testes não-paramétricos foram propostos para lidar com a análise dos casos assintóticos apresentados na tabela 1.1. Os testes adotados encontram-se devidamente especificados na seqüência e são válidos apenas para a análise de situações que envolvem efeitos fixos.

### 3.1 Testes (Brunner e Puri, 2001)

Brunner e Puri (2001) consideram um modelo não-paramétrico para realizar inferências no caso em que o número de tratamentos e pontos no tempo é fixo e o número de repetições por tratamento é moderado ou tende a infinito.

Os autores representam cada indivíduo aninhado a um tratamento por vetores de observações  $\mathbf{Y}_{ik} = (Y_{i1k}, \dots, Y_{ibk})'$  independentes, onde  $i = 1, \dots, a$  e  $k = 1, \dots, n_i$ . De modo que,  $Y_{ijk}$  segue uma distribuição arbitrária  $F_{ij}$ .

As hipóteses adotadas são puramente não-paramétricas. São construídas com base na decomposição de distribuições  $F_{ij}$  proposta por Akritas e Arnold (1994), ou seja,  $F_{ij}(x) = M(x) + A_i(x) + B_j(x) + C_{ij}(x)$ , onde  $\sum_i A_i(x) = \sum_j B_j(x) = \sum_i C_{ij}(x) = \sum_j C_{ij}(x) = 0$ .

De modo que,  $M(x) = \bar{F}_{..}(x)$ ,  $A_i(x) = \bar{F}_{i.}(x) - \bar{F}_{..}(x)$ ,  $B_j(x) = \bar{F}_{.j}(x) - \bar{F}_{..}(x)$  e  $C_{ij}(x) = \bar{F}_{ij}(x) - \bar{F}_{i.}(x) - \bar{F}_{.j}(x) + \bar{F}_{..}(x)$ . Os postos e as distribuições das observações

$Y_{ijk}$  são apresentados na tabela 3.1.

Tabela 3.1: Postos e Distribuições das observações  $Y_{ijk}$ .

Tratamentos	Indivíduos	Postos				Distribuições			
		$t_1$	$t_2$	$\dots$	$t_b$	$t_1$	$t_2$	$\dots$	$t_b$
1	1	$R_{111}$	$R_{121}$	$\dots$	$R_{1b1}$	$F_{11}$	$F_{12}$	$\dots$	$F_{1b}$
	2	$R_{112}$	$R_{122}$	$\dots$	$R_{1b2}$	$F_{11}$	$F_{12}$	$\dots$	$F_{1b}$
	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
	$n_1$	$R_{11n_1}$	$R_{12n_1}$	$\dots$	$R_{1bn_1}$	$F_{11}$	$F_{12}$	$\dots$	$F_{1b}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	
a	1	$R_{a11}$	$R_{a21}$	$\dots$	$R_{ab1}$	$F_{a1}$	$F_{a2}$	$\dots$	$F_{ab}$
	2	$R_{a12}$	$R_{a22}$	$\dots$	$R_{ab2}$	$F_{a1}$	$F_{a2}$	$\dots$	$F_{ab}$
	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
	$n_a$	$R_{a1n_a}$	$R_{a2n_a}$	$\dots$	$R_{abn_a}$	$F_{a1}$	$F_{a2}$	$\dots$	$F_{ab}$

Segundo Wang (2004), “a idéia principal das hipóteses não-paramétricas é descrever os efeitos envolvidos por meio da decomposição de funções de distribuição”<sup>1</sup>. Assim, as hipóteses de interesse consistem nos seguintes aspectos:

1. Avaliar a ausência de efeito de tratamento (fator A), ou seja, verificar se as médias  $\bar{F}_i = \frac{1}{b} \sum_{j=1}^b F_{ij}$  das distribuições  $F_{ij}$  ao longo dos  $b$  pontos no tempo são as mesmas para todos os grupos  $i = 1, \dots, a$ . Basicamente, a hipótese de interesse pode ser descrita por  $H_0(A)$ : todo  $A_i = 0$ ;
2. Verificar se a resposta dos  $b$  pontos no tempo são homogêneas entre os grupos de tratamento, ou seja, testar a hipótese de ausência do efeito da interação entre tempo e tratamento,  $H_0(C)$ : todo  $C_{ij} = 0$ .

Para testar as hipóteses apresentadas, Brunner e Puri (2001) trabalham com versões baseadas em postos das formas quadráticas WTS (*Wald-type Statistics*) e ATS (*ANOVA-type Statistics*) e com uma matriz de covariâncias não estruturada. Uma breve descrição<sup>2</sup> dessas estatísticas é apresentada na seqüência:

<sup>1</sup>De acordo com Wang (2004), “as hipóteses usuais baseadas na decomposição de médias podem ser deduzidas das hipóteses não-paramétricas”.

<sup>2</sup>Os resultados assintóticos encontram-se devidamente apresentados em Brunner e Puri (2001).

1. WTS: sob  $H_0(C)$ , a estatística WTS ( $Q_n(C)$ ) apresenta distribuição assintótica  $\chi^2_{(a-1)(b-1)}$ . Já sob  $H_0(A)$ , a estatística em questão converge para a distribuição  $\chi^2_{(a-1)}$ .
2. ATS: sob  $H_0(C)$  e  $H_0(A)$ , as estatísticas ATS ( $F_n(C)$  e  $F_n(A)$ ) apresentam distribuição assintótica F com um número aproximado de graus de liberdade para o numerador ( $\hat{f}_1$ ) e denominador ( $\hat{f}_0$ ). As aproximações para os graus de liberdade são encontradas em Brunner e Puri (2001, p.41 e 43).

A estatística WTS apresenta como vantagem o fato de sua distribuição assintótica sob  $H_0$  ser uma função conhecida. Além disso, essa estatística converge extremamente devagar para a distribuição assintótica, o que resulta em decisões liberais para tamanhos de amostra pequenos ou moderados.

Já a ATS apresenta a desvantagem de sua distribuição assintótica sob  $H_0$  possuir quantidades desconhecidas, as quais são obtidas pela aproximação de Box (1954). A principal vantagem dessa estatística é que a aproximação para a distribuição assintótica funciona bem para amostras pequenas, moderadas e grandes.

### 3.2 Testes (Wang, 2004)

Wang (2004) utiliza um modelo não-paramétrico para analisar dados longitudinais heterocedásticos com  $a$  fixo,  $n_i$  pequeno ou grande e  $b \rightarrow \infty$ . A autora representa cada indivíduo aninhado a um nível de fator por uma série temporal  $Y_{ik} = (Y_{i1k}, \dots, Y_{ibk})'$  com  $Y_{ijk} = \mu + \alpha_i + d_{jk(i)} + \beta_j + \gamma_{ij} + \epsilon_{ijk}$ , onde  $i = 1, \dots, a$ ;  $j = 1, \dots, b$  e  $k = 1, \dots, n_i$ . As suposições abaixo são consideradas:

- $\sum_{i=1}^a \alpha_i = \sum_{j=1}^b \beta_j = \sum_{i=1}^a \gamma_{ij} = \sum_{j=1}^b \gamma_{ij} = 0$ ;
- $E(\epsilon_{ijk}) = 0$ ;
- $Y_{ijk} \sim F_{ij}$  arbitrária;
- $(Y_{ijk}, Y_{i'j'k'}) \sim F_{ij}F_{i'j'}$  para  $k \neq k'$  ou  $i \neq i'$ , o que caracteriza a independência entre séries temporais provenientes de indivíduos distintos;

- As séries de tempo satisfazem a condição  $\alpha$ -mixing. Basicamente, essa condição implica que a correlação existente entre as observações de um mesmo indivíduo em determinado nível de fator decresce à medida que a defasagem de tempo  $m$  aumenta. De acordo com essa condição, para alguma seqüência  $\alpha_m \rightarrow 0$ ,

$$|P(A \cap B) - P(A)P(B)| \leq \alpha_m \quad (3.1)$$

vale para todo  $A \in \sigma(Y_{i1k}, \dots, Y_{ilk})$ ,  $B \in \sigma(Y_{i,l+m,k}, Y_{i,l+m+1,k}, \dots)$  e todo  $i, k$ ; de modo que,  $\sigma(\cdot)$  caracteriza o espaço  $\sigma$ -álgebra gerado pela variáveis.

É importante notar que  $\alpha_m \rightarrow 0$  corresponde a  $Y_{ilk}$  e  $Y_{i,l+m,k}$  aproximadamente independentes para  $m$  grande.

As seguintes notações são utilizadas na obtenção dos resultados de interesse:

$$N = b \sum_{i=1}^a n_i,$$

$$\tilde{Y}_{i..} = \bar{Y}_{i..} = \frac{1}{bn_i} \sum_{j=1}^b \sum_{k=1}^{n_i} Y_{ijk},$$

$$\tilde{Y}_{.j.} = \frac{1}{a} \sum_{i=1}^a \frac{1}{n_i} \sum_{k=1}^{n_i} Y_{ijk},$$

$$\bar{Y}_{.j.} = \frac{1}{n} \sum_{i=1}^a \sum_{k=1}^{n_i} Y_{ijk},$$

$$\bar{Y}_{...} = \frac{1}{N} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} Y_{ijk} \text{ e}$$

$$\tilde{Y}_{...} = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b \frac{1}{n_i} \sum_{k=1}^{n_i} Y_{ijk}.$$

Para facilitar o entendimento, é importante ter em mente alguns conceitos. Por exemplo, a definição de efeito de fator como o conjunto de alterações que ocorrem na variável resposta como consequência de mudanças ocasionadas nos níveis desse fator.

Assim, apesar das suposições apresentadas acima caracterizarem uma situação inteiramente não-paramétrica, a representação da ausência de certos efeitos é conduzida de forma usual através da utilização das hipóteses apresentadas abaixo:

$$H_0(\alpha) : \quad \text{todo } \alpha_i = 0 \quad (\text{ausência de efeito do tratamento}),$$

$$H_0(\gamma) : \quad \text{todo } \gamma_{ij} = 0 \quad (\text{ausência de efeito da interação tempo x tratamento}),$$

$$H_0(\phi) : \quad \phi_{ij} = \alpha_i + \gamma_{ij} = 0, \quad \text{para todo } i \text{ e } j \quad (\text{ausência de efeito simples}).$$

É importante ressaltar que o teste de ausência de efeito simples é adotado para averiguar a homogeneidade de distribuições ao longo do tempo.

Wang et al. (2008) utiliza esse teste na comparação de grupos de curvas, ou seja, na comparação de diferentes tratamentos ao longo do tempo. Para testar se as curvas são diferentes, a autora utiliza o teste em questão e verifica quando as distribuições variam entre os tratamentos para qualquer ponto no tempo.

### 3.2.1 Teste baseados nas observações originais

Wang (2004) não trabalha com uma estrutura de simetria composta e considera  $\sigma_{ijj'} = Cov(Y_{ijk}, Y_{ij'k})$  e  $\sigma_{ijj} = \sigma_{ij}^2 = Var(Y_{ijk})$ . Sugere ainda alterações nos quadrados médios de tempo, interação e efeito simples, uma vez que os valores esperados desses quadrados não são iguais ao valor esperado do quadrado médio do erro sob a hipótese nula de interesse. As modificações são apresentadas abaixo:

$$MS\gamma = \frac{1}{(a-1)(b-1)} \sum_{i=1}^a \sum_{j=1}^b (\bar{Y}_{ij.} - \tilde{Y}_{i..} - \tilde{Y}_{.j.} + \tilde{Y}_{...})^2 \quad (3.2)$$

$$MSE = \frac{1}{a(b-1)} \sum_{i=1}^a \sum_{j=1}^b \frac{1}{n_i(n_i-1)} \sum_{k=1}^{n_i} (Y_{ijk} - \bar{Y}_{ij.} - \bar{Y}_{i.k} + \tilde{Y}_{i..})^2 \quad (3.3)$$

$$MS\phi = \frac{1}{(a-1)b} \sum_{i=1}^a \sum_{j=1}^b (\bar{Y}_{ij.} - \tilde{Y}_{.j.})^2 \quad e \quad (3.4)$$

$$MSE_\phi = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} \frac{(Y_{ijk} - \bar{Y}_{ij.})^2}{n_i(n_i-1)}. \quad (3.5)$$

Com valores esperados dados por:

$$E(MSE) = \frac{1}{a(b-1)} \sum_{i,j} \frac{\sigma_{ijj}}{n_i} - \frac{1}{ab(b-1)} \sum_{i=1}^a \sum_{j=1}^b \sum_{j'=1}^b \frac{\sigma_{ijj'}}{n_i},$$

$$E(MS\gamma) = E(MSE), \quad \text{sob } H_0(\gamma),$$

$$E(MSE_\phi) = \frac{1}{ab} \sum_{i,j} \frac{\sigma_{ijj}}{n_i}, \quad e$$

$$E(MS\phi) = E(MSE_\phi) \quad \text{sob } H_0(\phi).$$

Logo, as seguintes estatísticas podem ser utilizadas para avaliar as hipóteses  $H_0(\gamma)$  e  $H_0(\phi)$ , respectivamente:

$$F_\gamma = \frac{MS\gamma}{MSE} \quad e \quad F_\phi = \frac{MS\phi}{MSE_\phi}$$



A teoria assintótica das estatísticas dos testes mencionados é apresentada no teorema abaixo.

**Teorema 3.2.1.** (WANG, 2004) *Assuma que para cada grupo  $i$  e indivíduo  $k$ ,  $Y_{ijk}$ ,  $j=1,2,\dots$ , é  $\alpha$ -mixing com  $\alpha_m = O(m^{-5})$ . Além disso, considere que  $\limsup_j E[(Y_{ijk} - \mu_{ij})^{32}] < \infty$ , onde  $\mu_{ij} = E(Y_{ijk})$ . Então para  $b \rightarrow \infty$  enquanto  $a$  permanece fixo, os limites dos termos apresentados abaixo existem.*

$$\varsigma_1 = \frac{2}{a^2 b} \sum_{j=1}^b \sum_{j'=1}^b \sum_{i=1}^a \frac{\sigma_{ijj'}^2}{n_i(n_i - 1)}, \quad (3.6)$$

$$\varsigma_2 = \frac{2}{a^2 b} \sum_{j=1}^b \sum_{j'=1}^b \sum_{i \neq i'}^a \frac{\sigma_{ijj'} \sigma_{i'jj'}}{n_i n_{i'}} \quad (3.7)$$

De modo que,  $\sigma^2 = \lim_{b \rightarrow \infty} E(MSE) = \lim_{b \rightarrow \infty} E(MSE_\phi)$  e  $\sigma_*^2 = \lim_{b \rightarrow \infty} E(n(a)MSE) = \lim_{b \rightarrow \infty} E(n(a)MSE_\phi)$ .

(1) para  $n_i \geq 2$  fixo,

sob  $H_0(\gamma)$ ,  $\sqrt{b}(F_\gamma - 1) \xrightarrow{d} N\left(0, \frac{\tau_\gamma^2}{\sigma^4}\right)$ , onde  $\tau_\gamma^2 = \lim_{b \rightarrow \infty} \left(\varsigma_1 + \frac{\varsigma_2}{(a-1)^2}\right)$ ;

sob  $H_0(\phi)$ ,  $\sqrt{b}(F_\phi - 1) \xrightarrow{d} N\left(0, \frac{\tau_\phi^2}{\sigma^4}\right)$ .

(2) se  $n_i \rightarrow \infty$ , sob as condições adicionais  $\max_i\{n_i\}/\min_i\{n_i\} = O(1)$  e  $n(a) = \min\{n_i, 1 \leq i \leq a\}$ , são obtidos os seguintes resultados:

sob  $H_0(\gamma)$ ,  $\sqrt{b}(F_\gamma - 1) \xrightarrow{d} N\left(0, \frac{\tau_{\gamma*}^2}{\sigma_*^4}\right)$ , onde  $\tau_{\gamma*}^2 = \lim_{b \rightarrow \infty} n^2(a) \left(\varsigma_1 + \frac{\varsigma_2}{(a-1)^2}\right)$ ;

sob  $H_0(\phi)$ ,  $\sqrt{b}(F_\phi - 1) \xrightarrow{d} N\left(0, \frac{\tau_{\phi*}^2}{\sigma_*^4}\right)$ .

### 3.2.2 Testes baseados em postos

Conforme apresentado em Wang (2004), testes baseados nas observações originais são sensíveis a *outliers* e podem apresentar baixo desempenho para dados que não sejam provenientes da distribuição normal.

Desse modo, Wang sugere a construção de testes baseados em postos, ou seja, as observações originais  $(Y_{ijk})$  são substituídas pelos respectivos postos  $(R_{ijk})$  nas estatísticas dos testes de interesse. A substituição é possível devido à invariância dos postos à transformações monótonas<sup>1</sup>. Por meio de simulações, esses testes são comparados aos apresentados na subseção 3.2.1. São gerados diferentes cenários, especificamente com dados provenientes das distribuições normal, lognormal e Cauchy.

A probabilidade de erro do tipo I é computada para os testes de ausência dos seguintes efeitos: tempo, tratamento, interação entre tempo e tratamento e efeito simples.

Além disso, ela apresenta os resultados das simulações de poder dos testes. De modo que, os testes baseados em postos apresentam desempenho superior para os cenários das distribuições lognormal e Cauchy.

Considere  $MS\gamma_R$ ,  $MS\phi_R$ ,  $MSE_R$  e  $MSE_{\phi_R}$  os quadrados médios apresentados nas equações de (2.2) a (2.6). O processo consiste, então, em colocar  $R_{ijk}$  no lugar de  $Y_{ijk}$ . E desse modo, obter as estatísticas indicadas abaixo.

$$F_{R,\gamma} = \frac{MS\gamma_R}{MSE_R} \text{ e } F_{R,\phi} = \frac{MS\phi_R}{MSE_{\phi_R}}$$

Os resultados assintóticos são descritos na seqüência.

**Teorema 3.2.2.** (WANG, 2004) *Assuma que para cada grupo  $i$  e indivíduo  $k$ ,  $Y_{ijk}$ ,  $j=1,2,\dots$ , é  $\alpha$ -mixing com  $\alpha_m = O(m^{-5})$ . Seja  $R_{ijk}=H(Y_{ijk})$  e  $\tilde{\sigma}_{ijj'}=cov(R_{ijk},R_{ij'k})$ . De modo que  $\tilde{\tau}_\beta$ ,  $\tilde{\tau}_\gamma$ ,  $\tilde{\sigma}^4$ ,  $\tilde{\tau}_{\beta*}$ ,  $\tilde{\tau}_{\gamma*}$  e  $\tilde{\sigma}_*^4$  são definidos conforme o Teorema 3.2.1 com a substituição de  $\sigma_{ijj'}$  por  $\tilde{\sigma}_{ijj'}$ . Assim, para  $b \rightarrow \infty$  enquanto  $a$  permanece fixo, são observados os seguintes resultados:*

(1) para  $n_i \geq 2$  fixo,

$$\text{sob } H_0(\tilde{\gamma}), \sqrt{b}(F_{R,\gamma} - 1) \xrightarrow{d} N\left(0, \frac{\tilde{\tau}_\gamma^2}{\tilde{\sigma}^4}\right);$$

$$\text{sob } H_0(\tilde{\phi}), \sqrt{b}(F_{R,\phi} - 1) \xrightarrow{d} N\left(0, \frac{\tilde{\tau}_\gamma^2}{\tilde{\sigma}^4}\right).$$

---

<sup>1</sup>Transformações que preservam a ordem original do conjunto de dados.

(2) se  $n_i \rightarrow \infty$ , sob as condições adicionais  $\max_i\{n_i\}/\min_i\{n_i\} = O(1)$  e  $n(a) = \min\{n_i, 1 \leq i \leq a\}$ , são obtidos os seguintes resultados:

$$\text{sob } H_0(\tilde{\gamma}), \sqrt{b}(F_{R,\gamma} - 1) \xrightarrow{d} N\left(0, \frac{\tilde{\gamma}_{\gamma^*}^2}{\tilde{\sigma}_*^4}\right);$$

$$\text{sob } H_0(\tilde{\phi}), \sqrt{b}(F_{R,\phi} - 1) \xrightarrow{d} N\left(0, \frac{\tilde{\gamma}_{\gamma^*}^2}{\tilde{\sigma}_*^4}\right).$$

### 3.3 Testes (von Borries, 2008)

von Borries (2008) considera uma estrutura fatorial com  $a$  relativamente grande ou  $a \rightarrow \infty$ ,  $b$  fixo e  $n_i$  pequeno ( $\geq 2$  e limitado). De modo que, esses dados são denominados HDLLSS (*High Dimensional Longitudinal Low Sample Size* - Longitudinais Superdimensionados de Amostra Pequena).

O teste apresentado pelo autor consiste apenas na avaliação da hipótese de ausência de efeito simples, ou seja, aquela correspondente à comparações entre os níveis de um fator para um único nível de outro.

Assim, o modelo clássico ANOVA estabelece  $Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$  para  $i=1, \dots, a$ ;  $j=1, \dots, b$ ;  $k=1, \dots, n_i$ , onde  $\sum_{i=1}^a \alpha_i = \sum_{j=1}^b \beta_j = \sum_{i=1}^a \gamma_{ij} = \sum_{j=1}^b \gamma_{ij} = 0$ ,  $E(\epsilon_{ijk})=0$  e  $Y_{ijk} \sim F_{ij}$  arbitrária. De modo que, a hipótese paramétrica de ausência de efeito simples é descrita abaixo:

$$H_0(\phi) : \phi_{ij} = \alpha_i + \gamma_{ij} = 0, \quad \text{para todo } i, j \quad (3.8)$$

#### 3.3.1 Testes baseados nas observações originais

Assume-se a estrutura de covariância indicada abaixo, a qual implica que observações de um mesmo nível de fator e repetição são correlacionadas, enquanto que observações provenientes de níveis e/ou repetições diferentes não são.

Além disso, trabalha-se com a suposição de independência de  $Y_{ijk}$  para os conjuntos formados por uma repetição e seu respectivo nível de fator. Observa-se que a estrutura empregada é a mesma utilizada por Wang (2004).

$$\text{cov}(Y_{ijk}, Y_{i'j'k'}) = \begin{cases} \sigma_{ijj'} & \text{se } i = i', k = k' \\ 0 & \text{se } i \neq i', k \neq k' \end{cases}$$

As seguintes notações são adotadas na obtenção dos resultados:

$$\begin{aligned} \tilde{Y}_{.j} &= \frac{1}{a} \sum_{i=1}^a \bar{Y}_{ij}, \\ \bar{Y}_{ij} &= \frac{1}{n_i} \sum_{k=1}^{n_i} Y_{ijk}, \\ \bar{Y}_{...} &= \frac{1}{N} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} Y_{ijk}, \\ \tilde{Y}_{...} &= \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b \bar{Y}_{ij}, \text{ e} \\ N &= b \sum_{i=1}^a n_i. \end{aligned}$$

As estatísticas apresentadas abaixo utilizam observações originais como variável resposta. De modo que,  $E(\text{MS}\varphi) = E(\text{MSE})$  sob  $H_0(\phi)$ .

$$\text{MS}\varphi = \frac{1}{(a-1)b} \sum_{i=1}^a \sum_{j=1}^b (\bar{Y}_{ij} - \tilde{Y}_{.j})^2 \quad (3.9)$$

$$\text{MSE} = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} \frac{(Y_{ijk} - \bar{Y}_{ij})^2}{n_i(n_i - 1)} \quad (3.10)$$

A teoria assintótica para a estatística do teste de efeito simples é apresentada no teorema abaixo.

**Teorema 3.3.1.** ((VON BORRIES, 2008). *Teste de ausência de efeito simples*).  
Seja  $Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$ , onde  $\mu$  é a média global;  $\alpha_i$  ( $i = 1, \dots, a$ ) o efeito médio do fator;  $\beta_j$  ( $j = 1, \dots, b$ ) o efeito do ponto no tempo;  $\gamma_{ij}$  o efeito da interação fator-tempo; e  $\epsilon_{ijk}$  um erro aleatório com distribuição arbitrária  $F_{ij}$ , para todo  $k = 1, \dots, n_i$ .

Considere,  $H_0(\phi) : \phi_{ij} = \alpha_i + \gamma_{ij} = 0$  verdadeira. Se as observações  $Y_{ijk}$  têm momento central finito  $(2+\delta)(\delta > 0)$ , e o número de repetições é pequeno, com  $n_i \geq 2$  e limitado, observado para um número fixo  $b$  de pontos no tempo e com  $a \rightarrow \infty$ ,

$$F_\phi = \sqrt{ab}(\text{MS}\varphi - \text{MSE}) \xrightarrow{d} N \left( 0, \lim_{a \rightarrow \infty} \frac{2}{ab} \sum_{i=1}^a \frac{1}{n_i(n_i - 1)} \sum_{j=1}^b \sum_{j_1=1}^b \sigma_{ijj'}^2 \right) \quad (3.11)$$

### 3.3.2 Testes baseados em postos

Como os testes baseados nas observações originais apresentam as mesmas restrições apontadas por Wang (2004), é interessante considerar a versão baseada em postos.

Suponha então que  $R_{ijk}$  corresponde a representação em postos da observação  $Y_{ijk}$  e  $\tilde{R}_{.j} = \frac{1}{a} \sum_{i=1}^a \bar{R}_{ij}$ . As seguintes estatísticas são definidas:

$$MS\varphi_R = \frac{1}{(a-1)b} \sum_{i=1}^a \sum_{j=1}^b (\bar{R}_{ij} - \tilde{R}_{.j})^2 \quad (3.12)$$

$$MSE_R = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_i} \frac{(R_{ijk} - \bar{R}_{ij})^2}{n_i(n_i - 1)} \quad (3.13)$$

O resultado assintótico abaixo é semelhante ao apresentado no Teorema 3.3.1. As mesmas condições são consideradas.

$$F_{R,\phi} = \sqrt{ab}(MS\varphi_R - MSE_R) \xrightarrow{d} N\left(0, \lim_{a \rightarrow \infty} \frac{2}{ab} \sum_{i=1}^a \frac{1}{n_i(n_i - 1)} \sum_{j=1}^b \sum_{j_1=1}^b \sigma_{ijj_1}^2\right) \quad (3.14)$$

A demonstração da convergência acima não está formalizada na literatura. Todavia, Silva (2012) trabalha com a verificação dessa convergência através da utilização de diferentes cenários de simulação.

O autor utiliza o teste de Shapiro-Wilk para avaliar a hipótese apresentada abaixo e conclui com os resultados obtidos que o teste baseado em postos apresenta a mesma convergência que o baseado nas observações originais.

$$H_0 : F_{R,\phi'} = \frac{\sqrt{ab}}{\sqrt{n}}(MS\varphi_R - MSE_R) \xrightarrow{d} N(0, 1) \quad (3.15)$$

## 3.4 Testes (Zhang, 2008)

O avanço da tecnologia proporcionou um aumento considerável na quantidade de dados biológicos armazenados no últimos anos. Zhang (2008) utiliza um conjunto de testes não-paramétricos para analisar dados de aCGH caracterizados por um número elevado de tratamentos (genes) e um número fixo de repetições e pontos no tempo.

Segundo (Capelli e Nascimento, 2009), para a realização de um experimento aCGH, milhares de pequenas seqüências de DNA referentes a regiões específicas do genoma humano são dispostas na superfície de lâminas. As moléculas de DNA do indivíduo a ser avaliado quanto à presença de anomalias cromossômicas são divididas em grupos de teste e controle.

Assim, o aCGH possibilita que praticamente todo o genoma do paciente testado seja averiguado e que a localização genômica de eventuais alterações possa ser facilmente identificada.

De acordo com Maluf e Riegel (2011), o uso dessa tecnologia tem possibilitado o diagnóstico molecular de alterações cromossômicas submicroscópicas previamente não detectadas, principalmente em indivíduos com deficiência mental e/ou com múltiplas malformações congênitas de causa desconhecida.

Para analisar dados de aCGH, Zhang (2008) trabalha com as seguintes hipóteses:

1. Avaliar a ausência de efeito de tratamento, ou seja, verificar se o número local de cópias de DNA varia em uma determinada região genômica;
2. Detectar as regiões genômicas que apresentam o mesmo tempo de resposta, isto é, testar a hipótese de ausência do efeito da interação entre tempo e tratamento.

Para avaliar as hipóteses de interesse, Zhang (2008) adota um conjunto de testes não-paramétricos com base no mesmo modelo adotado por von Borries (2008). A suposição de normalidade não é assumida para a componente de erro  $\epsilon_{ijk}$  e a estrutura de covariância heterocedástica é dada por  $Cov(\epsilon_{ijk}, \epsilon_{ij'k}) = \sigma_{ijj'}$ .

Os teoremas assintóticos são apresentados na sequência. E as seguintes notações são adotadas pelo autor na obtenção dos resultados:

$$\sigma_{i,jj_1} = Cov(Y_{ijk}, Y_{ij_1k}) \text{ para qualquer } k \text{ e } \sigma_{i,jj} = Var(Y_{ijk}) = \sigma_{ij}^2$$

**Teorema 3.4.1.** (ZHANG, 2008). A estatística  $F_Y(\alpha) = \frac{MS\alpha}{MSE\alpha}$  é adotada pelo autor no teste de ausência do efeito de tratamento, de modo que os quadrados médios são dados por

$$MS\alpha = \frac{1}{a-1} \sum_{i=1}^a \sum_{j=1}^b (\tilde{Y}_{i..} - \tilde{Y}_{...})^2 \quad e \quad (3.16)$$

$$MSE\alpha = \frac{1}{ab} \sum_{i=1}^a \sum_{j,j_1}^b \frac{1}{n_i(n_i-1)} \sum_{k=1}^{n_i} (Y_{ijk} - \bar{Y}_{ij.})(Y_{ij_1k} - \bar{Y}_{ij_1.}) \quad (3.17)$$

Se  $Y_{ijk}$  apresenta quarto momento finito, então sob  $H_0(\alpha)$ , a estatística do teste possui a seguinte distribuição assintótica:

$$\frac{\sqrt{a}(F_Y(\alpha) - 1)}{V_\alpha} \xrightarrow{d} N(0,1) \quad \text{para } a \rightarrow \infty \quad (3.18)$$

De modo que, a componente de variância  $V_\alpha = \frac{\sqrt{\tau_\alpha}}{\sigma_\alpha}$  é caracterizada por

$$\tau_\alpha = \frac{1}{ab^2} \sum_{i=1}^a \frac{2}{n_i(n_i-1)} \sum_{j,j_1,j_2,j_3}^b (\sigma_{i,jj_1} \sigma_{i,j_2j_3}) \quad (3.19)$$

$$\sigma_\alpha = \frac{1}{ab} \sum_{i=1}^a \sum_{j,j_1}^b \frac{\sigma_{i,jj_1}}{n_i} \quad (3.20)$$

**Teorema 3.4.2.** (ZHANG, 2008). Para testar  $H_0(\gamma)$ : todo  $\gamma_{ij}=0$ , para  $i=1,\dots,a$  e  $j=1,\dots,b$ , o autor considera  $F_Y(\gamma) = \frac{MS\gamma}{MSE\gamma}$ . Tal que,

$$MS\gamma = \frac{1}{(a-1)(b-1)} \sum_{i=1}^a \sum_{j=1}^b (\tilde{Y}_{ij.} - \tilde{Y}_{i..} - \tilde{Y}_{.j.} + \tilde{Y}_{...})^2 \quad (3.21)$$

$$MSE\gamma = \frac{1}{a(b-1)} \sum_{i=1}^a \sum_{j=1}^b \frac{1}{n_i(n_i-1)} \sum_{k=1}^{n_i} (Y_{ijk} - \bar{Y}_{ij.})^2 \quad (3.22)$$

$$- \frac{1}{ab(b-1)} \sum_{i=1}^a \sum_{j,j_1}^b \frac{1}{n_i(n_i-1)} \sum_{k=1}^{n_i} (Y_{ijk} - \bar{Y}_{ij.})(Y_{ij_1k} - \bar{Y}_{ij_1.})$$

Se  $Y_{ijk}$  apresenta quarto momento finito, então sob  $H_0(\gamma)$ , a estatística do teste possui a distribuição assintótica indicada abaixo:

$$\frac{\sqrt{a}(F_Y(\gamma) - 1)}{V_\gamma} \xrightarrow{d} N(0,1) \quad (3.23)$$

De modo que,  $V_\gamma = \frac{\sqrt{\tau_\gamma}}{\sigma_\gamma}$  é a componente de variância definida por

$$\tau_\gamma = \frac{2}{a(b-1)^2} \sum_{i=1}^a \left[ \frac{1}{n_i(n_i-1)} \sum_{j,j_1}^b \sigma_{i,jj_1}^2 + \frac{1}{b^2 n_i(n_i-1)} \right] \quad (3.24)$$

$$\sigma_\gamma = \frac{1}{a(b-1)} \sum_{i=1}^a \sum_{j=1}^b \frac{\sigma_{i,j}^2}{n_i} - \frac{1}{ab(b-1)} \sum_{i=1}^a \sum_{j,j_1}^b \frac{\sigma_{i,jj_1}}{n_i} \left[ \sum_{j,j_1,j_2,j_3}^b \sigma_{i,jj_1} \sigma_{ij_2j_3} - \frac{2}{bn_i(n_i-1)} \sum_{j,j_1,j_2}^b \sigma_{i,jj_1} \sigma_{i,jj_2} \right] \quad (3.25)$$

O autor também define versões baseadas em postos dos testes apresentados no teorema anterior.

De modo geral, as estatísticas dos testes baseados em postos são idênticas às anteriores, com exceção da substituição da resposta  $Y_{ijk}$  pelo respectivo posto  $R_{ijk}$ . De forma que, a convergência permanece a mesma.

### 3.5 Algoritmo de agrupamento por partição

O p-valor obtido no teste de ausência de efeito simples baseado em observações originais é adotado como medida de similaridade no algoritmo de agrupamento PPCLUS-TEL, também desenvolvido pelo autor.

Silva (2012) trabalha com a versão baseada em postos desse algoritmo, a qual é denominada PPCLUSTEL-R e está baseada no teste apresentado na Subseção 3.3.2.

A idéia do algoritmo é testar iterativamente se grupos que contêm vários níveis de fator apresentam similaridade. De modo que, um grupo é particionado em dois menores quando o teste de efeito simples é rejeitado. Caso contrário, o grupo permanece intacto. O algoritmo pára quando não existem mais similaridades abaixo do limiar (nível de significância do teste).

Segundo von Borries (2008), fatores que não podem ser alocados a nenhum dos grupos criados são classificados como grupo 0. Isso significa que esses fatores resultam da rejeição da hipótese nula de ausência de efeito simples em todos os outros grupos criados.

Seja  $g$  o índice do grupo no qual o teste está sendo aplicado,  $D1$  contém as observações do grupo  $g$  e  $nf$  representa o número de fatores em  $D1$ . O algoritmo é descrito abaixo em sete passos.

1. Seja  $g = 1$ ,  $D1 = \text{Dados}$ .



2. Calcule a mediana para cada fator em  $D1$ .
3. Ordene os fatores em  $D1$  pelas suas medianas.
4. Teste  $D1$ .
  - 4.1. Se  $H_0$  não é rejeitada: algoritmo termina.
  - 4.2. Se  $H_0$  é rejeitada: vá para o Passo 5.
5. Retire um subconjunto de  $D1$  e chame de  $D2$ .
6. Calcule o número de fatores em  $D2$  e chame de  $nf$ .
  - 6.1. Se  $nf = 1$ :
    - 6.1.1. Aloque fator em  $D2$  para o grupo 0.
    - 6.1.2. Remova os fatores em  $D2$  de  $D1$ .
    - 6.1.3. Se  $nf$  em  $D1 = 0$ , então algoritmo termina.
    - 6.1.4. Se  $nf$  em  $D1 > 0$ , então faça  $D2 = D1$  e vá para o Passo 7.
7. Teste  $D2$ .
  - 7.1 Se  $H_0$  não é rejeitada:
    - 7.1.1 Atribua fatores de  $D2$  para o grupo  $g$ .
    - 7.1.2 Faça  $g = g + 1$ .
    - 7.1.3 Remova os fatores de  $D1$  em  $D2$ .
    - 7.1.4. Se  $nf$  em  $D1 = 0$  então algoritmo termina.
    - 7.1.5. Se  $nf$  em  $D1 > 0$  então faça:
      - A. Teste se cada fator em  $D1$  pertence ao novo grupo assinalado. Remova o fator de  $D1$  quando  $H_0$  não é rejeitada e o coloque nesse novo grupo.
      - B. Faça  $D2 = D1$  para os fatores remanescentes em  $D1$  e retorne ao Passo 7.

7.2 Se  $H_0$  é rejeitada:

7.2.1 Retire um subconjunto de  $D2$  e chame de  $D3$ .

7.2.2 Retorne para  $D1$  todos os fatores que não estão em  $D3$ .

7.2.3 Faça  $D2 = D3$  e remova  $D3$ .

7.2.4 Retorne ao Passo 7.

Segundo Silva (2012), conforme verificado em estudos de simulação, ao repetir o algoritmo diversas vezes com alteração dos limiares estabelecidos verifica-se que a partir de determinado limiar ocorre uma homogeneização dos grupos. Assim, a sugestão é a escolha do limiar que estabiliza os resultados.

# Capítulo 4

## Simulação e Aplicação de Algoritmos

Neste capítulo é conduzido um estudo de simulação Monte Carlo com o intuito de avaliar o poder dos testes de efeito simples apresentados por Wang (2004), von Borries (2008) e Zhang (2008).

Para a simulação dos dados experimentais e avaliação das curvas de poder dos testes em questão, foram desenvolvidos algoritmos no software estatístico SAS versão 9.2 com base no algoritmo PPCLUSTEL descrito em von Borries (2008).

Como o PPCLUSTEL trabalha exatamente com um dos testes de interesse, a macro SAS responsável pela realização do processo de agrupamento foi desmembrada e acrescida de outros comandos para que pudessem ser construídas as curvas de poder do teste que a compõe.

Para os testes apresentados por Wang (2004) e Zhang (2008), foram desenvolvidos algoritmos com base no mesmo delineamento fatorial adotado por von Borries (2008), ou seja, os dados foram gerados a partir do modelo  $Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$ .

De modo que, as macros WANG\_ES e ZHANG\_ES foram construídas para a obtenção dos resultados de interesse. Maiores detalhes podem ser vistos no Anexo C.

Para os testes de efeito simples propostos por Wang (2004) e von Borries (2008), o nível crítico foi calculado com base nos teoremas 3.2.2 e 3.3.1, respectivamente.

Além disso, é importante ressaltar que as curvas de poder do teste de Zhang (2008) foram obtidas a partir da definição dos testes de ausência de efeito de tratamento e de interação apresentados no teorema 3.4.1.

Isso foi possível já que os níveis críticos obtidos nesses testes puderam ser utilizados na construção de uma regra de decisão para o teste de ausência de efeito simples. A tabela 4.1 apresenta o critério de decisão adotado.

Tabela 4.1: Critério de decisão para o teste de efeito simples com base na metodologia de Zhang (2008).

Tratamento	Interação	Efeito simples
Não rejeita	Não rejeita	Não rejeita
Não rejeita	Rejeita	Rejeita
Rejeita	Não rejeita	Rejeita
Rejeita	Rejeita	Rejeita

A geração das respostas b-dimensionais  $Y_{ik} = (Y_{i1k}, Y_{i2k}, \dots, Y_{ibk})$  foi conduzida no ambiente SAS/ IML (*Interactive Matrix Language*). O módulo RANDNORMAL foi adotado na geração das respostas provenientes da distribuição normal multivariada com vetor de médias  $\mu_j \in \mathbb{R}^b$  e matriz de covariância  $\Sigma$ , simétrica e positiva definida. As respostas foram produzidas a partir da seguinte função densidade de probabilidade:

$$f(\mathbf{y}; \mu_j, \Sigma) = \frac{1}{(2\pi)^{b/2} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{y} - \mu_j)\Sigma^{-1}(\mathbf{y} - \mu_j)'}{2}\right) \quad (4.1)$$

Foram obtidos ainda, dados experimentais amostrados da distribuição t de Student multivariada com  $\nu=5$  graus de liberdade. O módulo RANDMVT com a função densidade descrita abaixo foi utilizado ao longo do processo.

$$f(\mathbf{y}; \nu, \mu_j, \Sigma) = \frac{\Gamma((\nu + b)/2)}{|\Sigma|^{1/2} (\pi\nu)^{b/2} \Gamma(\nu/2)} \left(1 + \frac{(\mathbf{y} - \mu_j)\Sigma^{-1}(\mathbf{y} - \mu_j)'}{\nu}\right)^{-(\nu+b)/2} \quad (4.2)$$

Os componentes do vetor de médias foram iguais a  $\mu_j = \cos(\pi(j + 1))$ , para  $j = 1, \dots, b$ . De modo que, para a matriz de covariância, os componentes seguiram a relação  $\sigma_{ijj_1} = 1 - |j - j_1| \times 0,2$ .

O número de replicações do processo de simulação foi igual a 400 para cada combinação de cenários. E os percentuais de observações contaminadas na média foram de 5% e 10%.

Foram simulados 864000 experimentos (2160 cenários x 400 experimentos por cenário). Os 2160 cenários foram resultantes das combinações entre fatores, testes avaliados, pontos no tempo, repetições por fator, desvios da média, distribuições de probabilidade e percentuais de observações contaminadas.

Assim, para que as curvas de poder do teste pudessem ser construídas, contabilizou-se em quantas das amostras geradas a hipótese nula de ausência de efeito simples foi rejeitada. Tal que, o poder empírico foi obtido pela quantidade de vezes que o teste rejeitou  $H_0$  dividida pelo número total de replicações do procedimento.

Os postos das observações originais foram utilizados como variável resposta ao longo de todo o processo de simulação, o que fez com que os testes fossem invariantes a transformações monótonas.

O resumo abaixo apresenta os diferentes cenários considerados no processo.

1. Objetivo: avaliar o poder dos testes de ausência de efeito simples com base nas metodologias desenvolvidas por Wang (2004), von Borries (2008) e Zhang (2008).
2. Distribuições: normal multivariada e t de Student multivariada com 5 graus de liberdade.
3. Percentual de observações contaminadas na média ( $p$ ): 5% e 10%.
4. Desvio da média ( $d$ ): 0, 0.4, 0.8 e 1.2.
5. Vetor de médias:  $\mu_j = \cos(\pi(j + 1))$ , para  $j = 1, \dots, b$ .
6. Matriz de covariância ( $\Sigma$ ) utilizada na geração dos dados: matriz com variâncias unitárias e covariâncias que decrescem para observações mais distantes no tempo seguindo  $\sigma_{ijj_1} = 1 - |j - j_1| \times 0.2$ .
7. Estrutura de dados longitudinais:
  - (a) Fatores ( $a$ ): 50, 100, 200, 400, 800.
  - (b) Pontos no tempo ( $b$ ): 10, 20, 40.
  - (c) Repetições ( $n_i$ ): 5, 10, 20.

É importante destacar que os p-valores foram obtidos por meio da estimação dos termos  $\sigma_{ijj_1}^2$  e  $\sigma_{ijj_1}$ , que compõem as estatísticas dos testes sob análise.

Basicamente, o ideal é que o processo de estimação fosse conduzido para todo  $i \neq i'$  com a estimação de  $a$  matrizes  $\Sigma_i$ . Entretanto, isso nem sempre é possível.

Considere, por exemplo, uma situação com  $n_i=4$  e  $b=10$ . Nesse caso, o processo de estimação dos 55 parâmetros de covariância<sup>1</sup> é prejudicado pelo número limitado de repetições por nível de tratamento.

Para solucionar o problema, Silva (2012) considerou uma mesma estrutura para todos os níveis de fator ( $\sigma_{1jj_1} = \dots = \sigma_{ajj_1}$ ), o que implicou na estimação de uma única matriz  $\Sigma$  por meio da utilização de todo o conjunto de dados<sup>2</sup>.

Assim, a estimação dessa estrutura foi efetuada via *jackknife*. A idéia central consistiu em repartir a amostra  $(y_1, \dots, y_n)$  em grupos mutuamente excludentes de igual tamanho e em seguida calcular os estimadores para cada um deles. De modo que, cada um dos elementos da amostra correspondeu a uma das linhas da tabela 2.1. Assim, o estimador *jackknife* de

$$\sigma_n^4 = \left( \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n} \right)^2 \quad (4.3)$$

foi dado por

$$\hat{\sigma}_{jack}^4 = n\sigma_n^4 - \frac{n-1}{n} \sum_{i=1}^n \sigma_{n(i)}^4 \quad (4.4)$$

Pawitan (2001) comparou os estimadores  $\sigma_n^4$  e  $\sigma_{jack}^4$  quanto ao viés e concluiu que o viés do primeiro é maior que o do segundo, motivo pelo qual  $\sigma_{jack}^4$  foi adotado na estimação de  $\sigma_{ijj_1}^2$ .

A macro S4JACKKNIFE foi adotada no processo de estimação da matriz  $\Sigma$ . Essa macro foi utilizada por Silva (2012) na comparação de diferentes métodos de estimação da matriz de covariância dos dados, inclusive as formas  $\Sigma$  e  $\Sigma_i$  mencionadas anteriormente. Maiores informações sobre a macro podem ser vistas no Anexo C.

<sup>1</sup>Como a matriz de covariância é simétrica, o número de parâmetros distintos dessa matriz é dado por  $\frac{b(b+1)}{2}$ .

<sup>2</sup>A matriz  $\Sigma$  possui dimensão definida pelo número de pontos no tempo ( $b$ ).

## 4.1 Resultados em dados simulados

As figuras 4.1 e 4.2 mostram o comportamento do poder do teste de ausência de efeito simples formulado por Zhang (2008). As curvas foram avaliadas para as duas distribuições multivariadas e para um percentual de contaminação igual a 5%. Foram consideradas apenas as situações com 20 repetições por tratamento e todas as possíveis quantidades de pontos no tempo.

A observação dessas figuras mostra que o teste apresenta uma boa capacidade de detecção de pequenos deslocamentos à medida que o número de tratamentos envolvidos aumenta. Aspecto que vai de encontro à própria definição do autor, já que o teste é descrito como adequado para situações que envolvem  $a$  relativamente grande ou  $a \rightarrow \infty$  e um número fixo de repetições por tratamento e pontos no tempo.

Para quantidades menores de fatores (50, 100 e 200), a diferença precisa ser um pouco maior para que o teste detecte deslocamentos nos dados que possuem cauda pesada (distribuição t de Student com 5 graus de liberdade).

Além disso, à medida que o número de pontos no tempo aumenta, observa-se uma melhoria no desempenho do teste em questão independente da distribuição utilizada. Um exemplo são os casos com 200 ou mais fatores e 40 pontos no tempo, para os quais as distribuições normal e t multivariadas apresentam curvas semelhantes.

Entretanto, isso não é suficiente para afirmar que o teste funciona para situações com uma grande quantidade de tratamentos e pontos no tempo. Até porque isso também depende da quantidade de repetições consideradas.

Observa-se que para os casos com 400 e 800 fatores, as diferenças entre as curvas foram mínimas. Comportamento que sugere uma possível estabilidade para um número de fatores próximos a 400 e um número de repetições próximo de 10.

Apesar de terem sido apresentados graficamente apenas os casos com 20 repetições por tratamento, é importante ressaltar que dos casos com  $n_i = 5$  e  $n_i = 10$  foram tiradas as mesmas conclusões.

As curvas de poder dos testes de Wang (2004) e von Borries (2008) apresentaram comportamento semelhante na comparação das diferentes distribuições. Os mesmos cenários das figuras 4.1 e 4.2 foram considerados. Os detalhes das comparações entre as diferentes metodologias são apresentados na próxima subseção.

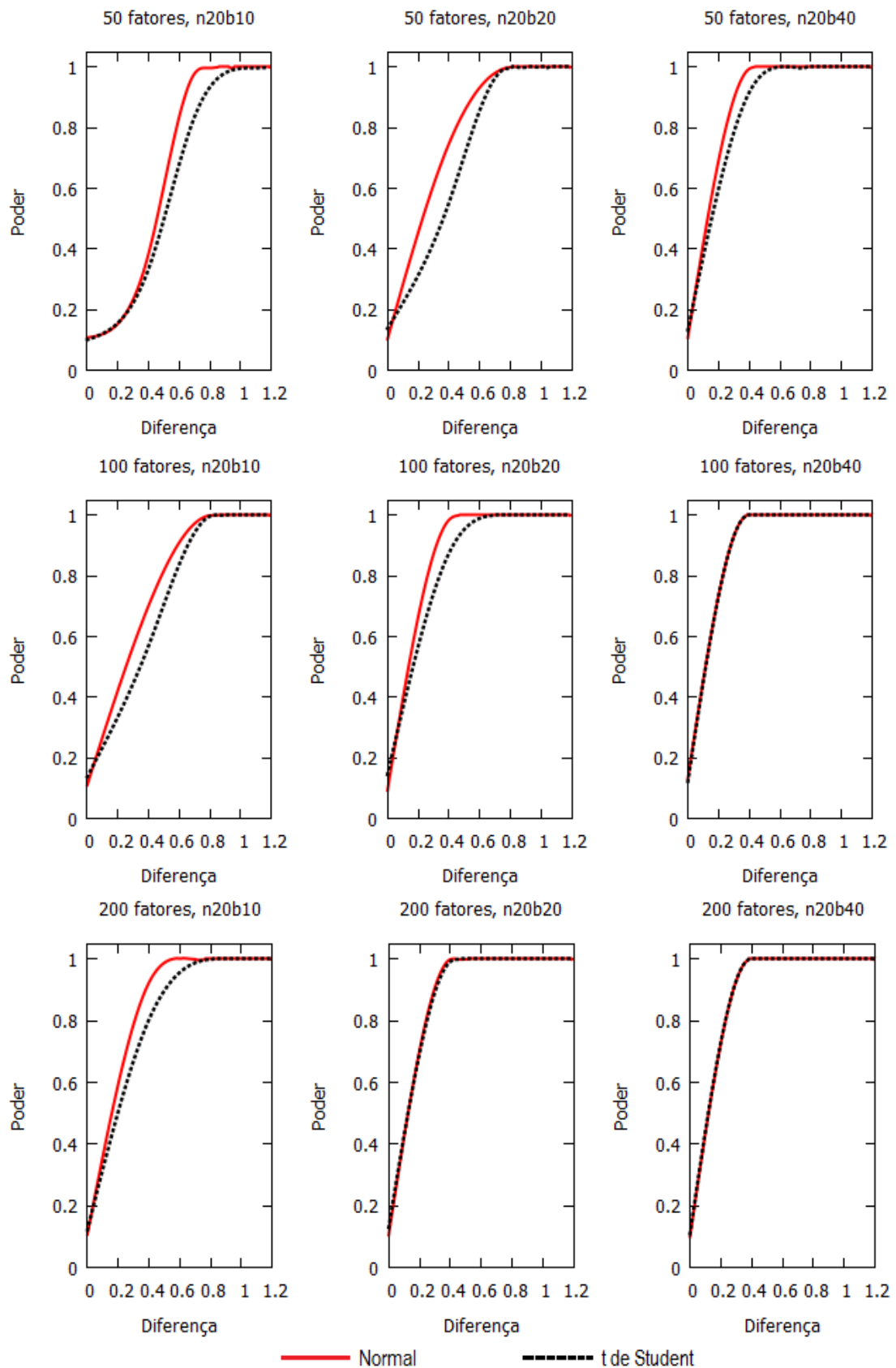


Figura 4.1: Poder do teste. 50, 100 e 200 fatores com 5% de contaminação e 20 repetições por tratamento.



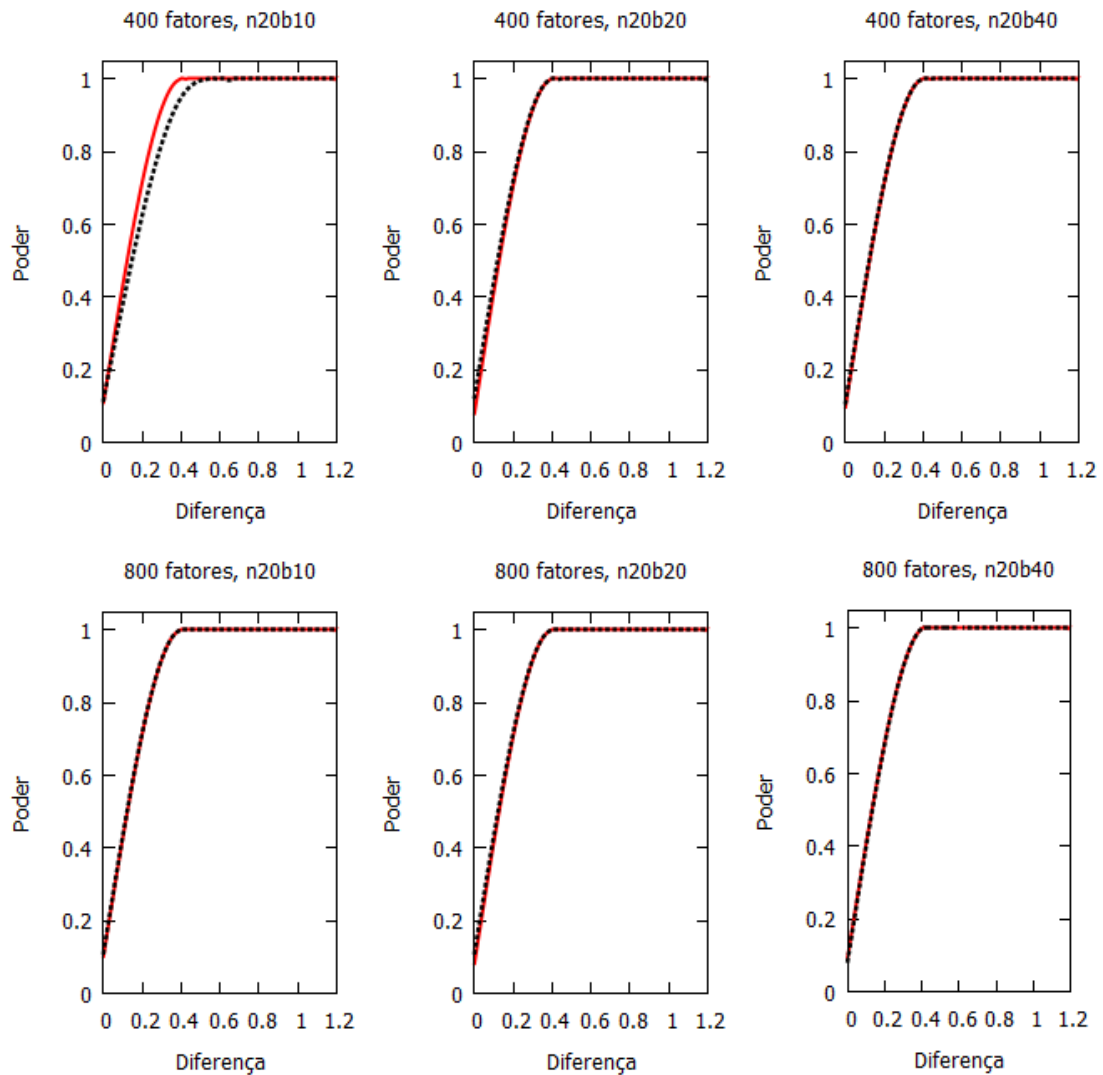


Figura 4.2: Poder do teste. 400 e 800 fatores com 5% de contaminação e 20 repetições por tratamento.

A figura 4.3 ilustra a curva de poder do teste de Zhang no caso com 800 fatores, 20 pontos no tempo, diferentes repetições por tratamento e 5% de contaminação.

É possível notar que quanto maior o número de repetições, maior é a facilidade desse teste detectar pequenos deslocamentos em uma fração razoavelmente pequena de dados contaminados.

Acentua-se que tal comportamento era esperado, já que quanto maior o número de repetições, melhores são as estimativas da matriz de covariância dos dados.

Ao avaliar o poder em função do percentual de contaminação da média, percebe-se que nos gráficos da figura 4.4 é necessário um deslocamento ligeiramente maior para

que o teste detecte a diferença em um menor percentual de dados contaminados.

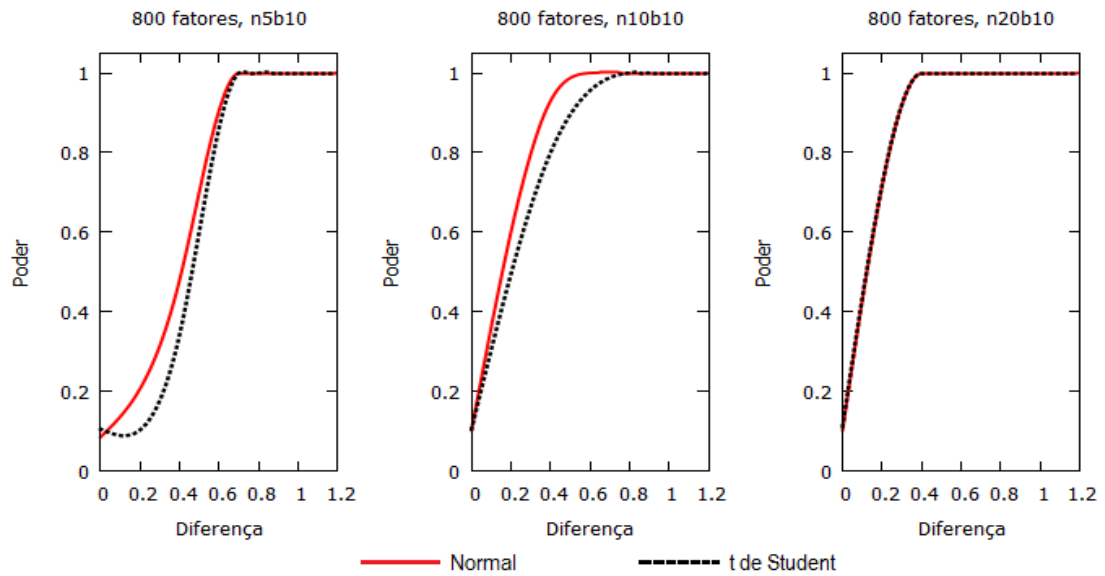


Figura 4.3: Poder do teste. 800 fatores e diferentes repetições por tratamento.

A figura 4.4 leva em consideração dados provenientes da distribuição normal com 100 fatores e 20 repetições por tratamento. De modo que, as curvas de 5% e 10% de contaminação se aproximam à medida que o número de pontos no tempo aumenta. Para quantidades de fatores mais elevadas ( $a = 400$ ), não foi verificada diferença entre as curvas dos diferentes níveis de contaminação.

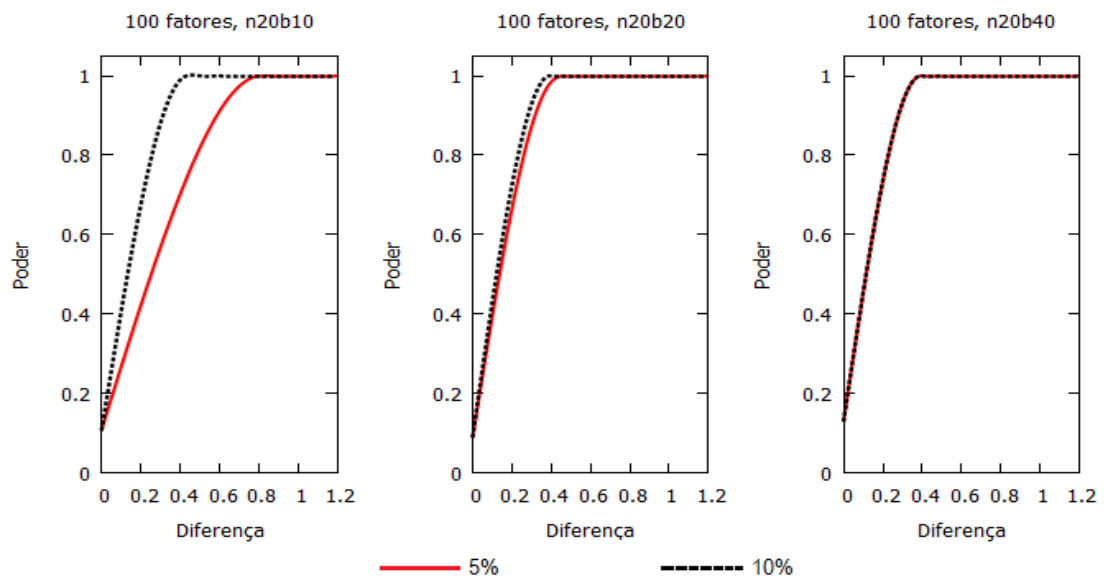


Figura 4.4: Poder do teste. 100 fatores e diferentes níveis de contaminação.

### 4.1.1 Comparação entre os algoritmos

As figuras 4.4 e 4.5 mostram as curvas de poder dos testes de efeito simples de Wang (2004), von Borries (2008) e Zhang (2008) para os casos de 50, 100, 200, 400 e 800 fatores. São considerados apenas os casos provenientes da distribuição normal multivariada com 20 repetições por tratamento.

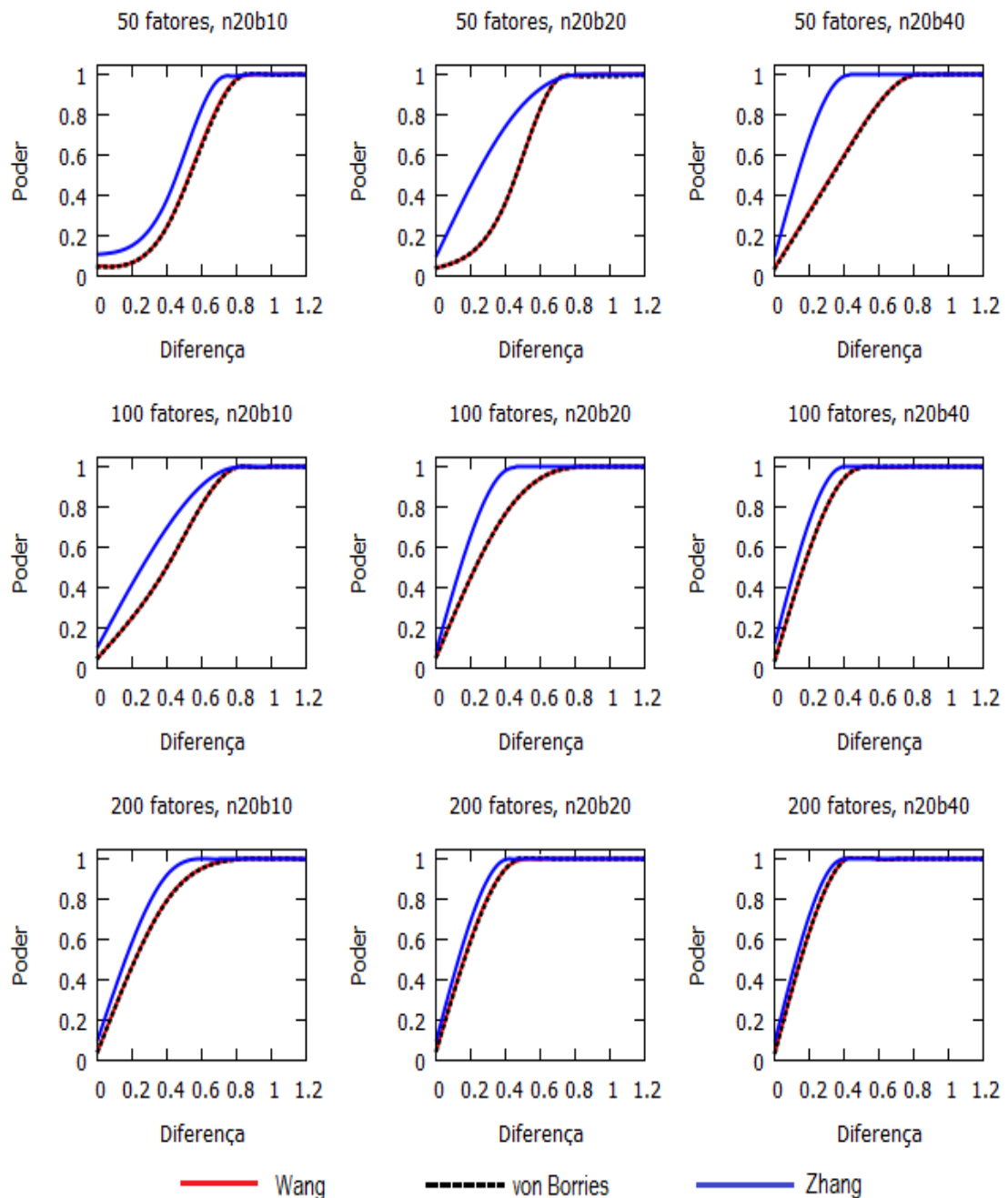


Figura 4.5: Poder do teste segundo diferentes metodologias.

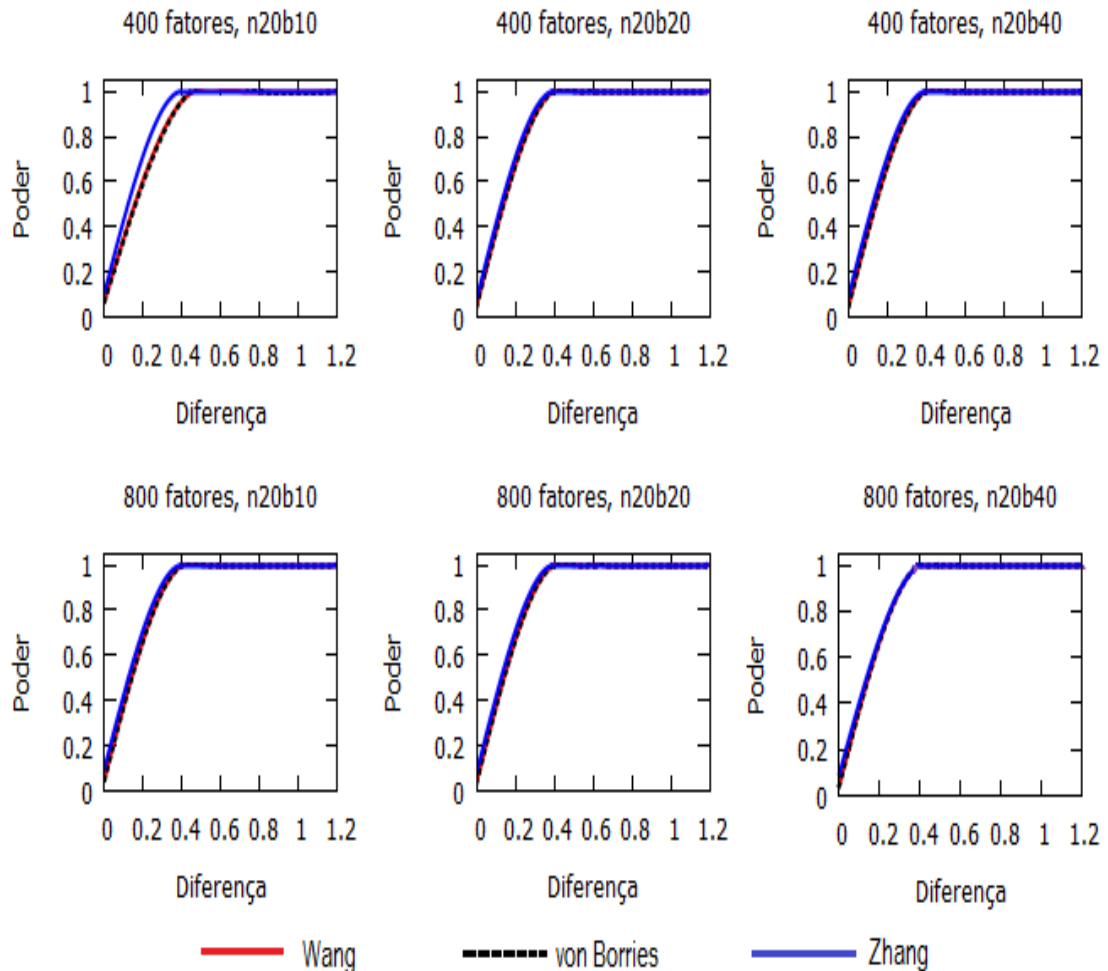


Figura 4.6: Poder do teste segundo diferentes metodologias.

Inicialmente, foram construídas superfícies de poder para reduzir a quantidade de gráficos produzidos devido à grande quantidade de parâmetros envolvidos. Entretanto, a visualização mostrou-se insatisfatória e, conseqüentemente, optou-se pela utilização de gráficos bidimensionais.

Nas figuras são avaliadas as situações com 5% (gráficos da primeira linha) e 10% (gráficos da segunda linha) dos fatores deslocados. De modo que, o maior percentual de contaminação detecta melhor as diferenças para um número menor de pontos no tempo.

Verifica-se ainda que à medida que o número de fatores e repetições por fator aumenta, a detecção de pequenas diferenças melhora. Tal que, esse comportamento é observado em todas as metodologias.

De modo geral, o teste de Zhang (2008) apresenta melhores resultados se com-

parado aos demais, ou seja, precisa de um deslocamento menor para detectar as diferenças nas mais diversas situações ilustradas nas figuras em questão.

É importante destacar que por mais que os testes de Wang (2004) e von Borries (2008) tenham sido descritos pelos autores como adequados para situações totalmente distintas, esse comportamento não foi visualizado nessas simulações.

Pelo contrário, para todos os níveis de fator, os testes de efeito simples apresentados pelos autores tiveram comportamento semelhante. E por esse motivo, a curva do teste de von Borries (2008) foi construída com a linha tracejada.

Ressalta-se que os gráficos de poder do teste foram feitos no software Gnuplot versão 4.6. O ambiente MULTIPLOT foi utilizado na construção dos gráficos em questão. E que os resultados das simulações realizadas encontram-se devidamente apresentados no Anexo A.

## 4.2 Resultados em dados reais

### 4.2.1 Eletroencefalografia

O aumento na capacidade de processamento dos computadores tem possibilitado a utilização de ferramentas poderosas de análise de sinais biopotenciais<sup>1</sup>.

Um exemplo são as pesquisas com interfaces cérebro-computador<sup>2</sup> que tornam-se cada vez mais frequentes. É possível, por exemplo, utilizar sinais relacionados ao controle de movimentos no comando de dispositivos eletromecânicos de assistência a indivíduos portadores de deficiências físicas.

Nesse contexto, o registro da atividade elétrica cerebral via eletroencefalografia (EEG) surge como uma técnica não-invasiva capaz de caracterizar processos fisiológicos relacionados a uma grande variedade de disfunções do sistema nervoso central. E torna-se, portanto, parte integrante do processo de controle das interfaces mencionadas.

Basicamente, a EEG é realizada através da colocação de eletrodos sobre o couro cabeludo (figura 4.7). Um gel condutor é empregado ao longo do processo para fixá-los

---

<sup>1</sup>Medida da eletricidade emitida por determinado órgão ao realizar certa atividade.

<sup>2</sup>São sistemas computacionais capazes de controlar dispositivos externos através da utilização de sinais provenientes da atividade cerebral.

e melhorar a captura dos sinais elétricos da atividade cerebral.

Além disso, os sinais registrados podem ser obtidos através da aplicação ou não de estímulos. No primeiro caso, os sinais são denominados Potenciais Evocados (Evoked Potentials - EPs) e surgem involuntariamente quando o indivíduo é exposto a estímulos externos. Já no segundo, o sinal é controlado voluntariamente pelo indivíduo e caracterizado por diferentes faixas de frequência e estados mentais. Maiores detalhes podem ser vistos em (SÖRNMO; LAGUNA, 2005, p.34).

Cinco voluntários em plenas condições físicas e mentais foram selecionados pelo MSPL<sup>1</sup> e expostos a diferentes estímulos ao longo do processo de captura do sinal de EEG. Especificamente, foram utilizados onze estímulos visuais (figura 4.8) e três sonoros, os quais foram apresentados aos indivíduos de forma completamente aleatória. Foram tomadas medidas da cabeça de cada participante para o correto posicionamento de uma touca com 128 eletrodos adotada na captura do sinal.

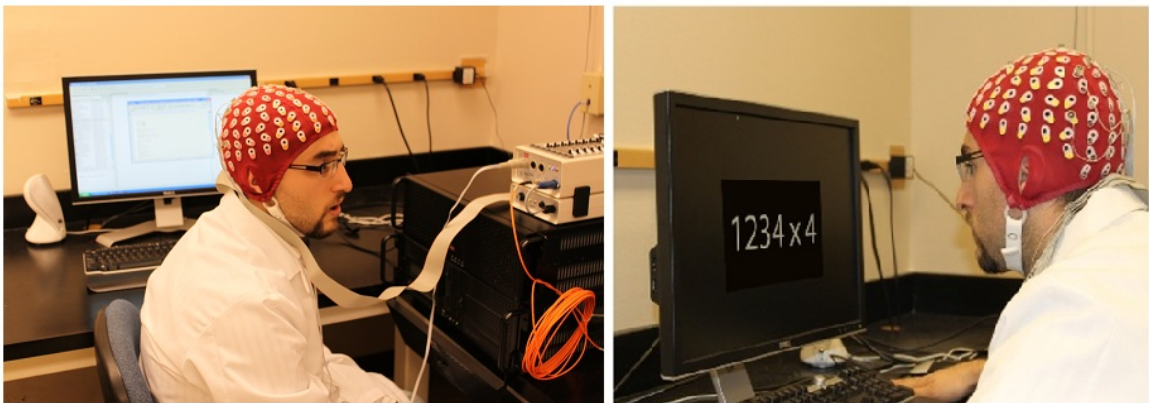


Figura 4.7: Coleta de dados de EEG no MSPL com o uso de estímulos visuais.

Os estímulos foram repetidos quatro vezes para cada indivíduo e o sinal elétrico registrado por quatro segundos. A transformada de Fourier (*Fourier Transform* - FT) foi utilizada para definir a frequência máxima ( $f_{max}$ ) do sinal captado. Assim, a taxa de reamostragem foi determinada conforme indicado abaixo:

$$T_x \leq \frac{1}{2f_{max}} \quad (4.5)$$

<sup>1</sup>Laboratório *Multi-Sensing-Processing and Learning* (MSPL) do Departamento de Engenharia Elétrica da Universidade do Texas em El Paso (UTEP).

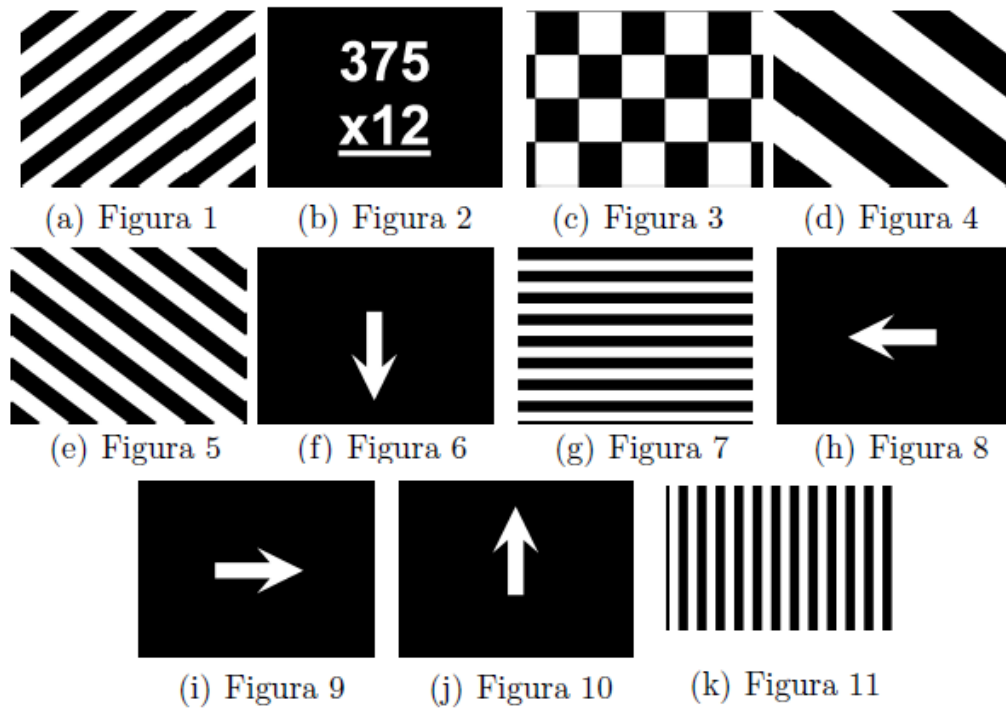


Figura 4.8: Estímulos visuais adotados na coleta dos sinais de EEG no MSPL.

Segundo Frondana (2012), essa taxa corresponde ao intervalo de tempo adotado para registrar os pontos do sinal, o que possibilita a representação de um sinal contínuo por um conjunto discreto de pontos. Além disso, a  $f_{max}$  fornecida pela transformada determina o número mínimo de pontos necessários para representar o sinal sem que haja perda de informações essenciais.

Assim, diferentes bases de dados podem ser obtidas através da aplicação da taxa de reamostragem. Especificamente, foi utilizada uma base com  $b = 168$  pontos no tempo. Os tratamentos são caracterizados por todas as possíveis combinações entre estímulos e eletrodos. A estrutura dos dados é apresentada na tabela 4.2, com exceção das repetições aleatórias dos estímulos. Destaca-se que os estímulos são independentes entre si e que a mesma suposição é considerada para os eletrodos.

A proposta de estudo consiste na aplicação de algoritmos de agrupamento em um conjunto de sinais de EEG. Supondo que os procedimentos aplicados sejam eficazes em identificar o real padrão dos grupos obtidos, torna-se possível, então, classificar um sinal desconhecido em alguma das classes previamente estabelecidas. De modo que, isso pode ser conduzido por meio de uma amostra de treinamento com o acréscimo desse sinal.

Tabela 4.2: Estrutura da base de dados de EEG.

Níveis de fator	Indivíduo	Ponto no tempo			
		$t_1$	$t_2$	...	$t_b$
Estímulo 1 - Eletrodo 1	1	$Y_{[1]11}$	$Y_{[1]21}$	...	$Y_{[1]b1}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	5	$Y_{[1]15}$	$Y_{[1]25}$	...	$Y_{[1]b5}$
Estímulo 1 - Eletrodo 2	1	$Y_{[2]11}$	$Y_{[2]21}$	...	$Y_{[2]b1}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	5	$Y_{[2]15}$	$Y_{[2]25}$	...	$Y_{[2]b5}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
Estímulo 1 - Eletrodo 128	1	$Y_{[128]11}$	$Y_{[128]21}$	...	$Y_{[128]b1}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	5	$Y_{[128]15}$	$Y_{[128]25}$	...	$Y_{[128]b5}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
Estímulo 14 - Eletrodo 1	1	$Y_{[1665]11}$	$Y_{[1665]21}$	...	$Y_{[1665]b1}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	5	$Y_{[1665]15}$	$Y_{[1665]25}$	...	$Y_{[1665]b5}$
Estímulo 14 - Eletrodo 2	1	$Y_{[1666]11}$	$Y_{[1666]21}$	...	$Y_{[1666]b1}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	5	$Y_{[1666]15}$	$Y_{[1666]25}$	...	$Y_{[1666]b5}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
Estímulo 14 - Eletrodo 128	1	$Y_{[1792]11}$	$Y_{[1792]21}$	...	$Y_{[1792]b1}$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
	5	$Y_{[1792]15}$	$Y_{[1792]25}$	...	$Y_{[1792]b5}$

Em uma primeira análise, Silva (2012) verificou, para diferentes limiares, que o algoritmo PPCLUSTEL-M<sup>1</sup> sempre retornava um subgrupo de cinco estímulos que se destacava dos demais: som com frequência de 100 Hz, som com frequência de 2000

<sup>1</sup>Algoritmo de agrupamento criado por Silva (2012) com base no teste de ausência de efeito simples de Wang (2004). A estimação da matriz de covariância é feita pelo método  $\Sigma_G$ , o qual considera que um grupo de variáveis possuem um mesmo comportamento ao longo do tempo.



Hz, figura 1, figura 5 e figura 7. Na figura 4.8 do trabalho do autor são ilustrados os sinais de EEG de cada um desses estímulos e é possível verificar a completa ausência de padrão dos sinais.

Assim, com o intuito de verificar se os algoritmos PPCLUSTEL-RG, PPCLUSTEL-RW e PPCLUSTE-RZ conseguem detectar diferenças entre os estímulos citados, apenas o subgrupo encontrado por Silva (2012) foi considerado nos processos de agrupamento.

Esses algoritmos incorporam a correlação existente entre os pontos no tempo e utilizam o p-valor resultante do teste de ausência de efeito simples como medida de similaridade.

A tabela 4.3 mostra os resultados que foram obtidos com o PPCLUSTEL-RW. Esse algoritmo conseguiu detectar bem três dos estímulos em questão, o som de 100 Hz (estímulo 1), a figura 1 (estímulo 3) e a figura 7 (estímulo 5). E apesar de ter utilizado o método de estimação  $\Sigma$ , obteve os mesmos resultados que o método  $\Sigma_G$  de Silva (2012).

A tabela 4.4 mostra os resultados obtidos com a aplicação do PPCLUSTEL-RG. Destaca-se que esse algoritmo detectou os mesmos grupos do PPCLUSTEL-RW.

É importante ressaltar que o algoritmo PPCLUSTEL-RZ não conseguiu realizar o agrupamento, já que ao final do processo todos os sinais estavam classificados no grupo 0. Como o comportamento do teste adotado não foi simulado para um número tão elevado de pontos no tempo, a suspeita é que o teste de ausência de efeito simples desse autor não funcione nas circunstâncias consideradas.

Tabela 4.3: Resultado do agrupamento com o algoritmo PPCLUSTEL-RW. Foi considerado o limiar  $\alpha=0,1$ .

Freq Linha(%) Coluna(%)	Grupos					Total
	1	2	3	4	5	
Estímulo 1	103,00	0,00	0,00	11,00	14,00	128
	80,87	0,00	0,00	8,52	10,94	
	93,64	0,00	0,00	4,58	12,28	
Estímulo 2	4,00	68,00	7,00	49,00	0,00	128
	3,13	53,13	5,47	38,28	0,00	
	3,64	53,13	14,58	20,42	0,00	
Estímulo 3	0,00	25,00	0,00	103,00	0,00	128
	0,00	19,53	0,00	80,47	0,00	
	0,00	19,53	0,00	42,92	0,00	
Estímulo 4	0,00	33,00	39,00	52,00	4,00	128
	0,00	25,78	30,47	40,63	3,13	
	0,00	25,78	81,25	21,67	3,51	
Estímulo 5	3,00	2,00	2,00	25,00	96,00	128
	2,34	1,56	1,56	19,53	75,00	
	2,73	1,56	4,17	10,42	84,21	

Tabela 4.4: Resultado do agrupamento com o algoritmo PPCLUSTEL-RW. Foi considerado o limiar  $\alpha=0,1$ .

Freq Linha(%) Coluna(%)	Grupos					Total
	1	2	3	4	5	
Estímulo 1	0,00	112,00	2,00	9,00	5,00	128
	0,00	87,50	1,56	7,03	3,91	
	0,00	97,39	1,59	3,83	4,35	
Estímulo 2	6,00	3,00	72,00	46,00	1,00	128
	4,69	2,34	56,25	35,94	0,78	
	12,24	2,61	57,14	19,57	0,87	
Estímulo 3	2,00	0,00	23,00	103,00	0,00	128
	1,56	0,00	17,97	80,47	0,00	
	4,08	0,00	18,25	43,43	0,00	
Estímulo 4	41,00	0,00	28,00	52,00	7,00	128
	32,03	0,00	21,88	40,63	5,47	
	83,67	0,00	22,22	22,13	6,09	
Estímulo 5	0,00	0,00	1,00	25,00	102,00	128
	0,00	0,00	0,78	19,53	79,69	
	0,00	0,00	0,79	10,64	88,70	

### 4.2.2 Microarranjo

Dados de expressão gênica são organizados em uma matriz onde cada linha corresponde a um gene e cada coluna a uma condição. As condições podem ser representadas por exemplo por pontos no tempo ou tecidos.

A análise de microarranjo permite que tais expressões sejam simultaneamente avaliadas e apresenta como principais objetivos o agrupamento e a classificação de genes de acordo com a sua expressão e a expressão de outros genes com padrão conhecido.

Em um experimento típico, são observados milhares ou dezenas de milhares de genes, quantidade que extrapola consideravelmente o número de indivíduos dos quais esses dados são coletados.

Em termos estatísticos, essa extrapolação constitui um grande problema, uma vez que o número de variáveis observadas (genes) é bem maior que o tamanho de amostra disponível.

Silva (2012) trabalhou com um exemplo apresentado em Gillespie et al. (2011), no qual são observados 10.928 genes de leveduras em 5 pontos no tempo, com 3 replicações e 2 tipos de condições experimentais: cepas de leveduras sem acréscimo de temperatura e cepas com alteração de temperatura.

O estudo apresentou como objetivo a identificação de genes de levedura com expressão gênica semelhante ao longo do tempo. Especificamente, as expressões das leveduras foram amostradas inicialmente a 23 °C e, então, 1, 2, 3 e 4 horas após, para um acréscimo de 7 °C nessa temperatura.

Além disso, foram consideradas cepas de leveduras dos tipo selvagem e com mutação. De modo que, as expressões gênicas das leveduras com mutação foram padronizadas pela mediana das expressões das cepas do tipo selvagem.

Para efetuar a análise, Silva (2012) não trabalhou com a redução da dimensionalidade dos dados conforme feito por Gillespie et al. (2011). O autor utilizou os algoritmos de agrupamento PPCLUSTEL-R e PPCLUSTEL-M, baseados nos testes de ausência de efeito simples de von Borries (2008) e Wang (2004) com uma forma distinta de estimação da matriz de covariância.

Aqui são considerados os mesmos dados utilizados por Silva (2012) e o mesmo

método de padronização da variável resposta.

Os algoritmos PPCLUSTEL-RG, PPCLUSTEL-RW e PPCLUSTEL-RZ foram aplicados com os limiares  $10^{-1}$ ,  $10^{-2}$ , ...,  $10^{-8}$ . De modo que, o padrão dos grupos passou a ser o mesmo a partir de  $10^{-4}$ ,  $10^{-4}$  e  $10^{-2}$ , respectivamente.

Para os algoritmos PPCLUSTEL-RG e PPCLUSTEL-RW a partir de  $10^{-4}$  sempre apareciam três grupos predominantes, o maior deles com aproximadamente 60% dos genes e os outros dois menores com cerca de 20%. Para o algoritmo PPCLUSTEL-RZ a partir de  $10^{-2}$  também passaram a aparecer três grupos de maior prevalência. Entretanto o maior deles era constituído por aproximadamente 40% dos genes e os menores por 25%.

Na execução dos algoritmos, os grupos foram identificados pelos números 0, 1, 2 e assim sucessivamente. Os codificados com os menores valores foram formados pelos genes de menor expressão, os maiores foram constituídos pelos genes de maior expressão e os intermediários pelos que não se expressaram. Para facilitar a visualização, os grupos centrais contendo os genes que não se expressaram ao longo do tempo foram retirados dos gráficos abaixo. Os *heatmaps* podem ser vistos no Anexo B.

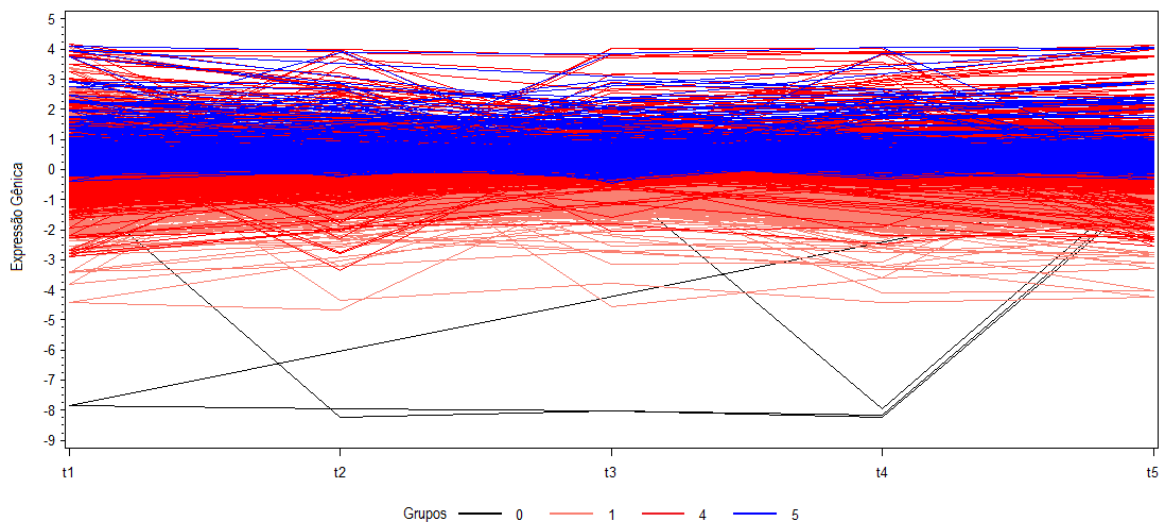


Figura 4.9: PPCLUSTEL-RG com limiar  $\alpha = 10^{-4}$ .

Observa-se que os algoritmos PPCLUSTEL-RG e PPCLUSTEL-RW separaram os genes com expressões acima e abaixo de zero. De modo que, esse comportamento não foi perceptível na aplicação do PPCLUSTEL-RZ. Esse algoritmo parece ter separado

os genes nos grupos com expressões entre -1 e 1 e com resposta fora desse intervalo.

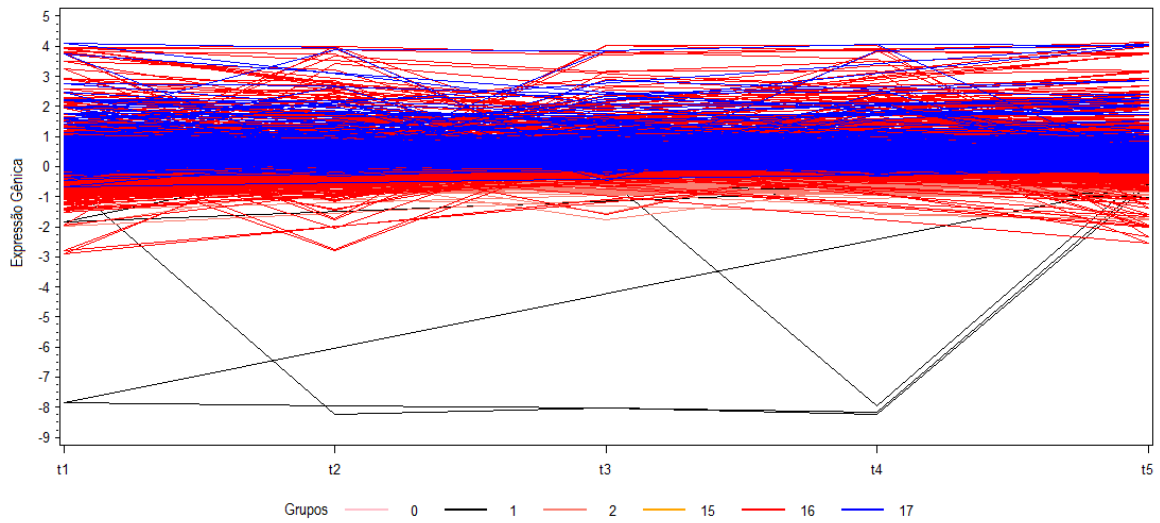


Figura 4.10: PPCLUSTEL-RW com limiar  $\alpha = 10^{-4}$ .

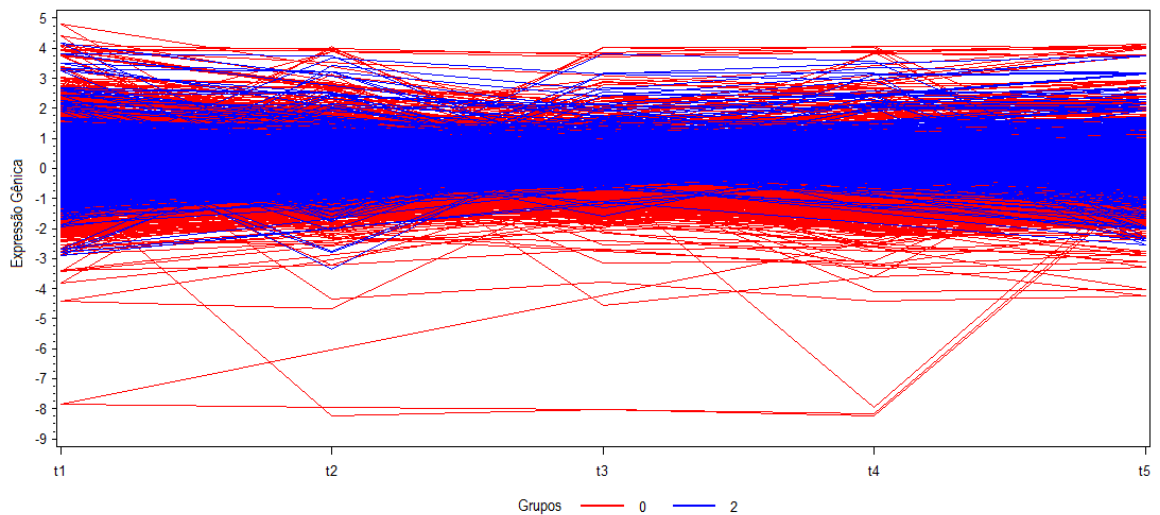


Figura 4.11: PPCLUSTEL-RZ com limiar  $\alpha = 10^{-8}$ .

A avaliação do desempenho dos algoritmos nos conjuntos de dados selecionados foi feita por meio de um índice de validação externa denominado Índice de Rand Corrigido (*Adjusted Rand Index* - ARI). Uma macro para calcular o ARI foi implementada por von Borries (2008) e foi adotada na obtenção de resultados de qualidade dos agrupamentos.

Entre os agrupamentos PPCLUSTEL-RG e PPCLUSTEL-RW o ARI obtido foi igual a 0,683. Já entre os algoritmos de von Borries/Zhang e Wang/Zhang foram

obtidos índices iguais a 0,326 e 0,319, respectivamente.

Os índices obtidos para os algoritmos PPCLUSTEL-RG e PPCLUSTEL-RW foram próximos de 1, o que indica que as partições formadas por esses algoritmos foram similares. Já os valores obtidos para o PPCLUSTEL-RZ em comparação ao demais algoritmos foram próximos de 0, o que mostra que o nível de semelhança da partição formada pelo algoritmo de Zhang foi menor em relação às partições formadas por von Borries e Wang.

### 4.2.3 Eletromiografia

Músculos esqueléticos são constituídos por fibras musculares que se contraem quando estimuladas por um motoneurônio. Motoneurônios, por sua vez, se originam na medula espinhal e apresentam em sua estrutura um único axônio. Ao conjunto composto por um motoneurônio, seu respectivo axônio e todas as fibras musculares inervadas a ele, dá-se o nome de unidade motora (HUSSAIN et al., 2006).

Assim, em condições normais, em resposta a estimulação de um motoneurônio, surge um campo elétrico próximo a cada fibra muscular inervada. A eletromiografia (EMG) é comumente adotada na captura dos sinais emitidos por esse campo.

Essa técnica corresponde basicamente ao registro da atividade elétrica das membranas fibromusculares em resposta à ativação fisiológica dos músculos. E, portanto, é uma importante ferramenta em pesquisas de áreas como: fisioterapia, medicina, educação física e terapia ocupacional.

Segundo Hussain et al. (2006), a combinação dos potenciais elétricos das fibras musculares de uma simples unidade motora corresponde ao Potencial de Ação da Unidade Motora (*Motor Unit Action Potential* - MUAP), o qual pode ser detectado por um eletrodo não-invasivo ou invasivo. De modo que, um sinal mioelétrico é formado por um conjunto de MUAPs, ou seja, é uma série temporal dos estímulos à unidade motora.

Para a realização da eletromiografia em músculos superficiais são adotados eletrodos de superfície (não-invasivos), uma vez que não causam desconforto durante a coleta dos sinais. Já para a eletromiografia de músculos profundos, é necessária a utilização de eletrodos invasivos para evitar a interferência dos músculos superficiais

na coleta dos dados.

Após coletado, o sinal é processado para que possa ser interpretado. Duas importantes características desse sinal são a amplitude e a frequência. De modo que, a amplitude representa a magnitude da atividade muscular e surge como consequência de variações na atividade das unidades motoras.

Em geral, entre os fatores capazes de influenciar sinais mioelétricos, destacam-se os seguintes: aquisição do sinal EMG de músculos vizinhos (*cross-talk*) e problemas com artefatos eletromecânicos (por exemplo, movimentação do equipamento). De forma que, tais fatores podem provocar o surgimento de ruídos, ou seja, sinais não desejados ao longo do sinal de EMG captado.

Correia (2007) trabalha com a análise dados de eletromiografia (EMG) do Laboratório *Multi-Sensing-Processing and Learning* (MSPL). Aqui é considerada uma base do mesmo Laboratório, composta por 25 indivíduos, 14 homens e 11 mulheres, com informações referentes à avaliação de sinais mioelétricos dos músculos *splenius capitis*, *sternocleidomastoid* e *trapezius* (figura 4.12). Os dados disponíveis totalizam cerca de 80GB de informação.

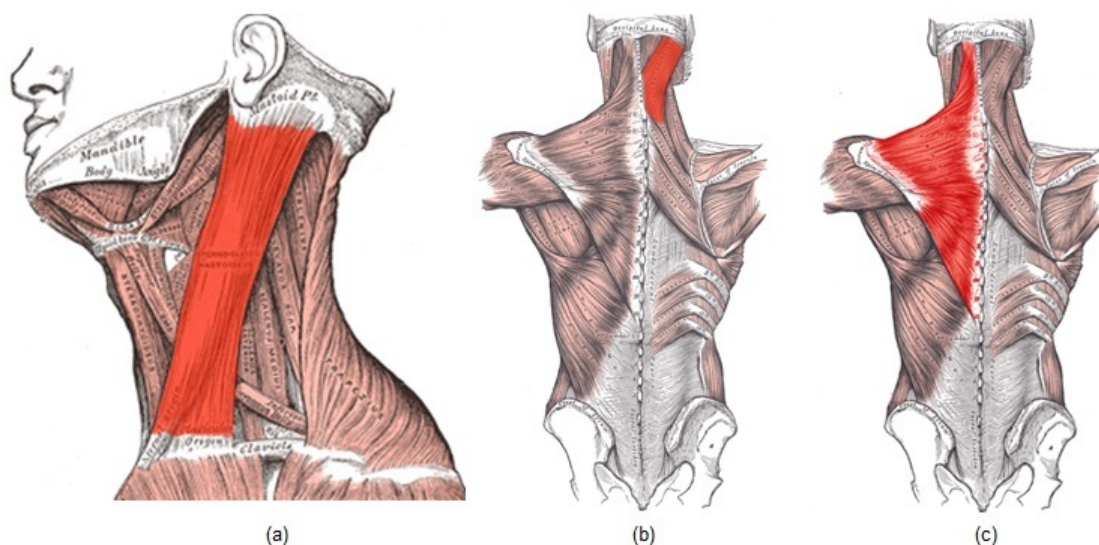


Figura 4.12: Músculos envolvidos na análise de fadiga muscular. (a) *Sternocleidomastoid* - Ação: flexionar e rotacionar lateralmente a espinha cervical. (b) *Splenius capitis* - Ação: estender e rotacionar a espinha cervical. (c) *Trapezius* - Ação: estabilizar, elevar, retrain e rotacionar a escápula. (Imagens disponíveis em: <<http://en.wikipedia.org/wiki/>> Acesso em: 06/04/2012).



Com o objetivo de identificar padrões de fadiga<sup>1</sup> nos músculos citados, foram captadas informações de amplitude e frequência do sinal mioelétrico. De modo que, essas mensurações foram realizadas a cada milissegundo em um período de 8 horas (H0-H7).

Os indivíduos foram expostos a diferentes cargas de peso (A,B,C,D,E) e as medidas foram obtidas para os lados esquerdo e direito de cada músculo. A estrutura geral dos dados é apresentada na tabela 4.5, onde as letras E e D representam os lados esquerdo e direito, respectivamente.

Tabela 4.5: Estrutura da base de dados de EMG.

Sexo	Indivíduo	<i>Sternocleidomastoid</i>			<i>Splenius capitis</i>			<i>Trapezius</i>		
		E/D			E/D			E/D		
		H0	...	H7	H0	...	H7	H0	...	H7
Masculino	1	-	-	-	-	-	-	-	-	-
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
	14	-	-	-	-	-	-	-	-	-
Feminino	1	-	-	-	-	-	-	-	-	-
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
	11	-	-	-	-	-	-	-	-	-

De acordo com a descrição apresentada anteriormente, o número de pontos no tempo excede, e muito, a quantidade de indivíduos dos quais as medidas foram tomadas.

Assim, para aplicar a metodologia proposta por Brunner e Puri (2001), os tratamentos foram representados pelos músculos e cada hora (H1-H7) passou a ser representada por um único valor que correspondeu à mediana das amplitudes registradas por milissegundo em cada um dos lados. De modo que, a média aritmética entre os dois lados foi tomada como a medida final de atividade muscular. Além disso, dos 25 indivíduos, foram considerados apenas 9 homens e 9 mulheres que haviam sido submetidos a uma das cargas de peso.

<sup>1</sup>Falha na manutenção de um nível desejado de rendimento ou trabalho durante uma atividade repetitiva ou sustentada.

Para a obtenção dos resultados de interesse, foram adotados diferentes procedimentos do programa SAS. Primeiro, definiu-se a variável resposta, baseada em postos, por meio do procedimento RANK. E após a especificação dos postos, os dados foram ordenados segundo repetições e pontos no tempo para a execução do procedimento MIXED.

Para analisar dados longitudinais com o PROC MIXED foi preciso especificar todos os fatores de parcela e subparcela no comando MODEL, bem como a estrutura da matriz de covariância no comando REPEATED por meio da opção TYPE=. O comando SUB= foi adotado para definir as observações relativas a cada indivíduo. De modo que, a suposição de independência entre os vetores de observações foi assumida.

Além disso, através da opção GROUP= foi possível listar os níveis de fator ou combinações desses níveis em que diferentes matrizes de covariância são permitidas, o que aumenta consideravelmente o número de parâmetros a serem estimados.

A opção ANOVAF foi adicionada na linha de comando do PROC MIXED para computar a estatística tipo ANOVA e os p-valores resultantes. Para evitar problemas computacionais pela utilização do método de máxima verossimilhança restrita (REML), na estimação da matriz de covariância não estruturada, Brunner e Puri (2001) recomendaram o uso do método de estimação quadrática não-viesada de variância mínima pela adição da opção METHOD=MIVQUE0 logo após o comando ANOVAF.

A estatística tipo Wald e os respectivos p-valores puderam ser obtidos pelo acréscimo da opção /CHISQ na linha do comando MODEL.

Os resultados obtidos para os testes de ausência do efeito do fator músculo e da interação entre músculo e tempo são apresentados na tabela abaixo.

Tabela 4.6: P-valores dos testes de ausência de efeitos principais e de interação.

Fonte de Variação	WTS			ATS			
	$Q_n(*)$	d.f.	p-valor	$F_n(*)$	$\hat{f}_1$	$\hat{f}_0$	p-valor
Músculo (A)	13,78	2	0,001	6,65	1,98	50,2	0,0013
Tempo (B)	36,61	6	<.0001	3,59	3,85	146	0,0069
Músculo*Tempo (C)	11,91	12	0,4532	0,53	6,57	146	0,8036

Conforme mencionado, foram observados dezoito indivíduos para cada tratamento. Todavia, há indícios que dezoito não seja um número considerável de repetições para a adoção da estatística WTS (AHMAD, 2008). Por esse motivo, são analisados apenas os p-valores obtidos para a *ANOVA-type Statistics*.

De acordo com a tabela 4.6, a hipótese nula de ausência do efeito da interação não é rejeitada ao nível de significância de 5% ( $\hat{\alpha}=0,8036>0,05$ ), o que faz com que os p-valores obtidos para os efeito principais tenham sentido.

Assim, há argumentos estatísticos contra a hipótese inicial de que os efeitos dos músculos sejam iguais ( $\hat{\alpha}=0,0013$ ) ao nível de significância considerado. Logo, conclui-se que a atividade muscular de pelo menos um dos músculos apresenta comportamento diferente dos demais.

Além da análise segundo a metodologia de Brunner e Puri (2001), foi conduzida também a análise de agrupamento com os algoritmos PPCLUSTEL-RG, PPCLUSTEL-RW e PPCLUSTEL-RZ. Para isso, foram considerados seis tratamentos, resultantes das combinações entre os níveis dos fatores músculo e lado, e uma base de dados com 56 pontos no tempo, já que foram definidos oito percentis<sup>2</sup> para cada hora (H1-H7).

Na execução de todos os algoritmos foram formados dois grupos. O PPCLUSTEL-RG e o PPCLUSTEL-RW retornaram o mesmo resultado, ambos com um grupo constituído pelos tratamentos *splenius* esquerdo, *splenius* direito, *trapezius* esquerdo e *trapezius* direito. Já no PPCLUSTEL-RZ, um dos grupos foi formado pelos tratamentos *splenius* direito, *splenius* esquerdo e *trapezius* esquerdo. Os limiares adotados foram iguais a  $10^{-8}$ ,  $10^{-8}$  e  $10^{-11}$ , respectivamente.

---

<sup>2</sup>P1, P5, P10, P25, P50, P75, P90 e P99

# Capítulo 5

## Conclusões

As técnicas tradicionais de análise de dados longitudinais mostram-se falhas no estudo de estruturas assintóticas. Um exemplo são os problemas que surgem na especificação das estatísticas dos testes de ausência de efeito de tempo e tratamento quando o número de graus de liberdade dessas fontes de variação tende a infinito.

Recentemente, diferentes metodologias foram propostas com o intuito de solucionar tal problema. Entre elas está o trabalho de Wang (2004), no qual são realizadas inferências sobre os efeitos dos níveis de fator, pontos no tempo e interação entre ambos nos casos em que  $b \rightarrow \infty$ ,  $n_i$  é pequeno ou grande e  $a$  é fixo.

Destacam-se, da mesma forma, os testes apresentados por Zhang (2008) e von Borries (2008), os quais são descritos como adequados para lidar com estruturas de dados que possuem  $a$  relativamente grande ou  $a \rightarrow \infty$ ,  $b$  fixo e  $n_i$  pequeno.

Evidencia-se que, ao contrário do que é exigido pelos testes tradicionais de análise de variância, os testes assintóticos desenvolvidos pelos autores mencionados não requerem que as medidas repetidas no tempo sejam contínuas ou homocedásticas. Entretanto, apesar dessa característica, esses testes apresentam como desvantagem a subjetividade dos termos pequeno, fixo e infinito adotados em suas definições.

Este trabalho buscou avaliar o desempenho desses testes em cenários de simulação previamente estabelecidos. As verificações foram conduzidas sob diferentes circunstâncias com o intuito de melhor compreender a dimensão dos termos utilizados.

Foram revisados e comparados numericamente testes de ausência de efeito simples para que posteriormente pudessem ser utilizados na construção dos algoritmos de agrupamento PPCLUSTEL-RW e PPCLUSTEL-RZ. O algoritmo PPCLUSTEL-RG

serviu de base para a implementação dos dois últimos e pode ser visualizado em von Borries (2008).

Conforme descrito em Silva (2012), as curvas de poder do teste de fato apresentaram melhores resultados para as distribuições normal e para 10% de fatores contaminados na média. Todavia, à medida que o número de fatores ou réplicas por tratamento aumentou, as curvas das distribuições normal e t multivariadas se aproximaram. Ressalta-se que o mesmo comportamento foi visto para os diferentes percentuais de contaminação adotados.

Em relação aos testes citados, a análise das curvas de poder mostrou que as abordagens de Wang (2004) e von Borries (2008) apresentaram comportamento muito semelhante em todos os cenários que foram investigados. Salienta-se a importância de tal resultado visto que, inicialmente, ambos os testes foram descritos como adequados para a análise de situações com características totalmente distintas. O teste de von Borries (2008) foi superior em apenas alguns casos e com uma diferença no valor do poder do teste na terceira casa decimal.

Obteve-se resultado igualmente relevante a partir do teste de ausência de efeito simples construído com base nos testes de ausência de efeito de tratamento e interação definidos por Zhang (2008). Esse teste obteve melhor desempenho em todos os cenários considerados. Mesmo naqueles com quarenta pontos no tempo, nos quais se esperava um melhor comportamento do teste desenvolvido por Wang (2004).

Portanto, conclui-se que os testes desenvolvidos por von Borries (2008) e Wang (2004) parecem não depender dos intervalos de convergência para os quais foram definidos. Fatores como a distribuição dos dados e a qualidade de estimação da matriz de covariância parecem exercer maior influência no desempenho dessas abordagens.

No entanto, é preciso que as simulações sejam conduzidas para outros cenários. Principalmente, para o teste construído com base no trabalho de Zhang (2008), já que ele não retornou resultados na aplicação em dados de eletroencefalografia. Nesse caso, especificamente, a desconfiança é que outros limiares deveriam ter sido testados. Foram considerados apenas os limiares de  $10^{-1}$ ,  $10^{-2}$ , ...,  $10^{-8}$ .

É importante ressaltar que não foi possível trabalhar com computadores de iguais capacidades de processamento na realização das simulações. Entretanto, mesmo assim, notou-se que o tempo de execução foi diretamente influenciado pela estimação

da matriz de covariância dos dados. E que, além disso, o método jackknife mostra-se extremamente dispendioso computacionalmente, principalmente quando o número de observações disponíveis para a estimação dos parâmetros é elevado.

Na execução dos algoritmos de agrupamento PPCLUSTEL-RG, PPCLUSTEL-RW e PPCLUSTEL-RZ foi utilizado o mesmo computador. De todos os procedimentos, o PPCLUSTEL-RG foi o mais rápido, seguido do PPCLUSTEL-RW. O algoritmo desenvolvido com base no trabalho de Zhang (2008) falhou ao ser executado. O software estatístico SAS versão 9.2 informou que a quantidade de memória disponível<sup>1</sup> era insuficiente para a realização do agrupamento. Desse modo, foi preciso utilizar um computador com mais memória para que o algoritmo funcionasse.

## 5.1 Estudos futuros

Abaixo são apresentadas algumas sugestões para pesquisas futuras.

- Todos os testes foram comparados para uma única forma de estimação da matriz de covariância dos dados. Seria interessante trabalhar com outros métodos de estimação e verificar o ganho que se tem com tais alterações;
- Conforme mencionado, é importante que os testes sejam simulados para outros cenários. Por exemplo, para situações que envolvam um número extremamente elevado de tratamentos e pontos no tempo. Surgiram dúvidas sobre a real capacidade do algoritmo PPCLUSTEL-RZ, já que ele não conseguiu classificar os tratamentos envolvidos no experimento de eletroencefalografia;
- Verificar o desempenho dos testes não-paramétricos de ausência de efeito simples em outras distribuições contínuas e discretas;
- No caso do estudo com dados de microarranjo, os resultados obtidos para o ARI apenas indicaram que os resultados dos algoritmos PPCLUSTEL-RG e PPCLUSTEL-RW são similares. Assim, não é possível saber se os agrupamentos formados estão corretos. Um estudo dos grupos formados por esses algoritmos (em termos de função genética) poderia sugerir qual deles está no caminho certo.

---

<sup>1</sup>Foi utilizado um computador com processador Intel(R)Core(TM)2 Duo e memória RAM de 3GB.

# Referências Bibliográficas

- [1] AHMAD, R. M. Analysis of high dimensional repeated measures designs: the one- and two-sample test statistics. Dissertação, University of Göttingen, 2008.
- [2] AKRITAS, M.G.; ARNOLD, S. Fully nonparametric hypothesis for factorial designs i: Multivariate repeated measures designs. *Journal of the American Statistical Association*, 89:336–343, 1994.
- [3] ARAÚJO, T. C. *Agrupamento de Sinais de Eletroencefalografia por Mistura de Distribuições Normais*. Monografia (Departamento de Estatística), Universidade de Brasília, 2010.
- [4] BOX, G.E.P. Some theorems on quadratic forms applied in the study of analysis of variance problems i. effects of inequality of variance in the one-way classification. *Annals of the Mathematical Statistics*, 25:290–302, 1954.
- [5] BRUNNER. E.; PURI, M.L. Nonparametric methods in factorial designs. *Statistical Papers*, 42:1–52, 2001.
- [6] CAPELLI, L. P.; NASCIMENTO, R. M. P. *O mapa da mina: entendendo o mapeamento gênico. Genética na escola*. 2008.
- [7] CORREIA, L. T. *Análise de Resposta Muscular por Eletromiografia utilizando Medidas Repetidas/ Dados Longitudinais*. Monografia (Departamento de Estatística), Universidade de Brasília, 2007.
- [8] DAVIDIAN, M.; GALLANT, A. R. The nonlinear mixed effects model with a smooth random effects density. *Biometrika*, 80:475–488.
- [9] DAVIDIAN, M.; GILTINAN, D. M. *Nonlinear models for repeated measurement data*. New York: Chapman & Hall., 1995.

- [10] DAVIS, C. S. *Statistical Methods for the Analysis of Repeated Measurements*. Springer, 2002.
- [11] DIGGLE, P.J; LIANG, K-Y; ZEGER, S. L. *Analysis of Longitudinal Data*. Oxford University Press, 1994.
- [12] DONG, L. Nonparametric tests for longitudinal data. Dissertação (mestrado em estatística), Iowa State University, 2005.
- [13] Músculos envolvidos na análise de fadiga muscular. Disponível em: <<http://en.wikipedia.org/wiki>>. Acesso em 06/04/2012.
- [14] FAN, J.; LIN, S-K. Test of significance when data are curves. *UC Los Angeles: Department of Statistics, UCLA.*, 1998.
- [15] FARAWAY, J. J. Regression analysis for a functional response. *Technometrics*, 39:254–261, 1997.
- [16] FITZMAURICE, G.; DAVIDIAN, M.; VERBEKE, G.; MOLENBERGHS, G. *Longitudinal Data Analysis*. CRC Press, 2008.
- [17] FRONDANA, I. M. Classificação de biopotenciais: via cadeias de markov ocultas. Dissertação (Mestrado em Estatística), Universidade de Brasília, 2012.
- [18] GILLESPIE, C. S.; LEI, G.; BOYS, R. J.; GREENALL, A.; WILKINSON, D. J. Analysing time course microarray data using bioconductor: a case study using yeast2 affymetrix arrays. *Oxford Bioinformatics*, page 3:81, 2011.
- [19] HAND, D.; CROWDER, M. *Practical Longitudinal Data Analysis*. Chapman e Hall, 1996.
- [20] HEDEKER, D. E GIBBONS, R. D. *Longitudinal Data Analysis*. J. Wiley & Sons, New York., 2006.
- [21] HUSSAIN, M.S.; REAZ, M. B. I.; MOHD-YASIN, F.; M.I. IBRAHIMY. Electromyography signal analysis using wavelet transform and higher order statistics to determine muscle contraction. *Expert Syst.*, 26(1):35–48, 2009.



- [22] HUYNH, H.; FELDT, L.S. Conditions under which mean square ratios in repeated measurements designs have exact f-distributions. *Journal of the American Statistical Association*, v.65:p.1582–1589, 1970.
- [23] JOHNSON, R. A.; WICHERN, D. W. *Applied Multivariate Statistical Analysis*. New Jersey: Prentice Hall, 2007.
- [24] KUEHL, R. O. *Diseño de Experimentos*. Thomson Learning, 2001.
- [25] LAIRD, N. M.; WARE, J. H. Random effects models for longitudinal data. *Biometrics*, 38:963–974, 1982.
- [26] LIANG, K.-Y.; ZEGER, S. L. Longitudinal data analysis using generalized linear models. *Biometrika*, 73:13–22, 1986.
- [27] MALLET, A. A maximum likelihood estimation method for random coefficient regression models. *Biomet*, 73:645–656, 1986.
- [28] MALUF, S.; RIEGEL, M. *Citogenética humana*. 2011.
- [29] MAUCHLY, J.W. Significance test for sphericity of a normal n-variate distribution. *Annals of Mathematical Statistics*, 11:204–209, 1940.
- [30] MILLIKEN, G. A.; JOHNSON, D. E. *Analysis of Messy Data*, volume 1: Designed Experiments. New York: Chapman & Hall., 2009.
- [31] PARIZZI, F. C. *Incidência de Fungos da Pré-colheita ao Armazenamento de Café*. Tese (Doutorado em Engenharia Agrícola), Universidade Federal de Viçosa, 2005.
- [32] PAWITAN, Y. *In all likelihood: statistical modeling and inference using likelihood*. Oxford, 2001.
- [33] PINHEIRO, J.C.; BATES, D.M. Approximations to the loglikelihood function in the nonlinear mixed effects model. *Journal of Computational and Graphical Statistics*, 4:12–35, 1995.
- [34] RAMSAY, J.O.; SILVERMAN, B.W. *Applied Functional Data Analysis: Methods and Case Studies*. Springer, 2002.

- [35] RENCHER, A. C. *Methods of Multivariate Analysis*. Wiley-Interscience, 2002.
- [36] RICE, J. A. Functional and longitudinal data analysis: Perspectives on smoothing. *Statistica Sinica*, 14:631–647, 2004.
- [37] ROCHA, L. C. S. *Análise de Componentes Independentes em Sinais de Eletroencefalografia*. Monografia (Departamento de Estatística), Universidade de Brasília, 2011.
- [38] SILVA, A. P. T. Implementação, análise e aplicação de algoritmos de agrupamento de dados superdimensionados, longitudinais e com amostras pequenas. Dissertação (Mestrado em Estatística), Universidade de Brasília, 2012.
- [39] SILVA, M. C. *Um novo método para classificação de dados de EEG: uma aplicação de Máquinas de Suporte Vetorial e Análise de Fourier*. Monografia (Departamento de Estatística), Universidade de Brasília, 2011.
- [40] SMITH, A.F.M.; ROBERTS, G.O. Bayesian computation via the gibbs sampler and related markov chain monte carlo methods. *Journal of the Royal Statistical Statistical Society*, 55:3–23, 1993.
- [41] SORNMO, L.; LAGUNA, P. *Bioelectrical Signal Processing in Cardiac and Neurological Application*. Elsevier Academic Press, 2005.
- [42] SOUZA, C. P. E. Testes de hipóteses para dados funcionais baseados em distâncias: um estudo usando splines. Master's thesis, Universidade Estadual de Campinas, 2008.
- [43] TEIXEIRA, E. C. M. *Identificação de Padrões em Fadiga Muscular*. Monografia (Departamento de Estatística), Universidade de Brasília - UnB, 2011.
- [44] VON BORRIES, G. *Partition Clustering of High Dimensional Low Sample Size Data based on p-values*. Tese (Doutorado em Estatística), Kansas State University, 2008.
- [45] VONESH, E. F. CHINCHILLI, V. M. *Linear and nonlinear models for the analysis of repeated measurements*. 1997.

- [46] WANG, H. *Testing in Multifactor Heteroscedastic ANOVA and Repeated Measures Designs with Large Number of Levels*. Tese (Doutorado em Estatística), The Pennsylvania State University, 2004.
- [47] WANG, H.; NEILL, J.; MILLER F. Nonparametric clustering of functional data. *Statistics and Its Inference*, 1:47–62, 2008.
- [48] WANG, L; AKRITAS, M. G. Two-way heteroscedastic anova when the number of levels is large. *Statistica Sinica*, 16:1387–1408, 2006.
- [49] XAVIER, L. H. Modelos univariado e multivariado para análise de medidas repetidas e verificação da acurácia do modelo univariado por meio de simulação. Master's thesis, Escola Superior de Agricultura Luiz de Queiroz, Universidade de São Paulo, 2000.
- [50] ZHANG, K. *Inference of Nonparametric Hypothesis Testing on High Dimensional Longitudinal Data and its Application in DNA Copy Number Variation and Microarray Data Analysis*. Tese (Doutorado em Estatística), Kansas State University, 2008.
- [51] ZHANG, K., WANG, H., BATHKE, A.C., HARRAR, S.W., PIEPHO, H., DENG, Y. Gene set analysis for longitudinal gene expression data. *BMC Bioinformatics*, 12:273, 2011.

# Apêndice A

## Poder do teste

Resultados das simulações do poder dos testes de ausência de efeitos simples para os cenários considerados no Capítulo 4.





























# Apêndice B

## *Heatmaps*

*Heatmaps* construídos para os dados de microarranjo analisados no Capítulo 4.

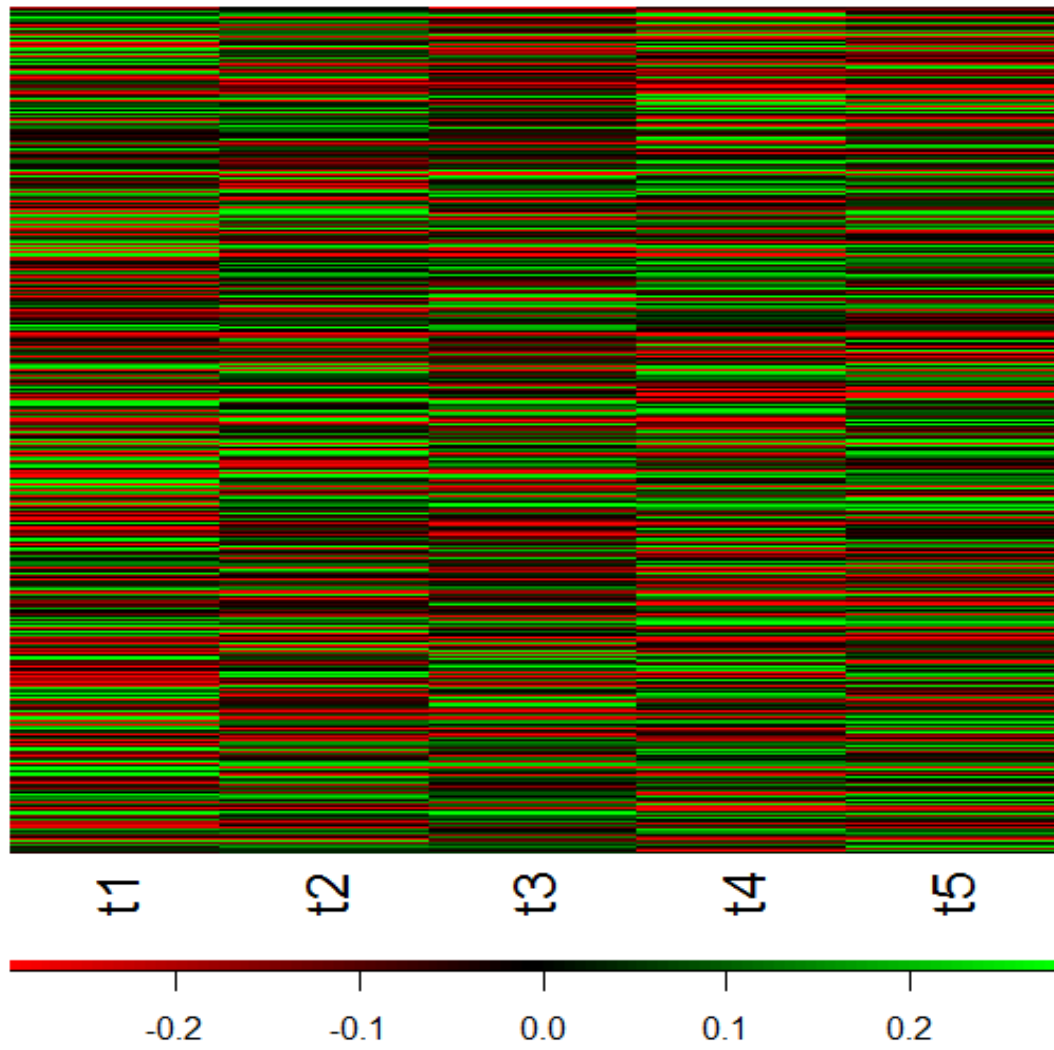


Figura B.1: Dados não ordenados

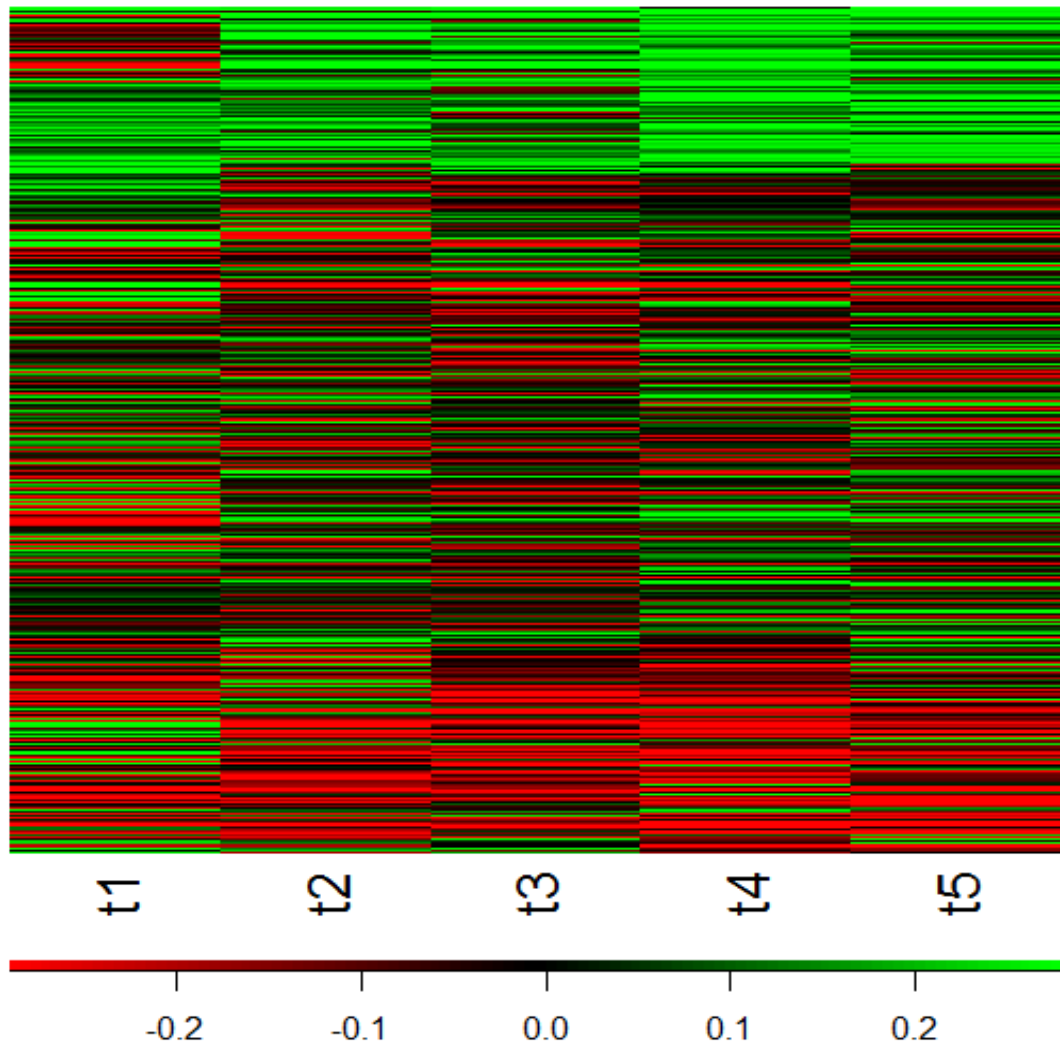


Figura B.2: PPCLUSTEL-RG com limiar  $\alpha = 10^{-4}$

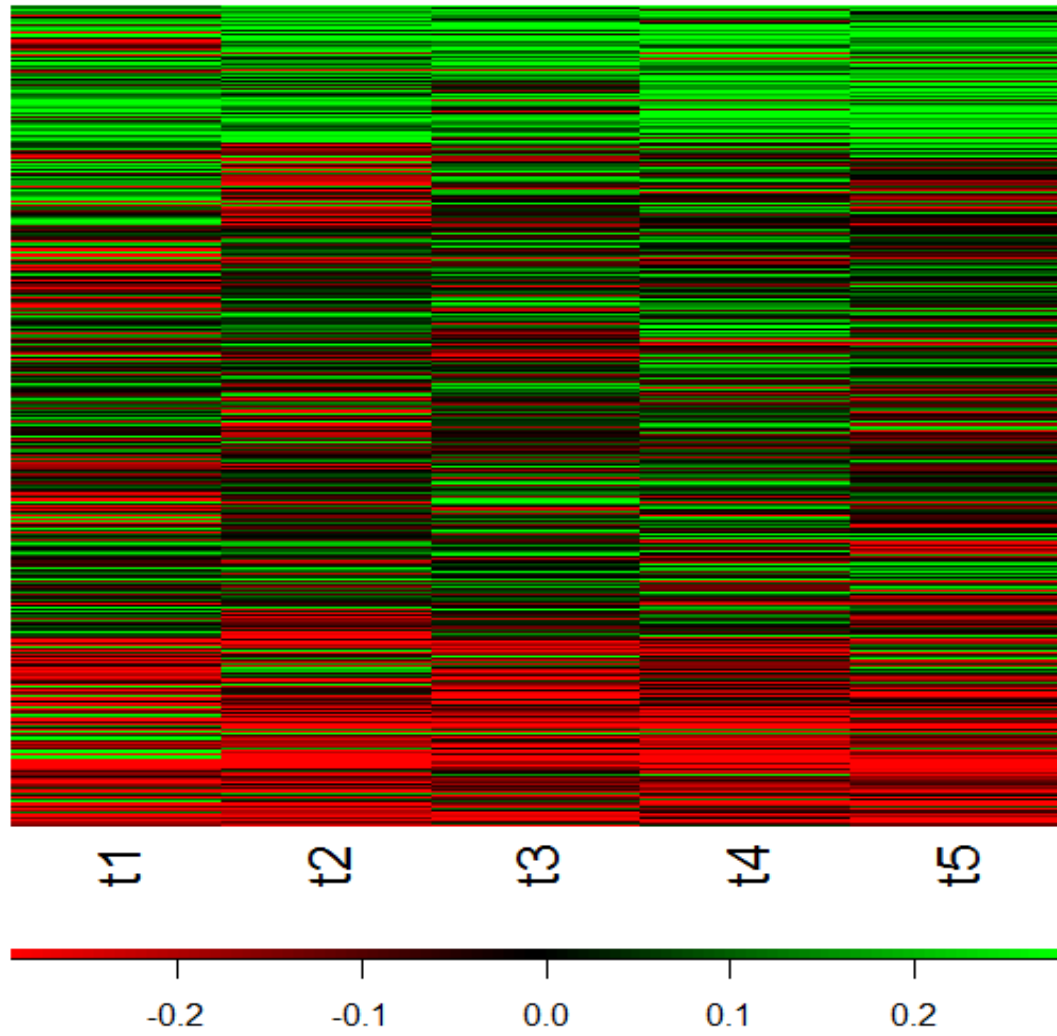


Figura B.3: PPCLUSTEL-RW com limiar  $\alpha = 10^{-4}$

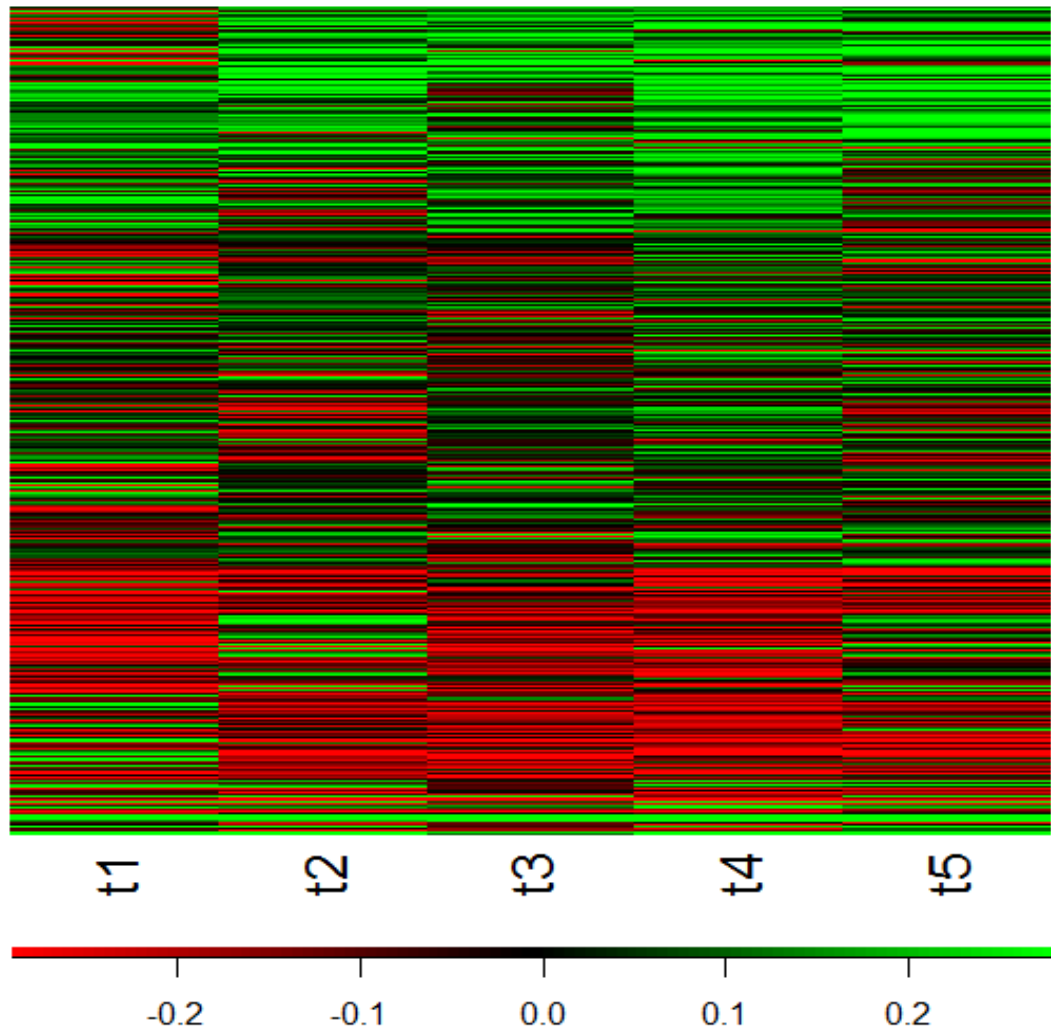


Figura B.4: PPCLUSTEL-RZ com limiar  $\alpha = 10^{-8}$

# Apêndice C

## Programações em SAS

Todas as programações utilizadas nas simulações de poder do teste e nos algoritmos de agrupamento deste trabalho estão disponíveis mediante contato com a autora pelo e-mail [thaysa.gsouza@gmail.com](mailto:thaysa.gsouza@gmail.com).