

UNIVERSIDADE ABERTA



Dissertação de Mestrado

Mestrado em Estatística, Matemática e Computação

Modelos de resposta à interação entre fármacos anestésicos: Análise de Regressão e Análise de Clusters

Ana Maria Figueiredo Nascimento Lopes dos Santos

dezembro de 2015

UNIVERSIDADE ABERTA



Dissertação de Mestrado

Modelos de resposta à interação entre fármacos anestésicos: Análise de Regressão e Análise de Clusters

Mestranda: Ana Maria Figueiredo Nascimento Lopes dos Santos

Orientadora: Catarina Sofia da Costa Nunes Duarte

Dissertação submetida á Universidade Aberta
para obtenção do grau de
Mestre em Estatística, Matemática e Computação

dezembro de 2015

Resumo

Os avanços tecnológicos e científicos, na área da saúde, têm vindo a aliar áreas como a Medicina e a Matemática, cabendo à ciência adequar de forma mais eficaz os meios de investigação, diagnóstico, monitorização e terapêutica. Os métodos desenvolvidos e os estudos apresentados nesta dissertação resultam da necessidade de encontrar respostas e soluções para os diferentes desafios identificados na área da anestesia. A índole destes problemas conduz, necessariamente, à aplicação, adaptação e conjugação de diferentes métodos e modelos das diversas áreas da matemática.

A capacidade para induzir a anestesia em pacientes, de forma segura e confiável, conduz a uma enorme variedade de situações que devem ser levadas em conta, exigindo, por isso, intensivos estudos. Assim, métodos e modelos de previsão, que permitam uma melhor personalização da dosagem a administrar ao paciente e por monitorizar, o efeito induzido pela administração de cada fármaco, com sinais mais fiáveis, são fundamentais para a investigação e progresso neste campo.

Neste contexto, com o objetivo de clarificar a utilização em estudos na área da anestesia de um ajustado tratamento estatístico, proponho-me abordar diferentes análises estatísticas para desenvolver um modelo de previsão sobre a resposta cerebral a dois fármacos durante sedação. Dados obtidos de voluntários serão utilizados para estudar a interação farmacodinâmica entre dois fármacos anestésicos. Numa primeira fase são explorados modelos de regressão lineares que permitam modelar o efeito dos fármacos no sinal cerebral BIS (índice bispectral do EEG – indicador da profundidade de anestesia); ou seja estimar o efeito que as concentrações de fármacos têm na depressão do eletroencefalograma (avaliada pelo BIS). Na segunda fase deste trabalho, pretende-se a identificação de diferentes interações com Análise de Clusters bem como a validação do respetivo modelo com Análise Discriminante, identificando grupos homogêneos na amostra obtida através das técnicas de agrupamento. O número de grupos existentes na amostra foi, numa fase exploratória, obtido pelas técnicas de agrupamento hierárquicas, e a caracterização dos grupos identificados foi obtida pelas técnicas de agrupamento k-means. A reprodutibilidade dos modelos de agrupamento obtidos foi testada através da análise discriminante.

As principais conclusões apontam que o teste de significância da equação de Regressão Linear indicou que o modelo é altamente significativo. As variáveis propofol e remifentanil influenciam significativamente o BIS e o modelo melhora com a inclusão do remifentanil.

Este trabalho demonstra ainda ser possível construir um modelo que permite agrupar as concentrações dos fármacos, com base no efeito no sinal cerebral BIS, com o apoio de técnicas de agrupamento e discriminantes. Os resultados desmontram claramente a interação farmacodinâmica dos dois fármacos, quando analisamos o Cluster 1 e o Cluster 3. Para concentrações semelhantes de propofol o efeito no BIS é claramente diferente dependendo da grandeza da concentração de remifentanil.

Em suma, o estudo demonstra claramente, que quando o remifentanil é administrado com o propofol (um hipnótico) o efeito deste último é potenciado, levando o sinal BIS a valores bastante baixos

Palavras-chave: Regressão linear; Modelos de resposta; Clustering; Anestesiologia.

Abstract

Mathematics has been playing an important role in the technological and scientific developments in the health area. When the areas of Medicine and Mathematics are combined science is most effective in linking research, diagnosis, monitoring and therapeutics. The developed methods and studies presented in this dissertation are a result of the search for solutions to different challenges identified in the area of anaesthesia. The nature of these problems leads, necessarily, to the development, adaptation and conjugation of diverse methods and models in the different areas of mathematics.

Induction of anaesthesia in patients, in a safe and reliable way, leads to a huge variety of situations that must be taken into account; therefore there is a demand for intensive studies. Methods and models of foreknowledge are crucial to research and improvement in this field, so as to allow for patient dosage's adaptation. Models may be used to help the clinician predict the individual drug dose required to induced a desired effect.

In this context, the aim is to develop a foreknowledge model towards the brain effect of two drugs during sedation. To this purpose statistical analysis will be used. Data obtained from volunteers will be used to study the pharmacodynamics interaction between the two anaesthetic drugs. In the first phase, linear regression models are explored, which allow to model the effect of the drugs on the brain signal BIS (bispectral index of the EEG – measure of depth of anaesthesia); that is to model of the drugs' concentration on the central nervous systems depression (as assess by BIS). In the second phase of this work, the different drug interactions are identified by means of Clusters analysis, as well as the validation of the corresponding model through Discriminative Analysis, identifying homogeneous groups in the obtained sample, through clustering techniques. On an exploratory phase, the number of groups in the sample was determined through hierarchical clustering and the characterization of the identified groups was defined using the k-means clustering. The reproducibility of the achieved clustering models was tested through discriminative analysis.

The main conclusions are that Linear Regression model is highly meaningful to estimate the effect of the hypnotic and analgesic drug on the brain signal BIS. propofol and remifentanil anaesthetic drugs influence BIS substantially, and the model improves with the inclusion of the remifentanil concentration.

This research also shows that it is possible to build a model with the support of mathematical techniques (clustering and discriminating); that allows the clustering of drugs concentration based on its' effect on the brain signal BIS. The results when we analyse the different clusters, clearly show the pharmacodynamics interaction of both drugs. For similar propofol concentrations, the effect on BIS is totally different and dependent on the level of the remifentanil concentration.

Keywords: Linear regression; Response models; Anaesthesia; Clustering.

Agradecimentos

O meu agradecimento às várias pessoas que de diferentes formas contribuíram prontamente para o desenvolvimento desta dissertação.

À minha orientadora, Professora Doutora Catarina Sofia da Costa Nunes Duarte, pelo muito que me ensinou, pelos contributos, apoio e confiança sempre presentes e sempre renovados, desde a génese do projeto. Expreço ainda o paradigma, enquanto pessoa e mestre.

À Universidade Aberta e aos meus professores, a instituição onde muito aprendi.

Ao meu marido, filhos e familiares que foram sempre companheiros de jornada, no apoio e incentivo.

Conteúdo

Resumo.....	i
Abstract.....	iii
Agradecimentos.....	v
Conteúdo.....	vii
Lista de Figuras.....	ix
Lista de Tabelas.....	ix
Capítulo 1.....	1
Introdução.....	1
1.1. Motivação.....	3
1.3.Estrutura da dissertação.....	4
Capítulo 2.....	5
Anestesia.....	5
2.1. Anestesia.....	7
2.2. Profundidade da anestesia (<i>Depth of Anesthesia - DoA</i>).....	8
2.3. Modelos.....	10
2.3.1 Modelo Farmacocinético.....	10
2.4 Target Controlled Infusion (TCI).....	12
Capítulo 3.....	15
Caracterização dos dados.....	15
3.1. Descrição da base de dados.....	17
3.2. Caracterização das variáveis anestésicas.....	20
Capítulo 4.....	25
Análise de Regressão.....	25
4.1. Introdução.....	27
4.2. Regressão e Correlação Linear.....	27
4.2.1 Testes de Hipóteses sobre o Coeficiente de Correlação.....	27
4.3. Modelo de Regressão Linear Simples.....	31
4.3.1. Modelo Teórico.....	31
4.3.2. Pressupostos do Modelo.....	32
4.3.3. Estimação dos Parâmetros do Modelo.....	33
4.3.3.1.Método dos Mínimos Quadrados.....	33
4.3.4. Teste de Hipóteses e Intervalos de Confiança para os Parâmetros do Modelo.....	37
4.3.4.1. Teste e intervalos de confiança para β_1	37

4.3.4.2. Teste e Intervalos de Confiança para β_0	38
4.4. Modelo de Regressão Linear Múltipla	41
4.4.1. Modelo Teórico e seus pressupostos	41
4.4.1.1. Interações	42
4.4.1.2. Pressupostos do modelo	42
4.4.1.3. Representação matricial do método de regressão linear múltipla.....	43
4.4.2. Estimação do Parâmetro do modelo.....	44
4.4.3. Análise da Variância (ANOVA) Aplicada à Regressão Linear Múltipla	45
4.5. Análise de Resíduos.....	48
4.5.1. Diagnóstico de Normalidade	48
4.5.2. Diagnóstico de Homoscedasticidade (Variâncias constantes)	50
4.5.3. Diagnóstico de Independência	51
4.5.4. Diagnóstico de Outliers e Observações Influentes.....	52
4.5.4.1. Observações Influentes	53
4.5.5. Colinearidade e Multicolinearidade	54
4.5.6. Análise de Variância Multivariada (MANOVA) aplicada à Regressão	56
4.6. Construção do Modelo de Regressão Linear Simples.....	59
4.7. Construção do Modelo de Regressão Múltipla.....	65
Capítulo 5	72
Análise de Clusters.....	72
5.1. Introdução.....	74
5.2. As variáveis	75
5.2.1. Seleção das variáveis.....	75
5.2.3. Escala das variáveis.....	76
5.3. Medidas de semelhança e medidas de dissemelhança	77
5.3.3. Medidas de semelhança e medidas de dissemelhança entre sujeitos	77
5.3.3.1. Dissemelhanças e distâncias – propriedades	77
5.3.3.2. Semelhanças – propriedades	78
5.3.3.3. Medidas de dissemelhança e de semelhança para variáveis quantitativas	78
5.3.3.4. Medidas de dissemelhança e de semelhança para variáveis qualitativas	80
5.3.3.4.1. Medidas de semelhança para variáveis nominais binárias. Medidas de associação ..80	
5.3.3.4.2. Medidas de semelhança para variáveis nominais com mais de dois níveis	82
5.3.3.4.3. Medidas de semelhança para variáveis ordinais	84
5.3.3.5. Coeficiente de semelhança para variáveis de diferentes tipos	84
5.3.3.6. Conversão das semelhanças em dissemelhanças	85
5.3.4. Medidas de semelhança entre variáveis	86
5.3.4.1. Medidas de semelhança entre variáveis quantitativas.....	86
5.3.4.2. Medidas de semelhança entre variáveis nominais binárias.....	87
5.3.4.3. Medidas de semelhança entre variáveis nominais com mais de dois níveis	87

5.3.4.4 Medida de semelhança entre variáveis ordinais	89
5.4. Métodos hierárquicos.....	89
5.4.1. Métodos de (des) agregação - Características	90
5.5. Escolha do número de <i>Clusters</i>	92
5.5.1. Análise do dendrograma.....	92
5.5.2. Coeficiente de fusão	93
5.6. Métodos não hierárquicos	93
5.7. Outros métodos.....	94
5.7.1. <i>TwoStep Cluster</i> (Análise de <i>Clusters</i> em duas fases)	94
5.8. Métodos hierárquicos / métodos não hierárquicos.....	96
5.9. Escolha da técnica a utilizar	97
5.10. Análise Discriminante	98
5.10.1. Critério Discriminante.....	99
5.10.2. Correlações Canónicas	100
5.10.3 Testes de Significância.....	102
5.11. Identificação de grupos homogéneos	104
5.12. <i>Clusters k-means</i> variáveis estandardizadas	115
5.13. <i>Clusters k-means</i> por Amplitude.....	120
5.14. Caracterização dos Clusters	124
Capítulo 6	125
6.1. Conclusões.....	127
Anexos A.....	135
Anexos B	136
Anexos C.....	143

Lista de Figuras

Figura 2.1: Monitor e sensor BIS. (Fotografias cedidas por Aspect Medical Systems™)

Figura 2.2: Escala do índice bispectral BIS. (Figura cedida por Aspect Medical Systems™)

BIS: índice bispectral; EEG: electroencefalograma

Figura 2.3: Modelo Farmacocinético Compartmentado

Figura 2.4: Farmacocinético e Farmacodinâmico

Figura 2.5: Diagrama de recolha de dados dos voluntários

Figura 3.2.1: Distribuição das variáveis (o eixo vertical representa a frequência, ou seja, o número de casos para os respectivos valores da variável)

Figura 3.2.2: Diagrama de extremos das variáveis propofol, remifentanil e BIS

Tabela 3.2.3: Voluntário 1 propofol Ce ($\mu\text{g.mL}$), remifentanil Ce (ng.mL), BIS e OAA/S

Figura 4.1: Classificação da correlação através do diagrama de dispersão

Figura 4.2: Interpretação geométrica dos parâmetros do modelo de regressão linear simples

Figura 4.3. - Representação gráfica dos resíduos

Figura 4.4: Hiperplano p-dimensional referente às variáveis explicativas

Figura 4.5: Normal p-p plot de resíduos

Figura 4.6: Confirmação da homoscedasticidade dos resíduos

Figura 4.7: Diagrama de dispersão

Figura 4.8: Gráficos dos Resíduos versus preditos; resíduos padronizados e da probabilidade normal dos resíduos

Figura 4.9: Q-Q plot

Figura 4.10: Gráficos dos Resíduos versus preditos; resíduos padronizados e da probabilidade normal dos resíduos

Figura 4.11: Gráficos dos Resíduos versus preditos; resíduos padronizados e da probabilidade normal dos resíduos (RLM)

Figura 4.12: Q-Q plot (RLM)

Figura 4.13: Superfície de Resposta ajustada entre remifentanil e propofol e BIS do modelo de 1º Ordem. e 2º Ordem

Figura 5.1: Dendrograma

Figura 5.2: Valor estandardizado (z-scores) da mediana para pelos 3 Clusters. O propofol Ce ($\mu\text{g/ml}$), remifentanil Ce($\eta\text{g/ml}$) e BIS

Figura 5.3: Função discriminantes para o agrupamento hierárquico com método de Ward's para 3 Clusters (z-scores)

Figura 5.4: Valor estandardizado (z-scores) da mediana para pelos 4 Clusters. O propofol Ce ($\mu\text{g/ml}$), remifentanil Ce($\eta\text{g/ml}$) e BIS

Figura 5.5: Valor real da mediana para pelos 4clusters. o propofol Ce (ug/ml), remifentanil Ce ($\eta\text{g/ml}$) e BIS

Figura 5.6: Representação das 2 principais funções discriminantes para o agrupamento hierárquico com método de Ward's para 4 Clusters

Figura 5.7: Valor estandardizado (z-scores) da mediana para pelos 5 clusters. o propofol Ce ($\mu\text{g/ml}$), remifentanil Ce($\eta\text{g/ml}$) e BIS

Figura 5.8: Valor real da mediana para pelos 5 Clusters. O propofol Ce (ug/ml), remifentanil Ce e BIS

Figura 5.9: Representação das 2 principais funções discriminantes para o agrupamento hierárquico com método de Ward's para 5 Clusters

Figura 5.10: Representação 3D dos 3 clusters Ward's em função dos valores de propofol Ce ($\mu\text{g/ml}$), remifentanil Ce ($\eta\text{g/ml}$) e BIS
Figura 5.9: Valor da mediana para propofol Ce (ug/ml), remifentanil Ce e BIS pelos 3 Clusters

Figura 5.11: Valor estandardizado (z-scores) da mediana para pelos 3 Clusters. o propofol Ce ($\mu\text{g/ml}$), remifentanil Ce($\eta\text{g/ml}$) e BIS

Figura 5.12: Funções discriminantes para o agrupamento por clusters kmeans

Figura 5.13: Valor por amplitude para propofol Ce ($\mu\text{g/ml}$), remifentanil Ce ($\eta\text{g/ml}$) e BIS pelos 3 Clusters.

Lista de Tabelas

Tabela 2.1: Escalas de Responsividade Observer's Assessment of Alertness/Sedation Scale-OAA/S

Tabela 3.2.1: Caracterização das variáveis anestésicas

Tabela 3.2.2: Estatísticas descritivas das características dos pacientes da base de dados da DoA induzida pela administração de propofol e remifentanil.

Tabela 4.1: Tabela da análise de variância (ANOVA)

Tabela 4.2: Tabela de decisão em função de d_L e d_U

Tabela 4.3: MANOVA – Teste Λ de Wilks

Tabela 4.4: Estatísticas descritivas

Tabela 4.5: -Resumo do modelo^b

Tabela 4.6: ANOVA^a (Variáveis preditoras: Constante)

Tabela 4.7: Coeficientes^a

Tabela 4.8: Testes de Normalidade

Tabela 4.9: Estatísticas de resíduos^a

Tabela 4.10: Verificação da multicolinearidade (Variáveis dependente: BIS)

Tabela 4.11: Diagnóstico de colinearidade (Variável Dependente: BIS)

Tabela 4.12: Variáveis Inseridas/Removidas^a

Tabela 4.13: Resumo do modelo

Tabela 4.14: ANOVA^a

Tabela 4.15: Coeficientes (Variáveis dependente: BIS)

Tabela 4.16: Teste de Kolmogorov-Smirnov de uma amostra

Tabela 4.17: Testes de Normalidade

Tabela 4.18: Verificação da multicolinearidade (Variáveis dependente: BIS)

Tabela 4.19: Diagnóstico de colinearidade (Variável Dependente: BIS)

Tabela 5. 1 Tabela de contingência

Tabela 5.2: Tabela de Contingência

Tabela 5.3: Reprodutibilidade do modelo de agrupamento para 3 Clusters

Tabela 5.4: Valores próprios e coeficientes canônicos das 2 funções discriminantes.

Tabela 5.5: Matriz de estrutura: coeficientes de correlação entre cada variável e as 2 funções discriminantes

Tabela 5.6: Coeficientes de funções discriminantes canônicas padronizados

Tabela 5.7: Resultados da classificação^a

Tabela 5.8 : Valores próprios

Tabela 5.9 : Matriz de estruturas

Tabela 5.10: Coeficientes de funções discriminantes canônicas padronizados

Tabela 5.11: Resultados da classificação^a

Tabela 5.12: Valores próprios

Tabela 5.13: Matriz de estruturas

Tabela 5.14: Coeficientes de funções discriminantes canônicas padronizados

Tabela 5.15: Número de iterações necessárias para a convergência da solução

Tabela 5.16: Distância Euclidiana entre os centros dos 3 Clusters

Tabela 5.17: Tabela de ANOVA

Tabela 5.18: Reprodutibilidade do modelo de agrupamento Clusters k-means

Tabela 5.19: Valores próprios e coeficientes canônicos das 2 funções discriminantes

Tabela 5.20: Matriz de estrutura

Tabela 5.21: Coeficientes canônicos das 2 funções discriminantes.

Tabela 5.22: Distância entre centros de Clusters finais

Tabela 5.23: Histórico de iteração^a

Tabela 5.18: Tabela de Resultados da classificação^a

Tabela 5.19: Valores próprios e coeficientes canônicos das 2 funções discriminantes

Tabela 5.19: Valores próprios e coeficientes canônicos das 2 funções discriminantes

Tabela 5.20: Matriz de estrutura

Tabela 5.21: Coeficientes de funções discriminantes canônicas padronizados

Capítulo 1

Introdução

1.1. Motivação

A anestesia geral é uma combinação de depressões de diferentes funções do sistema nervoso, incluindo a ocorrência do estado de inconsciência (amnésia e hipnose), relaxamento muscular, analgesia (ausência de dor) e controle dos reflexos simpáticos e parassimpáticos [1].

As dosagens de fármaco necessárias a administrar são, na maioria dos casos, ajustados em função da média populacional. Apesar da prática e conhecimentos dos médicos para ajustar as doses de fármacos necessárias, este procedimento não é adequado para a administração de uma dose de fármaco desejada individualizada. Sendo que uma das preocupações crescentes na área da anestesia geral está relacionada com a consciência durante a recuperação anestésica e as consequências que a anestesia geral pode ter a curto e longo prazo [2], o desenvolvimento de métodos mais eficazes para a administração de fármacos torna-se cada vez mais importante.

Com os avanços conseguidos nos últimos anos no campo de ferramentas como sensores e atuadores, assim como na modelação de sistemas biológicos, atingiu-se uma realidade onde é possível o controlo automático da anestesia. O desenvolvimento de métodos de administração automática e o controlo dos fármacos irão auxiliar o trabalho do anestesiológico, permitindo que este seja libertado de funções, podendo, assim, o seu foco estar apenas na supervisão do paciente, facilitando a capacidade de induzir a anestesia de forma mais segura e estável [3].

1.2 Objetivos

A presente dissertação tem como objetivos:

- Modelar o efeito dos fármacos no sinal cerebral BIS (índice bispectral do EEG – indicador da profundidade de anestesia);
- Modelar o sinal BIS a partir das concentrações dos fármacos;
- Estimar o efeito que as concentrações de fármacos têm na depressão do electroencefalograma (avaliada pelo BIS);
- Identificar se existem diferentes tipos de interações entre fármacos em relação ao seu efeito no sinal - BIS.

1.3.Estrutura da dissertação

A presente dissertação está dividida em seis capítulos.

O capítulo **1**, Introdução, é constituído pela motivação, objetivos, contribuições e estrutura da dissertação e tem como finalidade fazer a apresentação do trabalho.

O capítulo **2**, intitulado “Anestesia”, tem como finalidade dar uma perspetiva da anestesia geral e dos métodos utilizados para controlar manualmente a administração dos fármacos anestésicos na prática clínica.

O capítulo **3**, intitulado, “Caracterização dos dados”, é constituído por cinco subcapítulos que têm o propósito a caraterização e descrição dos dados recolhidos e a caraterização da amostra.

O capítulo **4** é denominado “Análise de Regressão Multivariada” Nos primeiros cinco subcapítulos abordarão a teoria da análise de regressão, procurando clarificar este conceito e através da revisão bibliográfica aprofundar conteúdos de interesse teórico para a fundamentação da metodologia adotada no nosso estudo. Os dois últimos subcapítulos são dedicados à aplicação da análise de regressão a um estudo que pretende modelar o efeito de fármacos anestésicos na prática clínica, onde será explicado o desenho do estudo, as variáveis selecionadas e os procedimentos que serão efetuados. Na análise dos dados, serão aplicados métodos de regressão a uma base de dados na área da Anestesiologia. Assim, procurar-se-á promover uma discussão a nível metodológico e dos próprios resultados do estudo.

O capítulo **5**, denominado, “Análise de Clusters”, é dividido em três subcapítulos. No primeiro subcapítulo abordaremos a teoria da análise Clusters e análise discriminante. No segundo subcapítulo são identificados grupos homogéneos de fármacos e o efeito no sinal cerebral BIS recorrendo a técnicas de agrupamento de dados (cluster analysis), confrontando os resultados. No terceiro subcapítulo são validados e caracterizados os grupos obtidos através de análise discriminante e técnicas estatísticas.

Por fim, o capítulo **6**, diz respeito às conclusões retiradas ao longo de todo o trabalho efetuado na dissertação e apresentação de sugestões para trabalhos futuros.

No anexo A, podem ser observadas as tabelas de correlação das variáveis estudadas. No anexo B, são apresentadas as saídas do SPSS nos modelos de regressão (Stepwise Forward,Enter). No anexo C, são apresentadas as saídas do SPSS da análise Clusters e Discriminante do estudo.

Capítulo 2

Anestesia

2.1. Anestesia

A anestesia pode ser definida como a falta de resposta a estímulos. Os fármacos utilizados na prática anestésica, com base no efeito fisiológico que induzem no paciente, são divididos em três categorias: hipnóticos, opióides e bloqueadores neuromusculares.

Os hipnóticos, classe de substâncias psicoativas, são administrados durante a cirurgia com a função de criar inconsciência e sedação ao paciente *Depth of Anesthesia* - DoA. Exemplos de hipnóticos são o isoflurano, o sevoflurano, e o benzodiazepínico (fármacos voláteis) e o propofol (fármaco injetável) [4, 5].

Opióides têm a capacidade de se ligarem a recetores específicos do sistema nervoso central (SNC) e do sistema nervoso periférico (SNP), levando à diminuição da reação e percepção da dor. Assim, os analgésicos são administrados com o intuito de suprimir a sensação de dor (analgesia). Exemplo deste tipo de fármaco é o remifentanil [4, 5].

Por último, a função dos bloqueadores neuromusculares (ou relaxantes musculares) é induzir o bloqueio neuromuscular (NMB) ou paralisia muscular. Deste último tipo de fármacos, os mais frequentemente utilizados na prática anestésica, fazem parte o atracurium e o rocuronium. Estes são agentes NMB não despolarizantes de duração intermédia e o seu efeito é medido numa escala de 0% a 100% e designado por nível de relaxamento neuromuscular ou nível de NMB [4, 5].

A quantidade de analgésico administrado é de extrema importância, pois não existe nenhum indicador claro do grau de dor. Os efeitos dos analgésicos e dos hipnóticos estão interligados, interagindo estes fármacos um com o outro de modo a atingir um nível adequado de profundidade da anestesia (DoA, do inglês *depth of anesthesia*) [5].

Os hipnóticos, propofol, e o opióide remifentanil são muito utilizados na sedação de pacientes adultos quando submetidos a intervenções assim como nos cuidados intensivos. São, atualmente, considerados os fármacos mais adequados para a anestesia totalmente intravenosa (TIVA) e para a indução e manutenção da DoA [31]. Os efeitos sedativos de um fármaco podem ser avaliados clinicamente através de escalas como a Observer's Assessment of Alertness/Sedation Scale-OAA/S, que classifica o nível de sedação do paciente em leve, moderada e profunda, de acordo

com a avaliação clínica da responsividade, fala, expressão facial e avaliação ocular dos pacientes, tabela 2.1 [7].

Tabela 2.1: Escala de sedação OAA/S

Score	Responsividade.
5	Resposta normal do nome.
4	Resposta lenta ao nome.
3	Responde somente quando chamado com tom de voz alto e repetidas vezes.
2	Responde após estímulo tátil.
1	Não responde ao estímulo tátil
0	Ausência de resposta ao estímulo

2.2. Profundidade da anestesia (*Depth of Anesthesia - DoA*)

Uma das principais funções do anestesiológico durante a cirurgia é controlar a profundidade da anestesia (DoA), no entanto, esta é difícil de medir com precisão, pois depende de vários fatores [8, 9, 10], como: o equilíbrio das concentrações de fármaco no plasma com concentrações do fármaco no local de efeito; a relação entre a concentração do fármaco e efeito do fármaco; a influência de estímulos nocivos. Na (figura 2.1) está representado o monitor e o sensor para controlar a profundidade da anestésica (DoA).

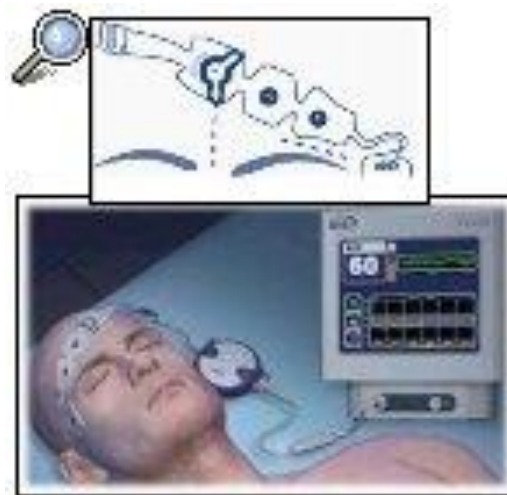


Figura 2.1: Monitor e sensor BIS. (Fotografias cedidas por Aspect Medical Systems™)

O efeito do hipnótico é avaliado através do índice bispectral (BIS) (figura 2.2) [11, 12, 13]. O índice bispectral (BIS) é utilizado para orientar e controlar a administração de hipnóticos e analgésicos. Na prática clínica, a pressão arterial e a frequência cardíaca também são utilizadas como orientações para determinar o nível de analgesia [6, 14, 8]. O sinal BIS é um parâmetro processado do EEG, utilizado como indicador da profundidade da hipnose, medindo o grau de depressão no sistema nervoso central [4, 6, 14]: um sujeito desperto tem um BIS de 100 enquanto um estado de ausência de atividade elétrica cerebral tem um BIS de zero (EEG isolelétrico), sendo geralmente de 97.7 o valor registrado antes do paciente ser anestesiado (na prática monitor nunca apresenta 100).

Num ambiente cirúrgico, o BIS deve ser mantido entre 40 e 60 para garantir um estado de anestesia adequado (anestesia geral) [4, 2].

Durante a cirurgia, se o valor do BIS é demasiado alto, o anestesiológista aumenta a dose de analgésicos e hipnótico com o intuito de aumentar a profundidade da anestesia, e se o valor do BIS é demasiado baixo, diminui essa dose.

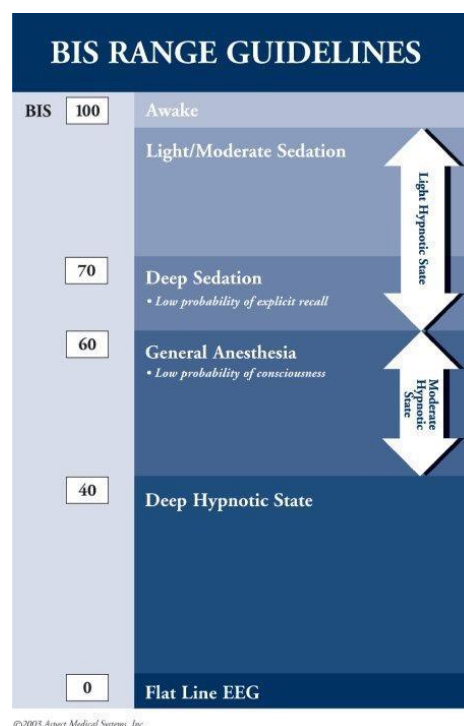


Figura 2.2: Escala do índice bispectral BIS. (Figura cedida por Aspect Medical Systems™). BIS: índice bispectral; EEG: electroencefalograma

2.3. Modelos

O efeito fisiológico induzido no paciente pela administração de uma dose, quantidade de fármaco por quilograma, depende, por exemplo, das características farmacocinéticas e farmacodinâmicas (PK/PD) [31].

A fase farmacocinética (PK) corresponde ao percurso do fármaco no organismo. Parâmetros como a biodisponibilidade (fração de fármaco que atinge a circulação numa forma inalterada e fica disponível para a circulação sistêmica), a ligação às proteínas plasmáticas, o volume de distribuição e a indução de enzimas afetam a farmacocinética do fármaco no organismo [15]. Interações farmacocinéticas podem surgir devido a alterações na absorção, distribuição, metabolismo ou eliminação do fármaco, sendo as interações que afetam a distribuição e o metabolismo do fármaco as mais importantes para os anestesiológicos [15, 16].

A fase farmacodinâmica (PD) corresponde à interação do fármaco com o local de ação conduzindo ao efeito farmacológico. Interações farmacodinâmicas são aquelas em que os efeitos de um fármaco são alterados pela presença de outro fármaco no seu local de ação ou quando a concentração de um fármaco é alterado, sendo que diferentes concentrações de fármacos vão produzir diferentes efeitos no corpo [15, 17].

2.3.1 Modelo Farmacocinético

O modelo farmacocinético dos fármacos anestésicos pode ser descrito como 3 compartimentos diferentes onde o fármaco propofol se dilui no corpo. Esta abordagem presume que cada compartimento tenha uma concentração uniforme do medicamento com permutas de fluxo com os outros compartimentos, como indicado na figura 2.3. Este modelo foi escolhido de forma a desenvolver o trabalho descrito em [18].

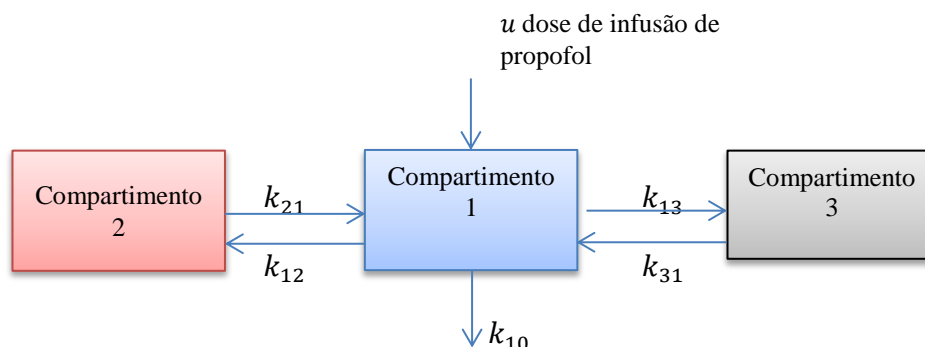


Figura 2.3: Modelo Farmacocinético Compartmentado

Este modelo pode ser expresso matematicamente pela equação (2.1).

$$\begin{cases} \dot{m}_1(t) = \frac{10^4}{3600} u(t) + k_{21}m_2(t) + k_{31}m_3(t) - k_{10}m_1(t) - k_{12}m_1(t) - k_{13}m_1(t) \\ \dot{m}_2(t) = k_{12}m_1(t) - k_{21}m_2(t) \\ \dot{m}_3(t) = k_{13}m_1(t) - k_{31}m_3(t) \end{cases} \quad (2.1.)$$

Aqui, a variável $u(t)$ [mL/h] é a taxa de infusão de propofol a ser normalizada para [$\mu g/s$] (assume-se que a diluição do propofol seja 10 mg/mL), $m_i(t)$ [μg] é a massa de propofol em cada um dos compartimentos no tempo t e $k_{ij}[s^{-1}]$ são as constantes que controlam a taxa de permutas entre os compartimentos. Para mais,

$$C_p^{pro}(t) = \frac{m_1(t)}{1000 \times v_1} \quad (2.2.)$$

e

$$v_1 = weight \times v_c \quad (2.3)$$

Onde v_1 é o volume do compartimento 1 dado em [L] e v_c é uma variável do paciente que representa o volume do primeiro compartimento por kilograma L/kg . Como $m_1(t)$ representa a massa no compartimento 1 e é dado em [μg] e v_1 em [L], $C_p^{pro}(t)$ é em [$\mu g/mL$].

Os parâmetros $v_c, k_{10}, k_{12}, k_{13}, k_{21}, k_{31}[min^{-1}]$ que fazem parte desta equação são aqueles que devem ser estimados. A figura 2.4 mostra o significado físico destas equações. Este modelo pode também ser descrito da seguinte forma:

$$\dot{x}(t) = Ax(t) + Bu(t) \quad (2.4)$$

$$y(t) = Cx(t) \quad x(t) = \begin{bmatrix} m_1(t) \\ m_2(t) \\ m_3(t) \end{bmatrix} \quad (2.5)$$

$$A = \begin{bmatrix} -k_{10} & -k_{12} & -k_{13} & k_{21} & k_{31} \\ & k_{12} & & -k_{22} & 0 \\ & k_{13} & & 0 & -k_{31} \end{bmatrix} \quad B = \begin{bmatrix} \frac{10^4}{3600} \\ 0 \\ 0 \end{bmatrix} \quad C = \begin{bmatrix} \frac{1}{1000 \times v_1} & 0 & 0 \end{bmatrix}$$

Para o caso do analgésico remifentanil a sua farmacocinética também é descrita por um modelo de 3 compartimentos seguindo a mesma estrutura matemática. Os parâmetros restantes $k_{e0}[min^{-1}]$, pertencem ao modelo farmacodinâmico:

$$C_e^{prop}(s) = \frac{1}{\frac{1}{k_{e0}}s + 1} C_p^{prop}(s) \quad (2.6)$$

As concentrações de efeito ($C_e^{pro}(t)$), representa a concentração de fármaco no local de efeito, no caso do hipnótico propofol seria o cérebro. Ao contrário do $C_p^{pro}(t)$ que representa a concentração plasmática de propofol. O modelo completo é mostrado na figura 2.4.

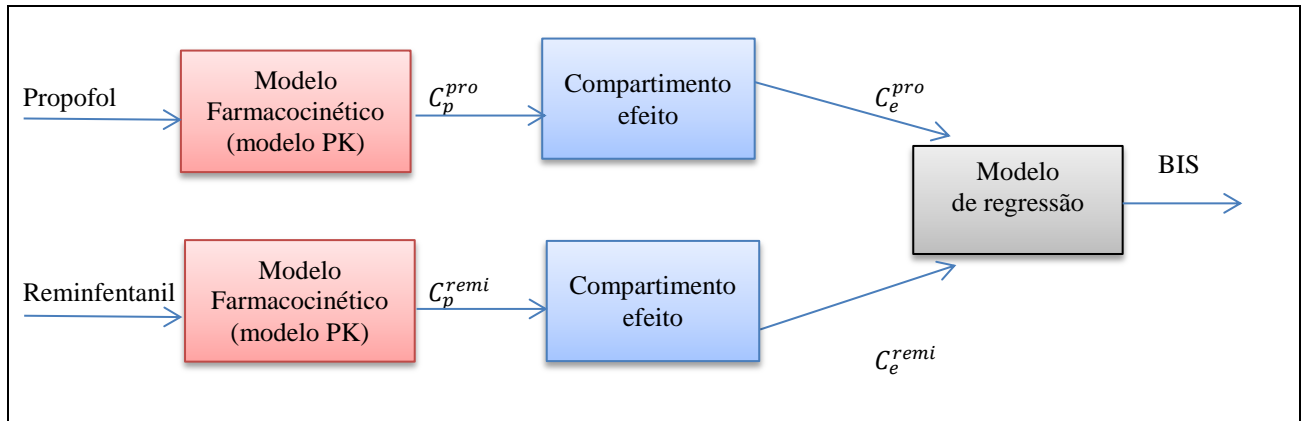


Figura 2.4: Farmacocinético e Farmacodinâmico

2.4 Target Controlled Infusion (TCI)

Target Controlled Infusion (TCI) é uma técnica de infusão computadorizada, que tem sido amplamente usada com uma variedade de fármacos, para controlar o plasma teórico ou concentrações de efeito local, definindo um alvo para um determinado efeito anestésico desejado [4, 6]. A administração segura e eficaz de fármacos anestésicos requer o conhecimento *apriori* das características farmacocinéticas e farmacodinâmicas (PK/PD) [19]. Os sistemas de TCI utilizam modelos farmacocinéticos, que descrevem matematicamente o processo de distribuição e eliminação do fármaco, para calcular a taxa de infusão de fármaco necessária para atingir a concentração de efeito desejada [20, 21, 22, 23].

A base de dados deste trabalho foi recolhida pelo sistema TCI RugLoopII, cujo software controla as doses de fármaco a administrar. No caso do propofol, usa o modelo PK/PD de Schnider [28,32,45]. e no caso do remifentanil (analgésico utilizado em alguns dos pacientes), usa o modelo PK/PD de Minto [28,32,45]. A (figura 2.5) representa o esquema dos voluntários ligados ao monitor de BIS, monitor DATEX (para pressões artérias, CO₂, O₂, etc...), TCI Fresenius

para as infusões em TCI de propofol e remifentanil. O computador recolhe os dados dos três equipamentos.

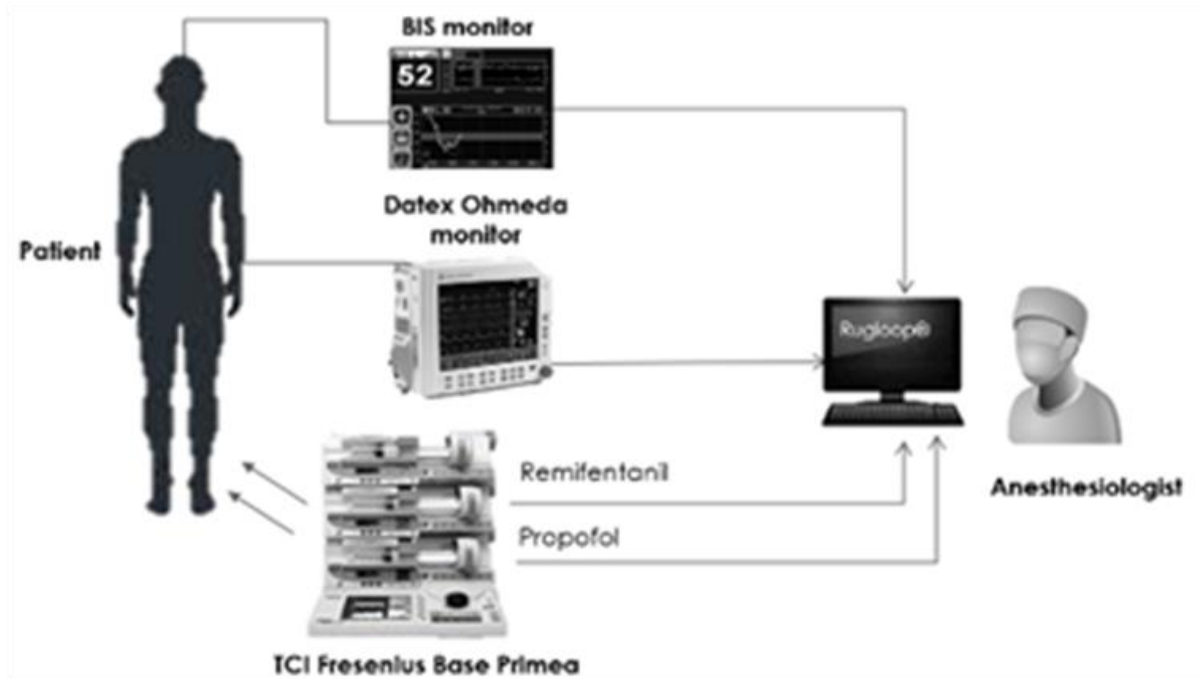


Figura 2.5:Diagrama de recolha de dados dos voluntários

Capítulo 3

Caracterização dos dados

3.1. Descrição da base de dados

Se os resultados obtidos com base em dados reais são importantes, não menos importante é a informação das características dessas bases de dados, pois a validade dos resultados está dependente das condições em que os mesmos são obtidos. A informação contida neste capítulo é relevante não só para validar os resultados mas também para se perceber alguns procedimentos nas metodologias desenvolvidas nos capítulos que se seguem.

Os dados reais disponibilizados para os trabalhos desenvolvidos nesta dissertação, diz respeito ao índice de inconsciência (DoA) induzido pela administração do hipnótico propofol e também do analgésico remifentanil.

Os dados foram recolhidos num estudo clínico realizado no início de 2008, no Hospital Geral de Santo António (HGSA). O estudo foi realizado com a autorização do Director do Bloco Operatório, do Director de Serviço de Anestesiologia e aprovação pelo Conselho de Administração, à data este era o procedimento obrigatório no HGSA. Todos os voluntários participaram no estudo após assinarem consentimento informado. Não foram pagos honorários a nenhum interveniente nem foi paga qualquer importância aos voluntários. Estiveram sempre presentes na recolha de dados três médicos anestesistas, um enfermeiro de anestesia e uma investigadora do Serviço de Anestesiologia. Todos os voluntários foram internos de anestesiologia, ou especialistas de anestesiologia, ou estudantes de medicina.

Esta base de dados, foi recolhida pelo sistema TCI RugLoopII, cujo software controla as doses de fármaco a administrar. No caso do propofol, usa o modelo PK/PD de Schnider [40] e no caso do remifentanil, usa o modelo PK/PD de Minto [41].

Esta base de dados, é composta por 8 voluntários saudáveis, submetidos a anestesia geral por administração de propofol e remifentanil. A DoA é monitorizada pelo sinal BIS® (monitor A2000 Xp da Aspect Medical). No que diz respeito às variáveis da intervenção clínica e que descrevem o processo de sedação como é o caso do volume total, da concentração de efeito, C_e , e da concentração pretendida e programada pelo anestesista, CT, para o propofol e o remifentanil os registos são obtidos de 5 em 5 segundos. Foi então iniciada a infusão de Propofol numa concentração inicial de 1,5µg/ml. Um minuto depois de atingida essa concentração de efeito, foi dado início a aumentos sucessivos da concentração de efeito em degraus de 0,5 µg/ml, tendo-se avaliado, um minuto após atingida cada concentração de efeito, o nível de sedação e registado os valores de dióxido de carbono, CO₂, expirado, frequência respiratória, BIS, saturação de oxigénio

no sangue, SpO₂, Tensão Arterial, TA, (medida a cada minuto), e Frequência Cardíaca, FC. O nível de sedação foi avaliado pela escala de *OAA/S*.

A concentração de efeito do propofol foi sendo aumentada até a obtenção de um grau de sedação de nível 3 na escala de *OAA/S*. Atingido esse nível, adicionou-se remifentanil na concentração de 20 ng/ml que foi infundido por TCI na concentração inicial de 1 ng/ml. Seguidamente, efetuaram-se aumentos da concentração do remifentanil de 0,5 ng/ml até a obtenção de um grau de sedação de nível 2 na escala de *OAA/S* e/ou CO₂ expirado superior a 50 mmHg e/ou SpO₂, inferior a 94%. Uma vez atingida uma destas três situações, procedeu-se à diminuição da concentração de efeito do propofol por degraus de 0,5 µg/ml até a concentração que permitisse manter um grau de sedação de nível 3 na escala de *OAA/S*, mantendo-se a infusão de remifentanil na mesma concentração atingida antes. Uma vez atingido o nível 3 na escala *OAA/S* interrompeu-se a infusão dos dois fármacos, mantendo-se a monitorização e o registo das variáveis até se atingir o nível 5 na escala *OAA/S*.

Se CO₂ expirado superior a 55 mmHg e/ou frequência respiratória inferior a 6 cpm, proceder-se-ia à redução imediata de remifentanil para a concentração de efeito, no caso de ocorrer apneia suspender-se-ia o remifentanil, ventilando-se manualmente com máscara facial e oxigénio através do circuito de um ventilador Siemens 900C e, se necessário, administrando-se Naloxona em doses fracionadas de 100 µg por via entra venosa; no caso de ocorrer bradicardia, administrar-se-ia Atropina na dose de 0,5 mg via entra venosa.

Quando atingida a concentração máxima de remifentanil era realizada uma colheita de sangue arterial a partir da artéria radial para determinação de gasimetria, nomeadamente dos valores de *Ph*, PaO₂ e sobretudo de PaCO₂.

Foram aplicados estímulos dolorosos, que consistiram em aplicação de estímulos eléctricos utilizando o neuroestimulador habitualmente usado para monitorização do bloqueio neuromuscular.

Foram aplicados eléctrodos sobre o trajecto do nervo cubital no pulso. Após atingida a concentração máxima de propofol foi aplicada estimulação com estímulos simples de intensidade crescente (30, 50 e 70 mA), observando-se a ocorrência de movimento de fuga do membro ou outros movimentos como reacção ao estímulo. Se não ocorresse reacção era aplicado um estímulo tetânico de 50 Hz com a duração de 5 segundos. Uma vez atingida a concentração máxima de remifentanil procedia-se de igual modo. A estimulação com estímulo tetânico e a compressão forte do rebordo da órbita foram usados também como estímulo. Depois de atingido

o nível 5 na escala OAA/S, o estudo foi considerado terminado, passando o voluntário à sala de recobro onde permaneceu monitorizado durante mais 30 minutos, recebendo soro. Terminado esse período, passou para um cadeirão onde ficou sentado durante 45 minutos, sem monitorização, bebendo um chá com açúcar e comendo bolachas. Depois pôde vestir-se, deambular um pouco e regressar a casa, tendo o cuidado de não conduzir nesse dia.

3.2. Caracterização das variáveis anestésicas

As variáveis propofol , remifentanil e BIS apresentam uma distribuição normal (teste Kolmogorov-Smirnov) (figura e tabela 3.2.1.).

Tabela 3.2.1: Caracterização das variáveis anestésicas				
		Propofol Ce (ug/ml)	Remifentanil Ce (ng/ml)	BIS
N	Válido	4214	4214	4214
	Ausente	0	0	0
Média		2,4130	,5814	75,0138
Mediana		2,4994	,3514	77,0000
Desvio Padrão		1,14581	,62501	15,92569
Variância		1,313	,391	253,628
Mínimo		,00	,00	27,00
Máximo		4,52	2,02	98,00
Percentis	25	1,5251	,0000	65,0000
	50	2,4994	,3514	77,0000
	75	3,4601	1,0108	85,0000

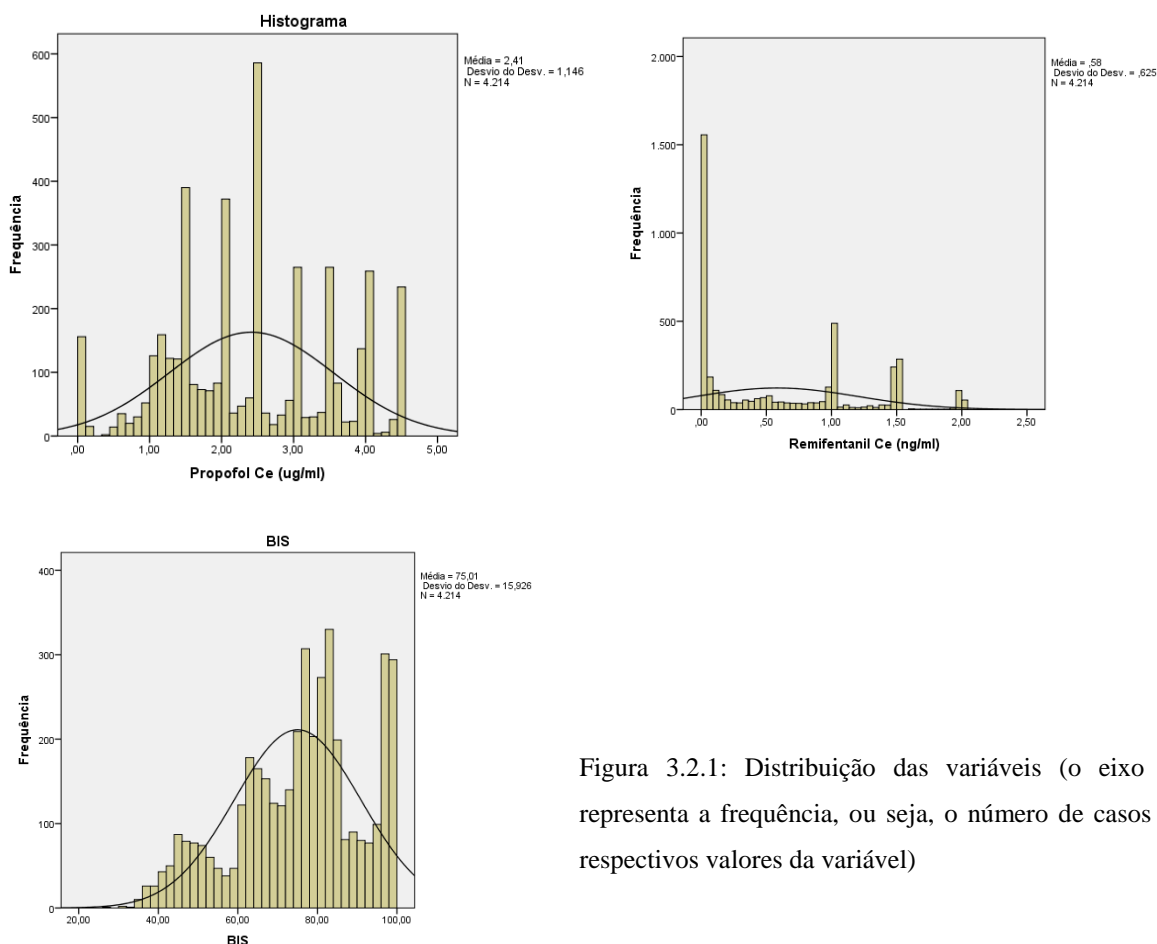


Figura 3.2.1: Distribuição das variáveis (o eixo vertical representa a frequência, ou seja, o número de casos para os respectivos valores da variável)

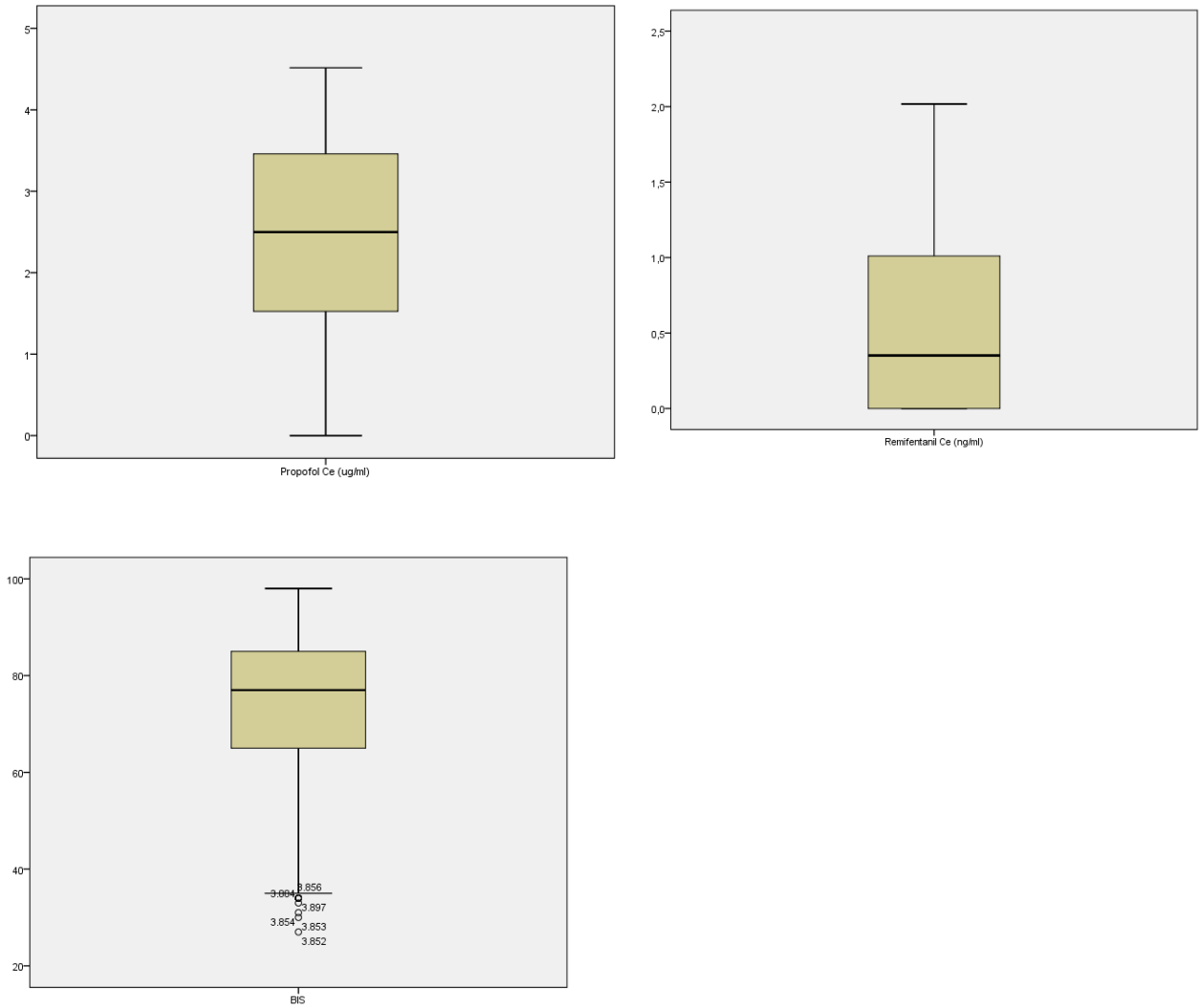


Figura 3.2.2 : Diagrama de extremos das variáveis propofol, remifentani e BIS.

O grupo de voluntários foi constituído com idades compreendidas entre os 27 ± 10 anos, 68.3 ± 14.1 Kg de peso, 174.3 ± 8.2 cm de altura e 6 eram do sexo masculino e dois do sexo feminino. A tensão arterial sistólica inicial/ de referência foi 135.8 ± 19 mmHg e a frequência cardíaca inicial foi 67.6 ± 14.4 bpm.

Com a administração de propofol foi possível, em todos os voluntários, atingir a pontuação 3 na escala OAA/S. A concentração máxima de propofol alcançada foi de 3.5 ± 0.7 µg/ml. Para esta concentração do fármaco, o valor do BIS foi de 68 ± 15.7 ; a tensão arterial sistólica foi de 125.6 ± 14.8 mmHg; a frequência cardíaca inicial de 64.3 ± 10 bpm e a SpO₂, foi de 99 ± 1.1 %. A estimulação com estímulos simples ou com tétano aplicada após atingida a concentração máxima de propofol, produziu reacção de fuga ou verbalização da sensação de dor em todos os voluntários. Com a administração de remifentanil foi possível, em todos os voluntários, atingir a pontuação 2 na escala OAA/S. A concentração máxima de remifentanil atingida foi de 1.4 ± 0.3 ng/ml. A concentração máxima foi de 1.03 ng/ml em dois voluntários, 1.53 ng/ml em cinco e 2.03 ng/ml em um voluntário. À concentração máxima de remifentanil correspondeu um valor de BIS de 60.8 ± 10.8 ; tensão arterial sistólica de 122.9 ± 13.8 mmHg; FC de 58.5 ± 11.2 bpm e SpO₂ de 98.7 ± 1.1 %.

Em relação aos valores registados no momento da concentração máxima de propofol, tensão arterial sistólica não apresentou variação significativa, a frequência cardíaca mostrou uma tendência para baixar.

Tabela 3.2.2: Estatísticas descritivas das características dos pacientes da base de dados da DoA induzida pela administração de propofol e remifentanil.

Dados Demográficos	
Idade (anos)	27 ± 10
Peso (Kg)	$68,3 \pm 14.1$
Altura (cm)	174.3 ± 8.2
Sexo	6 masculino /2 feminino

Valores expressos em Média \pm SD

Concentrações Máximas e Mínimas de propofol Ce ($\mu\text{g/mL}$) e remifentanil Ce ($\eta\text{g/mL}$)

Concentração máxima propofol Ce ($\mu\text{g/mL}$) 3.5 ± 0.7

Concentração máxima remifentanil Ce ($\eta\text{g/mL}$) 1.4 ± 0.3

Valores expressos em Média \pm SD

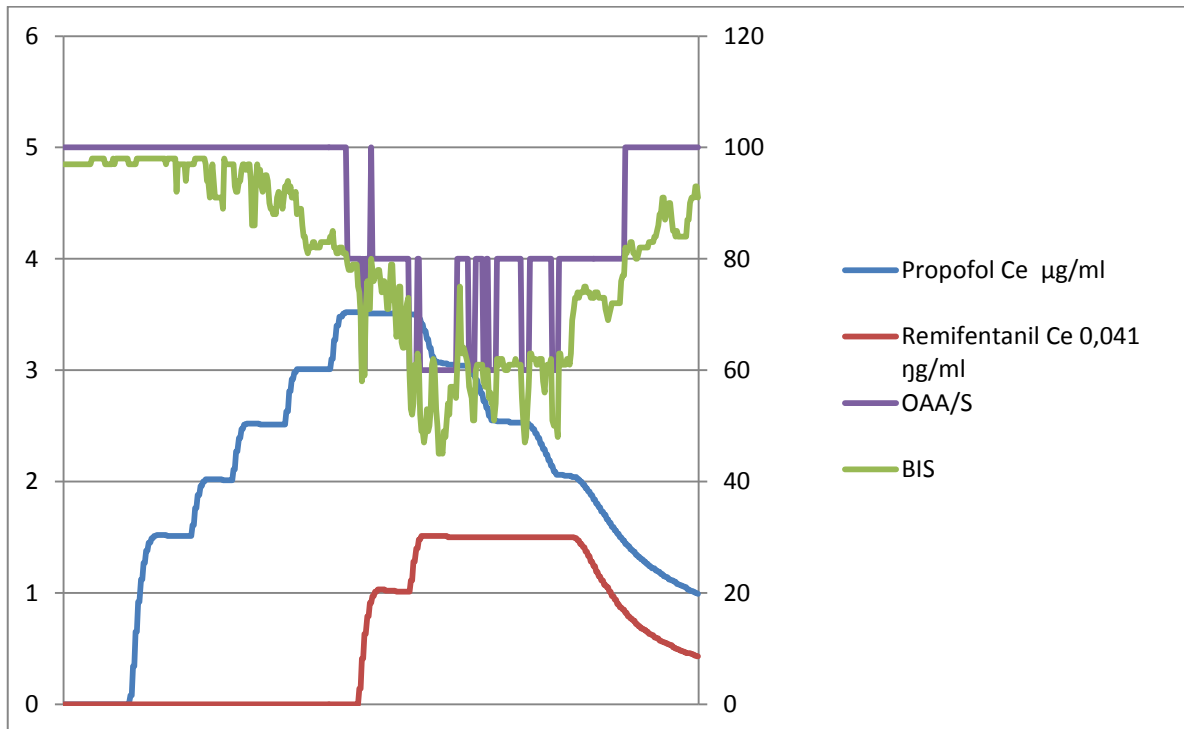


Figura 3.2.3: Voluntário 1 propofol Ce ($\mu\text{g. mL}$), remifentanil Ce ($\eta\text{g. mL}$), BIS e OAA/S

Capítulo 4

Análise de Regressão

4.1. Introdução

“O termo ‘regressão’ foi proposto pela primeira vez por Sir Francis Galton em 1885 num estudo onde demonstrou que a altura dos filhos não tende a reflectir a altura dos pais, mas tende sim a regredir para a média da população. Atualmente, o termo “Análise de Regressão” define um conjunto vasto de técnicas estatísticas usadas para modelar relações entre variáveis e prever o valor de uma ou mais variáveis dependentes (ou de resposta) a partir de um conjunto de variáveis independentes (ou preditoras).” [26]. O objetivo desta técnica é identificar e estimar uma função que descreva, o mais próximo possível, a relação entre essas variáveis e que assim irá permitir prever o valor que a variável dependente Y irá assumir para um determinado valor da variável independente X .

O modelo de regressão poderá ser escrito genericamente como:

$$Y = f(X_1, X_2, X_3, \dots, X_n) + \varepsilon$$

onde o termo ε representa uma perturbação aleatória na função ou o erro da aproximação. O número de variáveis independentes varia entre aplicações: quando se tem apenas uma variável independente, denomina-se Modelo de Regressão Simples; quando se tem mais de uma variável independente, denomina-se de Modelo de Regressão Múltipla. A forma da função também varia, podendo ser representada por uma equação linear, polinomial ou outro mesmo tipo de função (simples ou multivariada).

4.2. Regressão e Correlação Linear

4.2.1 Testes de Hipóteses sobre o Coeficiente de Correlação

A análise de correlação tem como objetivo a avaliação do grau de associação entre duas variáveis, X e Y , ou seja, mede a “força” de relacionamento linear entre as variáveis X e Y e designa-se por ρ .

O coeficiente de correlação linear, também chamado de covariância normalizada, é representado por [26,27,46,48]:

$$\rho_{X,Y} = \frac{\sigma_{X,Y}}{\sigma_X \sigma_Y}$$

Onde: $\sigma_{X,Y}$ é a covariância entre as variáveis X e Y

$\sigma_X \sigma_Y$ são os desvios padrão das variáveis X e Y

A covariância entre duas variáveis pode ser estimada pela equação:

$$S_{X,Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

Onde $S_{X,Y}$ é a covariância amostral entre as duas variáveis X e Y

\bar{x} e \bar{y} são as médias aritméticas de cada uma das variáveis

n o tamanho da amostra

x_i e y_i são as observações simultâneas das variáveis

Admitindo-se que a distribuição conjunta das variáveis é normal bivariada, para quantificar a relação entre as duas variáveis quantitativas, utiliza-se como medida da correlação o coeficiente de correlação de Pearson cujo estimador é dado por:

$$\hat{\rho} = \frac{S_{X,Y}}{s_X s_Y}$$

Onde $s_X = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$ e $s_Y = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$ são os desvios padrão das amostras.

Para se decidir sobre a existência de correlação e o sentido da variação da reta de regressão, calcula-se ρ e o erro de ρ , e seguidamente efetua-se um teste de t -Student, para as seguintes hipóteses:

$H_0: \rho = 0$, a reta de regressão em y é paralela ao eixo das abcissas.

$H_1: \rho \neq 0$, a reta de regressão em y é paralela ao eixo das abcissas.

A estatística do teste é $t_0 = \frac{\hat{\rho}\sqrt{n-1}}{\sqrt{1-\hat{\rho}^2}}$

Onde: t_0 é a estatística do teste

n o tamanho da amostra

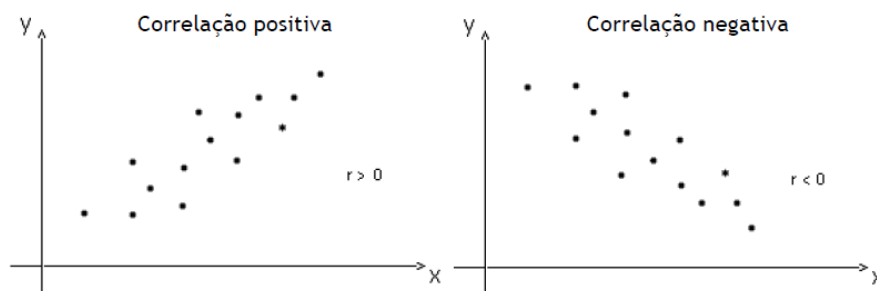
$\hat{\rho}$ é a estimativa do coeficiente de correlação linear

Para encontrar o $t_{critico}$ (t_c) consulta-se uma tabela de t -Student, e é interpretado conforme o seguinte critério:

$t < t_c$	t_c	$t > t_c$
t não é significativo		t é significativo
ρ não é significativamente diferente de 0 (a reta é paralela ao eixo dos xx)		ρ é significativamente diferente de 0 (a reta não é paralela ao eixo dos xx)

A partir da observação do diagrama de dispersão verificamos se a correlação entre as duas variáveis é mais ou menos forte, de acordo com a proximidade dos pontos em relação a uma reta. Na Figura 4.1., podemos observar alguns exemplos de gráficos de dispersão e a respectiva “classificação” da correlação. Essa correlação pode ser positiva (para valores crescentes de X há uma tendência a valores também crescentes de Y) ou negativa (para valores crescentes de X a tendência é observarem-se valores decrescentes de Y). A figura seguinte ilustra correlações lineares positivas e negativas.

O coeficiente de correlação linear de Pearson é adimensional e varia entre -1 e +1, o que não ocorre com a covariância. Assim, as unidades adotadas pelas variáveis não afetam o valor do coeficiente de correlação. Caso os dados se alinhem perfeitamente ao longo da reta com declive positivo teremos a correlação linear positiva perfeita com o coeficiente de Pearson igual a 1. A correlação linear negativa perfeita ocorre quando os dados se alinham perfeitamente ao longo de uma reta com declive negativo e o coeficiente de correlação de Pearson é igual a -1.



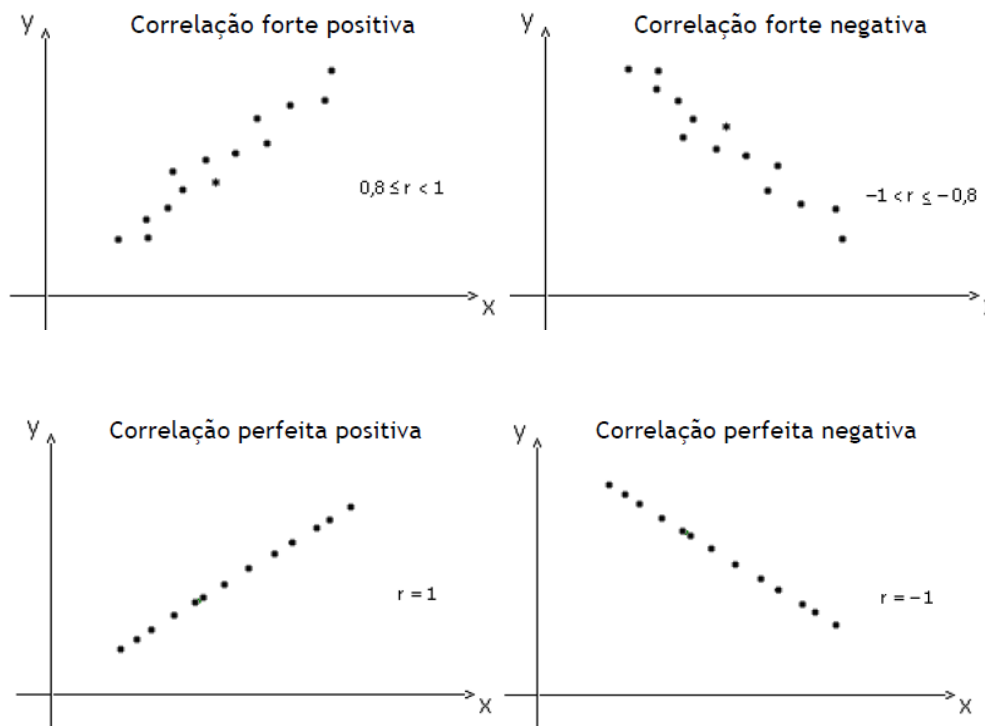


Figura 4.1: Classificação da correlação através do diagrama de dispersão, disponível em [58]

Quando a escala de medida é ordinal devemos utilizar o coeficiente de correlação de Spearman, pois este, ao contrário do coeficiente de correlação de Pearson, não requer a suposição que a relação entre as variáveis é linear, nem requer que as variáveis sejam medidas em intervalo de classe, podendo ser usado para as variáveis medidas no nível ordinal.

É importante realçar que as correlações ordinais não podem ser interpretadas da mesma maneira que as correlações de Pearson. Inicialmente não mostram tendência linear, mas podem ser consideradas como índices de monotonia, ou seja, permitem-nos avaliar as variações para aumentos positivos da correlação (aumentos no valor de X correspondem a aumentos no valor de Y) e para coeficientes negativos [26,27,46,48].

4.3. Modelo de Regressão Linear Simples

A análise de regressão linear estuda a relação entre a variável dependente ou variável resposta (Y) e uma ou várias variáveis independentes ou regressoras (X_1, \dots, X_p).

Esta relação representa-se por meio de um modelo matemático ou seja, por uma equação que associa a variável dependente (Y) com as variáveis independentes (X_1, \dots, X_p).

O Modelo de Regressão Linear Simples define-se como a relação linear entre a variável dependente (Y) e uma variável independente (X).

Enquanto o Modelo de Regressão Linear Múltiplo se define como a relação linear entre a variável dependente (Y) e várias variáveis independentes (X_1, \dots, X_p).

Neste capítulo vamos apenas debruçar-nos sobre o modelo de regressão linear simples e múltipla. Será apresentado o modelo teórico e os seus pressupostos, assim como a estimação dos parâmetros do modelo pelo método dos mínimos quadrados. Serão ainda construídos testes e intervalos de confiança para os parâmetros do modelo [26,27,46,48].

4.3.1. Modelo Teórico

A equação representativa do modelo de regressão linear simples é dado por:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n \quad (4.1)$$

onde:

- y_i representa o valor da variável resposta ou dependente, Y , na observação i , $i = 1, \dots, n$ (aleatória);
- x_i representa o valor da variável independente, X , na observação i , $i=1, \dots, n$ (não aleatória);
- $\varepsilon_i, i = 1, \dots, n$ são variáveis aleatórias que correspondem ao erro (variável que permite explicar a variabilidade existente em Y e que não é explicada por X);
- β_0 e β_1 correspondem aos parâmetros do modelo.

O parâmetro β_0 representa o ponto em que a recta regressora corta o eixo dos yy quando $X = 0$ e é chamado de intercepto ou coeficiente linear.

O parâmetro β_1 representa a inclinação da reta regressora, expressando a taxa de mudança em Y , ou seja, indica a mudança na média da distribuição de probabilidade de Y para um aumento de uma unidade na variável X .

Reta de Regressão

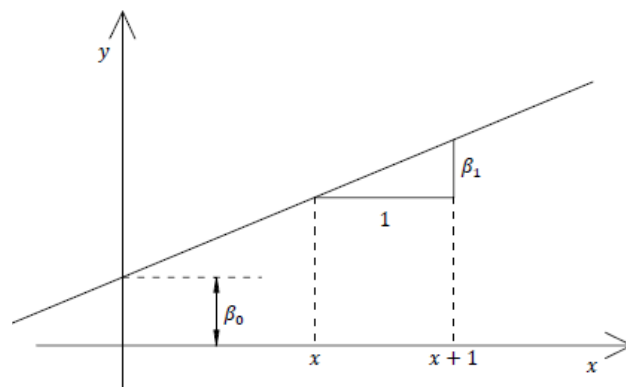


Figura 4.2: Interpretação geométrica dos parâmetros do modelo de regressão linear simples

4.3.2. Pressupostos do Modelo

Ao definir o modelo (4.1) estamos a pressupor que:

- a) A relação existente entre Y e X é linear.
- b) Os erros são independentes com média nula.

Pressupondo então que $E(\varepsilon_i) = 0$, tem-se:

$$\begin{aligned} E(y_i) &= E(\beta_0 + \beta_1 x_i + \varepsilon_i) \\ &= \beta_0 + \beta_1 x_i + E(\varepsilon_i) \\ &= \beta_0 + \beta_1 x_i. \end{aligned} \tag{4.2}$$

Podemos ainda afirmar que o erro de uma observação é independente do erro de outra observação, o que significa que:

$$\text{cov}(\varepsilon_i, \varepsilon_j) = E(\varepsilon_i \varepsilon_j) - E(\varepsilon_i)E(\varepsilon_j) = E(\varepsilon_i \varepsilon_j) = 0, \quad \text{para } i \neq j, \quad i, j = 1, \dots, n.$$

c) A variância do erro é constante, isto é $\text{var}(\varepsilon_i) = \sigma^2, i = 1, \dots, n.$

Tem-se então

$$\text{var}(\varepsilon_i) = E(\varepsilon_i^2) - \underbrace{[E(\varepsilon_i)]^2}_{=0} = E(\varepsilon_i^2) = \sigma^2,$$

$$\text{var}(y_i) = \text{var}(\beta_0 + \beta_1 x_i + \varepsilon_i) = \underbrace{\text{var}(\beta_0 + \beta_1 x_i)}_{\text{termo constante}} + \underbrace{\text{var}(\varepsilon_i)}_{=\sigma^2} = \sigma^2.$$

d) Os erros, $\varepsilon_i, i = 1, \dots, n,$ são normalmente distribuídos.

Concluimos pois, de b) e c), que

$$\varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n,$$

Logo

$$y_i \sim N(\beta_0 + \beta_1 x_i + \sigma^2), \quad i = 1, \dots, n.$$

4.3.3. Estimação dos Parâmetros do Modelo

Supondo que existe efetivamente uma relação linear entre X e Y , coloca-se a questão de como estimar os parâmetros β_1 e β_0 .

Karl Gauss entre 1777 e 1855 propôs estimar os parâmetros β_1 e β_0 visando minimizar a soma dos quadrados dos desvios, $e_i, i = 1, \dots, n,$ chamando este processo de método dos mínimos quadrados. Este método será descrito de seguida [26,27].

4.3.3.1. Método dos Mínimos Quadrados

O método dos mínimos quadrados consiste na obtenção dos estimadores dos coeficientes de regressão $\hat{\beta}_0$ e $\hat{\beta}_1$, minimizando os resíduos do modelo de regressão linear, calculados como a diferença entre os valores observados, y_i , e os valores estimados, \hat{y}_i , [26,27,46,48], isto é:

$$e_i = y_i - \hat{y}_i, \quad i = 1, \dots, n.$$

Em termos gráficos, os resíduos são representados pelas distâncias verticais entre os valores observados e os valores ajustados, como mostra a Figura 4.3.

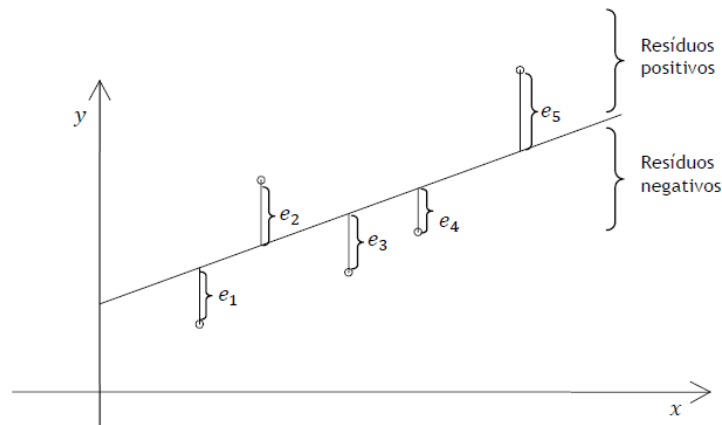


Figura 4.3. - Representação gráfica dos resíduos

O método dos mínimos quadrados propõe então encontrar os valores de β_0 e β_1 , para os quais a soma dos quadrados dos resíduos (SQE) é mínima. Tem-se então:

$$\begin{aligned}
 SQE &= \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\
 &= \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \\
 &= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2,
 \end{aligned} \tag{4.3}$$

com $\sum_{i=1}^n e_i = 0$ (daí o facto de ser considerado o quadrado de e_i , $i=1, \dots, n$).

Precisamos agora de calcular as derivadas parciais de SQE em ordem a β_0 e β_1 , obtendo-se:

$$\begin{cases} \frac{\partial SQE}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i), \\ \frac{\partial SQE}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i \end{cases}$$

Igualando estas derivadas a zero e substituindo β_0 e β_1 por $\hat{\beta}_0$ e $\hat{\beta}_1$, por forma a indicar valores concretos destes parâmetros, tem-se:

$$\begin{cases} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} n \hat{\beta}_0 \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_1 x_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 - \sum_{i=1}^n x_i y_i = 0 \end{cases} \quad (4.4)$$

$$\Leftrightarrow \begin{cases} \hat{\beta}_0 = \frac{\sum_{i=1}^n y_i}{n} - \frac{\hat{\beta}_1 \sum_{i=1}^n x_i}{n} \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases} \quad (4.5)$$

Em que $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ e $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$, representarem médias de X e Y , respetivamente.

A partir da 2ª equação de (4.4) para obter a expressão de $\hat{\beta}_1$. Ora

$$\hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i$$

$$\Leftrightarrow \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{i=1}^n x_i.$$

Como

$$\begin{aligned} & \sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i + \hat{\beta}_1 \bar{x} \sum_{i=1}^n x_i \\ &= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} + n \hat{\beta}_1 \bar{x}^2, \end{aligned}$$

Vem

$$\begin{aligned} \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} + n\hat{\beta}_1 \bar{x}^2 \\ \Leftrightarrow \hat{\beta}_1 \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \Leftrightarrow \\ \Leftrightarrow \hat{\beta}_1 &= \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \end{aligned} \quad (4.6)$$

$\hat{\beta}_0$ e $\hat{\beta}_1$, anteriormente determinados em (4.5) e (4.6), são designados como os Estimadores de Mínimos Quadrados de β_0 e β_1 .

São vantagens do método dos mínimos quadrados:

- Obter as melhores estimativas, pois elas não são enviesadas;
- Ter em conta os desvios maiores, diluindo o efeito dos maiores valores;
- Permitir realizar testes de significância na equação de regressão;
- A reta de regressão passa pelo ponto obtido pelo cálculo das médias das duas amostras.

De seguida serão apresentadas algumas propriedades do ajuste dos mínimos quadrados.

- Como vimos, os resíduos correspondem à diferença entre os valores observados, y_i , $i = 1, \dots, n$ e os correspondentes valores ajustados, \hat{y}_i , $i = 1, \dots, n$, isto é:

$$\begin{aligned} e_i &= y_i - \hat{y}_i \\ &= y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i), \end{aligned}$$

Como

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, i = 1, \dots, n.$$

- $\sum_{i=1}^n e_i^2 = 0$, o que significa que a soma dos resíduos é sempre nula;
- $\sum_{i=1}^n e_i^2$ é mínima;
- $\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$ o que significa que a soma dos valores observados y_i é igual à soma dos valores ajustados \hat{y}_i ;
- A reta obtida pelo método dos mínimos quadrados passa sempre pelo ponto (\bar{x}, \bar{y}) .

4.3.4. Teste de Hipóteses e Intervalos de Confiança para os Parâmetros do Modelo

Nesta secção construiremos testes de hipóteses e intervalos de confiança para β_0 e β_1 , considerando os pressupostos anteriormente referidos. Estes pressupostos levaram-nos a concluir que as observações

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2), \quad i = 1, \dots, n.$$

4.3.4.1. Teste e intervalos de confiança para β_1

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2}. \end{aligned}$$

Vimos também que a distribuição amostral de β_1 , para o modelo de regressão normal também é normal, uma vez que β_1 é uma combinação linear dos y_i , com:

$$E(\hat{\beta}_1) = \beta_1;$$

$$var(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

tem-se

$$(\hat{\beta}_1) \sim N\left(\beta_1; \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right).$$

Consideremos que pretendemos testar as hipóteses

$$\begin{cases} H_0: \beta_1 = \beta_1' \\ H_1: \beta_1 \neq \beta_1' \end{cases}$$

Logo, a estatística de teste será dada por

$$T = \frac{\hat{\beta}_1 - \beta_1'}{S_{\hat{\beta}_1}} \sim t_{(n-2)},$$

Com

$$S_{\hat{\beta}_1} = \sqrt{\frac{QME}{\sum_{i=1}^n (x_i - \bar{x})^2}} = \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

T tem distribuição t – *student* com $n - 2$ graus de liberdade, $t_{(n-2)}$, (ver[26,27]).

Consideremos agora que pretendemos testar

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases} \quad (4.7)$$

que são as hipóteses que queremos testar no modelo em questão. Neste caso, a estatística de teste poderá ser reescrita da seguinte forma

$$T = \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}} \sim t_{(n-2)}. \quad (4.8)$$

Logo, rejeita-se H_0 , para um nível de significância α , se $|T_{obs}| > t_{(1-\alpha/2, n-2)}$, onde T_{obs} representa o valor observado da estatística T e $t_{(1-\alpha/2, n-2)}$ o quantil de ordem $1 - \alpha/2$, da distribuição t com $n - 2$ graus de liberdade.

No que diz respeito ao intervalo de confiança, a $(1 - \alpha) \times 100\%$, para β_1 esse será dado por:

$$\left[\hat{\beta}_1 - t_{(1-\alpha/2, n-2)} S_{\hat{\beta}_1}; \hat{\beta}_1 + t_{(1-\alpha/2, n-2)} S_{\hat{\beta}_1} \right].$$

4.3.4.2. Teste e Intervalos de Confiança para β_0

Como vimos, o estimador pontual de β_0 é dado por:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Assumindo a normalidade das observações e visto que:

$$E(\hat{\beta}_0) = \beta_0$$

$$\text{var}(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

tem-se

$$(\hat{\beta}_0) \sim N \left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \right).$$

Consideremos as hipóteses

$$\begin{cases} H_0: \beta_0 = \beta_0' \\ H_1: \beta_0 \neq \beta_0' \end{cases}$$

Logo, a estatística de teste será dada por

$$T^0 = \frac{\beta_0 - \beta_0'}{S_{\hat{\beta}_0}} \sim t_{(n-2)},$$

com

$$\begin{aligned} S_{\hat{\beta}_0} &= \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \\ &= \sqrt{QME \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}. \end{aligned}$$

T^0 também segue uma distribuição t com $n - 2$ graus de liberdade, $t_{(n-2)}$.

Consideremos agora que pretendemos testar

$$\begin{cases} H_0: \beta_0 = 0 \\ H_1: \beta_0 \neq 0' \end{cases}$$

que são as hipóteses que queremos testar no modelo em questão. Neste caso a estatística de teste poderá ser reescrita da seguinte forma:

$$T^0 = \frac{\hat{\beta}_0}{S_{\hat{\beta}_0}} \sim t_{(n-2)}.$$

Assim, rejeita-se H_0 , para um nível de significância α , se $|T^0_{obs}| > t_{(1-\alpha/2, n-2)}$, onde T^0_{obs} representa o valor observado da estatística T^0 .

Quanto ao intervalo de confiança β_0 , com $(1 - \alpha) \times 100\%$, será dado por:

$$\left[\hat{\beta}_0 - t_{(1-\alpha/2, n-2)} S_{\hat{\beta}_0}; \hat{\beta}_0 + t_{(1-\alpha/2, n-2)} S_{\hat{\beta}_0} \right].$$

4.4. Modelo de Regressão Linear Múltipla

Muitas aplicações da análise de regressão envolvem situações com mais do que uma variável explicativa. Esse modelo de regressão recebe o nome de modelo de regressão múltipla (MRLM).

Na regressão linear múltipla assumimos que existe uma relação linear entre uma variável y (variável dependente) e p variáveis independentes (preditoras), (x_1, x_2, \dots, x_p) [26,27].

4.4.1. Modelo Teórico e seus pressupostos

O modelo de regressão linear múltipla com p variáveis explicativas é definido da seguinte forma:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad i=1, \dots, n, \quad (4.9)$$

em que:

- y_i representa o valor de variável resposta na observação $i = 1, \dots, n$;
- $x_{i1}, x_{i2}, \dots, x_{ip}$, $i = 1, \dots, n$ são os valores da i -ésima observação das p variáveis explicativas, (constantes conhecidas);
- $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ são os parâmetros ou coeficientes de regressão;
- ε_i , $i=1, \dots, n$ correspondem aos erros aleatórios.

Este modelo descreve um hiperplano p -dimensional referente às variáveis explicativas como mostra a Figura 4.4.

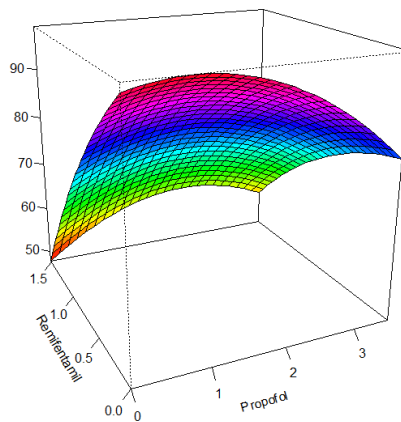


Figura 4.4: Hiperplano p -dimensional referente às variáveis explicativas. Fonte: elaborada pelo autor

Os parâmetros β_j , $j = 1, \dots, p$, representam a média esperada na variável resposta, y , quando a variável X_j , $j = 1, \dots, p$ sofre um acréscimo unitário, enquanto todas as outras variáveis X_k , $k \neq j$ são mantidas constantes.

Por esse motivo, os β_j , $j = 1, \dots, p$ são chamados de coeficientes parciais.

O parâmetro β_0 corresponde ao intercepto do plano de regressão. Se a abrangência do modelo incluir X_j , $j = 1, \dots, p$, então β_0 será a média de y nesse ponto. Caso contrário não existe interpretação prática para β_0 .

4.4.1.1. Interações

Vamos considerar o caso particular do modelo de regressão linear múltipla com duas variáveis explicativas X_1 e X_2 . Assim, o modelo será definido por:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon. \quad (4.10)$$

Sempre que considerarmos um modelo mais complexo, em que existe interação entre as variáveis explicativas, obtemos:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 \underbrace{X_1 X_2}_{\text{Interação}} + \varepsilon. \quad (4.11)$$

Neste caso, $X_1 X_2$ representa a interação existente entre as variáveis X_1 e X_2 . Se a interação existir e for significativa, o efeito de X_1 na resposta média depende do nível X_2 e vice-versa.

4.4.1.2. Pressupostos do modelo

Os pressupostos para o modelo de regressão linear múltipla são análogos ao do modelo de regressão linear simples. Assim tem-se:

- a) $E[\varepsilon_i] = 0$, $i = 1, \dots, n$;
- b) Os erros são independentes;
- c) $V[\varepsilon_i] = \sigma^2$, $i = 1, \dots, n$ (*variâncias constantes*),
- d) Os erros têm distribuição normal.

Logo destes pressupostos, concluímos que $\varepsilon_i \sim N(0, \sigma^2)$, $i = 1, \dots, n$ e, conseqüentemente, que y tem distribuição normal com varância σ^2 e, para o caso de modelo definido em (4.9),

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 \dots + \beta_p X_p.$$

4.4.1.3. Representação matricial do método de regressão linear múltipla

Anteriormente em (4.9), a expressão geral de i -ésima observação no modelo de regressão linear (sem interação) é dada por:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad i=1, \dots, n,$$

Este modelo pode ser reescrito em notação matricial da seguinte forma:

$$Y = \beta X + \varepsilon \tag{4.12}$$

onde

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1,p} \\ 1 & x_{21} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{n,p} \end{bmatrix} = [1x_1 \dots x_p]$$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}.$$

Concluimos, então, que:

- ε é um vetor de dimensão $n \times 1$ cujas componentes são os erros aleatórios, $\varepsilon_i, \quad i=1, \dots, n$;
- Y é um vetor $n \times 1$ cujas componentes correspondem às n respostas, y_1, \dots, y_n , constituído pelas observações da variável resposta;
- X é uma matriz $n \times (p + 1)$ denominada matriz do modelo, cujas colunas são constituídas pelos vetores $1 = (1, \dots, 1)'$, $x_j = (x_{1j}, \dots, x_{nj})'$, $j = 1, \dots, p$. A notação A' representa a transposta da matriz A .

- β é um vetor coluna $(p + 1) \times 1$ cujos elementos são os coeficientes de regressão, $\beta_0, \beta_1, \dots, \beta_p$.

Uma vez que ε é normalmente distribuído, tendo-se $\varepsilon_i \sim N(0, \sigma^2 I_n)$, com 0 o vetor nulo e I_n a matriz identidade de ordem n , Y será normalmente distribuído com $E(Y) = \beta X$ e a matriz de variâncias-covariâncias $cov(Y) = \sigma^2 I_n$, isto é:

$$Y \sim N(\beta X, \sigma^2 I_n).$$

4.4.2. Estimação do Parâmetro do modelo

De modo análogo à regressão simples, usando o método dos mínimos quadrados, pretendemos encontrar o vetor de estimadores $\hat{\beta}$, com componentes $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$, que minimiza

$$\begin{aligned} SQE &= \sum_{i=1}^n e_i^2 = e'e = (Y - X\beta)'(Y - X\beta) \\ &= YY' - Y'X\beta - \beta'X'Y + \beta'X'X\beta \\ &= Y'Y - 2\beta'X'Y + \beta'X'X\beta, \end{aligned}$$

uma vez que se tem $Y'X\beta = \beta'X'Y$, pois este produto é igual a um escalar. Derivando $\hat{\beta}$ obtemos:

$$\frac{\partial SQE}{\partial \beta} = 2X'Y + 2X'X\beta.$$

Igualando a derivada a zero e substituindo β por $\hat{\beta}$, obtemos:

$$\begin{aligned} -2X'Y + 2X'X\hat{\beta} &= 0 \\ \Leftrightarrow (X'X)\hat{\beta} &= X'Y \end{aligned}$$

$$\Leftrightarrow \hat{\beta} = (X'X)^{-1}X'Y, \tag{4.13}$$

onde X^{-1} representa a matriz inversa de X . De (4.12), concluímos que o modelo de regressão linear ajustado é

$$\hat{Y} = X\hat{\beta}$$

E o vetor dos resíduos

$$e = Y - \hat{Y}.$$

4.4.3. Análise da Variância (ANOVA) Aplicada à Regressão Linear Múltipla

A análise de variância é importante para a análise de regressão linear múltipla. Este tema não foi abordado na análise de regressão linear simples, uma vez que não traz novidades em termos de aplicação dos testes, já que o teste t e o teste F darão os mesmos resultados. Chega observar que o teste F é o quadrado do teste t .

Na análise de regressão múltipla, o teste F produz um teste mais geral. Através da sua utilização determina-se se qualquer das variáveis independentes no modelo tem poder de explicação. Cada variável pode então ser testada individualmente com o teste t para determinar se é uma das variáveis significativas.

A análise de variância, baseia-se na decomposição da soma dos quadrados total, SQT , (que corresponde à variação da variável resposta), na soma dos quadrados explicada, SQR , (que corresponde à variação da variável resposta que é explicada pelo modelo) e na soma dos quadrados dos resíduos, SQE , (que corresponde à variação da variável resposta que não é explicada pelo modelo).

Desta forma, podemos escrever,

$$SQT = SQR + SQE \Leftrightarrow$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Assim, no conceito de regressão linear múltipla, as hipóteses a testar serão

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \text{ (a relação entre as variáveis é não linear)}$$

$$H_1: \exists_i \beta_j \neq 0, \quad j = 1, \dots, p.$$

Para testar a hipótese H_0 utiliza-se a estatística do teste

$$F = \frac{\frac{SQR}{p}}{\frac{SQE}{n-p-1}} \frac{QMR}{QME} \sim F_{p, n-p-1},$$

Como $\frac{SQR}{\sigma^2} \sim \chi_p^2$, $\frac{SQE}{\sigma^2} \sim \chi_{n-p-1}^2$ e SQR e SQE independentes. Assim, sob H_0 , a estatística de F segue uma distribuição F central com p e $n - (p + 1)$ graus de liberdade, $F_{p,n-p-1}$.

Logo, se $F_{obs} > F_{(1-\alpha;p,n-p-1)}$ rejeita-se a hipótese H_0 , com F_{obs} o valor observado de estatística F e $F_{(1-\alpha;p,n-p-1)}$ o quantil $1 - \alpha$ de distribuição F central com p e $n - p - 1$ graus de liberdade. Ao rejeitarmos H_0 concluímos que pelo menos uma das variáveis explicativas contribui significativamente para o modelo.

Estas somas de quadrados podem ser apresentadas numa tabela como a que apresentamos de seguida.

Causas de Variação	Graus Liberdade	Soma de quadrados	Quadrados Médios	F
Regressão (modelo)	p	SQR	$QMR = \frac{SQR}{p}$	$F = \frac{QMR}{QME}$
Erro (resíduo)	$n - p - 1$	SQE	$QEM = \frac{SQE}{n-p-1}$	
Total	$n - 1$	SQT		

Tabela 4.1: Tabela da análise de variância (ANOVA)

Como vimos anteriormente, o coeficiente de determinação é igual ao quadrado do coeficiente de correlação de Pearson, que agora poderá ser reescrito da seguinte forma

$$R^2 = \frac{\text{variação explicada}}{\text{variação total}}$$

$$= \frac{SQR}{SQT} = 1 - \frac{SQE}{SQT}$$

Este coeficiente é usado para quantificar a capacidade explicativa do modelo, ou seja, é uma medida da proporção da variação da variável resposta Y que é explicada pela equação de regressão quando estão envolvidas as variáveis independentes X_1, X_2, \dots, X_p .

Como já foi referido anteriormente,

$$0 \leq R^2 \leq 1$$

Temos, no entanto, de ter atenção ao facto de que $R^2 \cong 1$ não significa que o modelo de regressão providencia um bom ajustamento aos dados, dado que a adição de uma variável aumenta sempre o valor deste coeficiente (mesmo que tenha muito pouco poder explicativo sobre a variável resposta).

Assim, quando R^2 é elevado em determinados modelos, leva-nos a interpretações erradas de novas observações ou estimativas pouco fiáveis do valor esperado de Y . Por isso, concluímos que R^2 poderá não ser um bom indicador do grau de ajustamento do modelo.

Assim sendo, é preferível utilizar o coeficiente de determinação ajustado, que é uma medida ajustada do coeficiente de determinação e que é “penalizada” quando são adicionadas variáveis pouco explicativas.

O coeficiente de determinação ajustado é definido por:

$$R_a^2 = 1 - \left(\frac{n-1}{n-(p+1)} \right) (1 - R^2)$$

Note-se que a inclusão de mais variáveis diminui o valor de R_a^2 , pois aumenta p , e não traz muito “incremento” a R^2 . Ou seja, ao contrário do coeficiente de determinação R^2 , o coeficiente de determinação ajustado, R_a^2 , não aumenta sempre quando adicionamos uma nova variável. Aliás, se adicionarmos variáveis com pouco poder explicativo este tende a decrescer. Pelo que, quando existe uma diferença significativa entre R^2 e R_a^2 , estamos perante uma situação em que provavelmente tenham sido incluídas no modelo variáveis estatisticamente não significativas.

4.5. Análise de Resíduos

Como vimos nos capítulos anteriores, os resíduos são dados pela diferença entre os valores da variável resposta observada e a variável resposta estimada, isto é,

$$e_i = y_i - \hat{y}_i, \quad i = 1, \dots, n.$$

A análise de resíduos pretende verificar se o modelo de regressão que está a ser utilizado é adequado. Portanto, os resíduos devem verificar os pressupostos

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon},$$

Como $\underline{Y} = (y_1, \dots, y_n)'$, \underline{X} a matriz do modelo, $\underline{\beta} = (\beta_1, \dots, \beta_p)'$ e $\underline{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_p)'$,

- a) ε_i , $i = 1, \dots, n$ são normalmente distribuídos;
- b) $var(\varepsilon_i) = \sigma^2$, $i = 1, \dots, n$, têm variância constante (homoscedasticidade);
- c) ε_i e $\varepsilon_{j,i} \neq j$, são independentes;
- d) Não existem *Outliers* influentes.

No caso da regressão linear múltipla, para além destes pressupostos, é preciso ainda verificar se existe colinearidade ou multicolinearidade entre as variáveis explicativas.

A seguir são expostas algumas “técnicas” por forma a verificar estes pressupostos.

4.5.1. Diagnóstico de Normalidade

A normalidade dos resíduos pode ser analisada quer através de gráficos, quer usando alguns testes, nomeadamente através do:

- i. gráfico P-P plot dos resíduos;
- ii. histograma dos resíduos estandardizados;
- iii. teste de Kolmogorov-Smirnov;
- iv. teste de Shapiro-Wilk.

Observemos:

i. Gráfico P-P plot dos resíduos;

Neste gráfico, vamos visualizar a distribuição de probabilidades dos valores observados com os valores esperados, representada por uma diagonal, segundo uma distribuição normal.

Caso a normalidade se verifique, as observações registadas aproximam-se dessa diagonal, sem nenhum afastamento significativo.

A Figura 4.5 mostra o gráfico p-p plot de resíduos. Nesta situação, a normalidade é verificada já que os pontos se aproximam da reta.

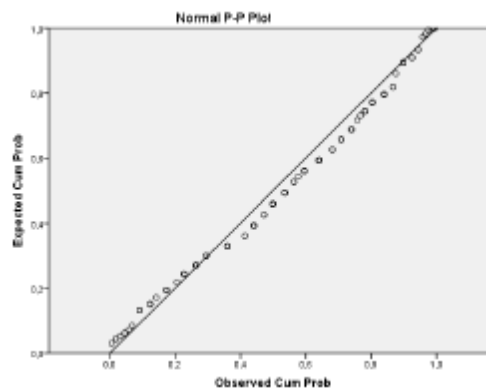


Figura 4.5: Normal p-p plot de resíduos

ii. Histograma dos resíduos estandardizados

Também se pode fazer um histograma dos resíduos no qual se procuram afastamentos evidentes em relação à forma simétrica e unimodal da distribuição normal. Este gráfico apenas deverá ser utilizado em amostras de dimensão elevada, já que quando se trabalha com amostras de dimensão reduzida o histograma não é muito conclusivo.

iii. Teste de Kolmogorov-Smirnov (K-S)

Neste caso, o teste de K-S é utilizado para testar as hipóteses:

$$\begin{cases} H_0: A \text{ distribuição é normal} \\ H_1: A \text{ distribuição não é normal,} \end{cases}$$

A estatística de teste, é dada por, ver [26],

$$D = \max\{\max(|F(x_i) - F_0(x_i)|); \max(|F(x_i - 1) - F_0(x_i)|)\}$$

em que $F(x_i) - F_0(x_i)$ representa a diferença entre a frequência acumulada de cada uma das observações e a frequência acumulada que essa observação teria, sendo a sua distribuição normal.

Este teste observa a máxima diferença absoluta entre a função de distribuição acumulada assumida pelos dados, neste caso da distribuição normal, e a função de distribuição empírica dos dados.

iv. Teste de Shapiro-Wilk (S-W)

Este teste sugere-nos preferência em relação ao teste de $K - S$ para amostras de pequenas dimensões ($n < 30$). Neste caso, as hipóteses a serem testadas são as definidas anteriormente para o teste de $K - S$.

A estatística de teste é definida da seguinte forma:

$$w = \frac{(\sum_{i=1}^n a_i x_i)^2}{(\sum_{i=1}^n x_i - \bar{x}_i)^2}$$

- a_i são constantes geradas a partir da média, variância e covariância de n ordens, ver [26].

4.5.2. Diagnóstico de Homoscedasticidade (Variâncias constantes)

Um dos pressupostos do modelo de regressão linear é a de que os erros devem ter variância constante. Esta condição é designada por homoscedasticidade.

A variância ser constante equivale a supor que não existem observações incluídas na variável residual cuja influência seja mais intensa na variável dependente.

Uma das técnicas, usadas para verificar a suposição de que os resíduos são homoscedásticos, é a análise do gráfico dos resíduos versus valores ajustados. Este gráfico deve apresentar pontos

dispostos aleatoriamente sem nenhum padrão definido, como se pode ver, por exemplo, na Figura 4.6

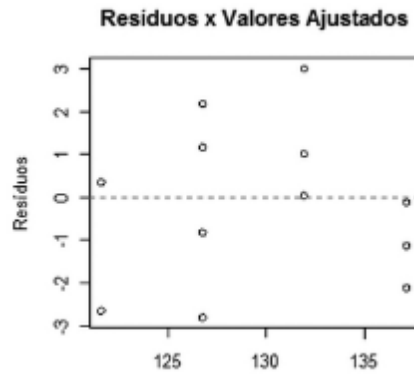


Figura 4.6: Confirmação da homoscedasticidade dos resíduos

4.5.3. Diagnóstico de Independência

Para testar o pressuposto da independência dos resíduos ou a presença de autocorrelação entre eles, pode utilizar-se o teste de Durbin-Watson (DW).

O teste de Durbin-Watson testa as hipóteses:

$$\begin{cases} H_0: \text{Não existe autocorrelação dos resíduos} \\ H_1: \text{Existe autocorrelação dos resíduos} \end{cases},$$

A estatística de teste é dada por:

$$dw = \frac{(\sum_{i=2}^n e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$$

e toma valores entre zero e quatro, $0 \leq dw \leq 4$.

Esta estatística mede a correlação entre cada resíduo e o resíduo correspondente à observação imediatamente anterior.

Podemos tomar a decisão comparando o valor de d_w com os valores críticos d_L e d_U da tabela de Durbin-Watson, disponível no Anexo 1.

A tabela seguinte dá-nos as decisões a tomar em função dos valores críticos, d_L e d_U .

Tabela 4.2: Tabela de decisão em função de d_L e d_U

dw	Zona de Rejeição e de não Rejeição de H_0				
	$[0; d_L[$	$[d_L; d_U[$	$[d_U; 4 - d_U[$	$[4 - d_U; 4 - d_L[$	$[4 - d_L; 4[$
Decisão	Rejeitar H_0	Nada se pode concluir	Não Rejeitar H_0	Nada se pode concluir	Rejeitar H_0
	Auto-correlação positiva				Auto-correlação negativa

4.5.4. Diagnóstico de Outliers e Observações Influentes

De acordo com [29], Outliers são observações extremas que se encontram de tal forma afastadas da maioria dos dados que surgem dúvidas sobre se elas poderão ou não ter sido geradas pelo modelo proposto para explicar essa maioria dos dados.

Os Outliers podem ser classificados em severos ou moderados consoante o seu afastamento em relação às restantes observações. Os Outliers moderados encontram-se fora do intervalo $[Q_1 - 1,5Q; Q_1 + 1,5Q]$ e os Outliers severos encontram-se fora do intervalo $[Q_1 - 3Q; Q_1 + 3Q]$, em que Q_1 representa o 1º quartil dos dados e Q a amplitude interquartil, isto é, é a diferença entre o 3º e o 1º quartil, $Q = Q_3 - Q_1$.

Se um Outlier for influente vai interferir sobre a função de regressão ajustada o que significa que a inclusão ou não desse ponto modifica substancialmente os valores ajustados. Assim, um ponto é influente se a sua exclusão na regressão ajustada provoca uma mudança substancial nos valores ajustados.

Uma medida que serve para diagnosticar Outliers é Leverage (LEV). Para uma dada observação, um Leverage elevado indica que essa observação se distancia do centro das observações, exercendo influência sobre o valor previsto. O Leverage varia entre 0 e 1.

Acontece que um elevado Leverage indica apenas que a observação poderá ser influente.

Considera-se um Leverage elevado quando, ver [30],

$$LEV > \frac{3(p + 1)}{n}, \text{ para amostras de dimensão reduzida}$$

$$LEV > \frac{2(p + 1)}{n}, \text{ para amostras grandes}$$

onde n é a dimensão da mostra e p o número de variáveis independentes.

4.5.4.1. Observações Influentes

As observações influentes são aquelas que individualmente ou em conjunto com as outras observações demonstram ter mais impacto do que as restantes no cálculo dos estimadores.

Nesta subsecção, apresentamos várias medidas que são utilizadas para identificar as observações influentes.

1) SDFFIT

É uma das medidas de utilização mais frequente para medir a influência de cada observação. SDFFIT trata-se de uma medida estandardizada que mede a influência que a observação i tem sobre o seu valor ajustado.

Considera-se que uma observação é influente se, ver [30].,

$$\text{SDFFIT} > 2 \sqrt{\frac{p + 1}{n - p - 1}}.$$

2) SDFBETA

A influência que uma observação tem sobre a estimação de cada um dos coeficientes de regressão pode ser calculada pelo SDFBETA. Trata-se de uma medida estandardizada que corresponde à alteração nos coeficientes estimados, $\hat{\beta}_j, j = 0, \dots, p$, quando se exclui essa observação.

Neste caso, a observação é influente quando

$$\text{SDFFIT} > 1,96, \text{ para amostras de dimensão reduzida}$$

$$\text{SDFFIT} > \frac{2}{\sqrt{n}}, \text{ para amostras grandes}$$

3) Para verificar se uma observação é influente também podemos usar a distância deCook que mede a influência da i -ésima observação sobre todos os n valores ajustados

$$\hat{y}_i, i = 1, \dots, n.$$

Uma distância de Cook elevada significa que o resíduo e_i é elevado, ou a Leverage para essa observação é elevada ou ambas as situações.

De tal forma que uma observação é influente quando, ver [30],

$$\text{COOK} > \frac{4}{n - p - 1}$$

em que n é a dimensão da amostra e o p o número de variáveis independentes. Considera-se que observações com Distância de Cook superior a 1 são excessivamente influentes.

4.5.5. Colinearidade e Multicolinearidade

O termo colinearidade é utilizado para expressar a existência de correlação elevada entre duas variáveis independentes, enquanto o termo multicolinearidade é utilizado quando se trata de mais do que duas variáveis independentes fortemente correlacionadas. No entanto, existem autores que definem colinearidade como a existência de relação linear entre duas variáveis independentes e multicolinearidade como a existência de relação linear entre uma das variáveis independentes e as restantes.

Se considerarmos duas quaisquer variáveis independentes, X_1 e X_2 , entre as quais existe uma elevada correlação, a proporção da variação total da variável dependente, explicada por X_1 é idêntica à proporção da variação total da variável dependente, explicada por X_2 .

A colinearidade poderá ser diagnosticada:

- verificando se a matriz de correlações das variáveis independentes demonstra correlações elevadas. Caso a correlação de duas variáveis seja muito próxima de 1, indica de facto um problema;
- Verificando se, ao se realizar a regressão de X_i em função das outras variáveis independentes, o valor de $R^2 \cong 1$.

Um indicador usado com frequência para detetar a multicolinearidade é o Variance Inflation Factor (VIF).

A variância de cada um dos coeficientes de regressão associados às variáveis independentes é dada por, ver [26, 27]

$$\text{var}(\hat{\beta}_i) = \sigma^2 \left(\frac{1}{1 - R_i^2} \right) \times \frac{1}{\sum_{j=1}^n (x_{ij} - \bar{x}_i)^2},$$

em que R_i^2 é o R^2 de regressão de X_i sobre as restantes variáveis explicativas.

Esta variância é tanto maior quanto maior for a correlação múltipla entre X_i e as variáveis independentes.

O termo $\frac{1}{1-R_i^2}$ designa-se, em concreto, por VIF para o coeficiente de regressão X_i associado à variável X_i .

Segundo [26], caso se obtenham valores de $VIF > 5$ conclui-se que estamos perante problemas com a estimação de β_i devido à presença de multicolinearidade nas variáveis independentes.

Suponhamos que temos a equação de regressão

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon,$$

em que X_1 e X_2 são altamente correlacionadas.

Numa situação deste género, devíamos eliminar uma das variáveis e reestimar o modelo.

Existem vários métodos que permitem, na regressão linear múltipla, fazer uma seleção das variáveis independentes que melhor explicam a variável resposta, nomeadamente:

. FORWARD – o método começa apenas com a constante e adiciona uma variável independente de cada vez. A primeira variável selecionada é a que apresenta maior correlação com a variável resposta (maior score statistic)

. BACKWARD – o método faz o “contrário” do método Forward. Neste caso, todas as variáveis independentes são incorporadas no modelo. Depois, por etapas, cada uma pode ser ou não eliminada.

STEPWISE - o método Stepwise é uma “modificação” do método Forward que permite resolver problemas de multicolinearidade. Consiste no seguinte: fazemos entrar no modelo a variável explicativa que apresenta maior coeficiente de correlação com a variável dependente.

Em seguida, calculam-se os coeficientes de correlação parcial para todas as variáveis que não fazem parte da primeira equação de regressão, para que, a próxima variável a entrar, seja a que apresenta maior coeficiente de correlação parcial.

Estima-se a nova equação de regressão e analisa-se se uma das duas variáveis independentes deve ser excluída do modelo.

No final, se ambas as variáveis apresentarem valores t significativos, novos coeficientes de correlação parcial são calculados para as variáveis que não entraram.

Este processo finda, assim que se chegue à situação em que nenhuma variável deva ser acrescentada à equação.

4.5.6. Análise de Variância Multivariada (MANOVA) aplicada à Regressão

A Análise de Variância Multivariada (MANOVA) é uma generalização da Análise da Variância (ANOVA), sendo utilizada para a comparação entre médias nos casos em que existem pelo menos duas variáveis critério, ao passo que a ANOVA apenas permite a comparação de uma variável critério.

Tal como na ANOVA, na MANOVA primeiro testa-se a hipótese de igualdade global entre as médias dos grupos, caso o resultado do teste seja significativo realizam-se testes complementares com o objetivo de explicar as diferenças entre os grupos.

Na análise simples da variância, ANOVA, as médias das variáveis são testadas separadamente, não sendo consideradas eventuais covariâncias entre as variáveis. Na análise multivariada, MANOVA, as hipóteses de testes têm em consideração um conjunto de variáveis, p , em simultâneo.

Assim, as hipóteses de teste têm por objetivo testar a igualdade de k vetores de médias, sendo dadas por:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k, \text{ com } \mu_j = \begin{bmatrix} \mu_{1j} \\ \mu_{2j} \\ \dots \\ \mu_{pj} \end{bmatrix} \text{ e } j = 1, 2, \dots, k$$

Todos os grupos têm vetores de médias iguais

$$H_1: \mu_i \neq \mu_j, \text{ com } i \neq j$$

Existem pelo menos dois grupos que não têm vetores de médias iguais.

Os pressupostos necessários para se poder aplicar a MANOVA são generalizações dos pressupostos para aplicação da ANOVA, isto é:

- As populações, de onde são retiradas as amostras, devem ter uma distribuição normal.
- Os grupos populacionais de onde são retiradas as amostras devem ter idêntica variância.

Acresce que, as amostras devem ser recolhidas de forma aleatória e ser independentes entre si.

A estatística de teste a utilizar na ANOVA encontra-se perfeitamente definida, sendo o quociente entre a média dos quadrados entre os grupos e a média dos quadrados dentro dos grupos, que segue uma distribuição F . Quanto à estatística de teste a utilizar na MANOVA, a resposta não é assim tão simples.

Existem várias estatísticas de teste possíveis de utilizar na MANOVA, sendo a mais comum e a mais antiga a estatística Λ de Wilks. Podem, no entanto, ser utilizados outros testes tais como Pillai, Hotelling-Lawley e o teste de Roy, podendo inclusive apresentar resultados diferentes para a mesma análise.

Dado ser a estatística Λ de Wilks a mais usada, sendo uma distribuição de probabilidade muito usada em testes de hipóteses multivariada, em especial nos que dizem respeito à razão de verossimilhança e análise de variância multivariada. Trata-se de uma generalização multivariada da distribuição F .

Resumo do teste Λ de Wilks no quadro seguinte.

Tabela 4.3: MANOVA – Teste Λ de Wilks

Fontes de Variação	Graus de Liberdade	Matriz de Soma dos Quadrados	Λ
Entre os Grupos	$K - 1$	$B^1 = \sum_j n_j (\bar{x}_j - \bar{x})(\bar{x}_j - \bar{x})'$	$\Lambda = \frac{ W }{ B + W }$ $\Lambda \cap \Lambda_{(p;n-k;k-1)}$
Dentro dos Grupos	$n - k$	$W^2 = \sum_j \sum_{\mu} (x_{j\mu} - \bar{x}_j)(x_{j\mu} - \bar{x}_j)'$	
Total	$n - 1$	$T = B + W = \sum_j \sum_{\mu} (x_{j\mu} - \bar{x})(x_{j\mu} - \bar{x})'$	

Sendo,

- p = número de variáveis.
- k = número de indivíduos
- n = número de grupos.

Usando o teste Λ de Wilks, na presença de diferenças sistemáticas entre os tratamentos, espera-se sempre obter Λ menor que um, sendo que este é tão mais significativo quanto menor for o seu valor.

A regra de decisão deste teste é a seguinte:

- Rejeita-se a hipótese nula ao nível de significância α se $\Lambda_{calculado} > \Lambda_{(p;n-k;k-1)}$
- Caso contrário, aceita-se a hipótese nula, dizendo que o teste é significativo com um nível de significância α .

¹ Between

² Within

4.6. Construção do Modelo de Regressão Linear Simples

Com vista a perceber se a concentração de propofol C_e , (ug/ml), influencia o sinal BIS, foi realizada uma análise de regressão linear simples.

O modelo de regressão linear simples que representa a relação entre a variável dependente, BIS, e a variável independente, propofol C_e (ug/ml), é dado pela seguinte equação:

$$BIS = \beta_0 + \beta_1 Propofol + \epsilon \tag{4.14}$$

Tabela 4.4: Estatísticas descritivas			
	N	Média	Desvio Padrão
BIS	4214	75,01	15,926
Propofol C_e (ug/ml)	4214	2,41304880	1,145806379
N válido (de lista)	4214		

A Tabela 4.4, Estatística descritiva, mostra o valor médio e o desvio padrão de Propofol e do sinal BIS. Concluimos que, nesta amostra, a concentração de propofol é $2,41304880 \pm 1$ (ug/ml), enquanto que o do sinal BIS é $75,01 \pm 16$.

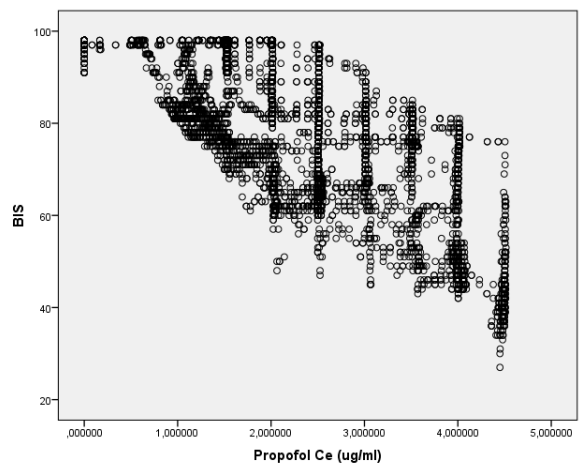


Figura 4.7: Diagrama de dispersão

Elaborámos um diagrama de dispersão com o intuito de perceber se a relação existente entre as duas variáveis é de facto linear.

De acordo com a observação do diagrama de dispersão (Figura.4.7) somos tentados a concluir que não existe uma relação linear entre o propofol e o sinal *BIS*.

Analisando a Tabela 4.5, podemos afirmar que a correlação existente entre as variáveis é positiva moderada ($R = ,736$).

Tabela 4.5:Resumo do modelo^b

Modelo	R	R quadrado	R quadrado ajustado	Erro padrão da estimativa	Durbin-Watson
1	,736 ^a	,541	,541	10,788	,036

a. Preditores: (Constante), Propofol Ce ($\mu\text{g/ml}$)

b. Variável Dependente: BIS

A tabela 4.5 apresenta o sumário do modelo, com a variável (Variável preditora: Constante), propofol Ce ($\mu\text{g/ml}$) cuja construção será feita posteriormente.

Neste modelo, encontramos o coeficiente de determinação ajustado $R^2_{ajust}=0,541$, donde podemos afirmar que 54,1% da variabilidade da variável dependente BIS é explicada pela variável independente do modelo ajustado, sendo a restante variabilidade 45,9% explicada por fatores não incluídos no modelo. Podemos considerar que este modelo não é um bom ajuste pois o valor de R^2_{ajust} não se encontra próximo de 1. . O valor do coeficiente de correlação múltipla é $R=0,736$. Existe pois uma relação positiva entre as variáveis.

Tabela 4.6: ANOVA^a (Variáveis predictoras: Constante)

Modelo		Soma dos Quadrados	df	Quadrado Médio	Z	Sig.
1	Regressão	578297,807	1	578297,807	4968,614	0,0E0
	Resíduo	490235,395	4212	116,390		
	Total	1068533,202	4213			

a. Variável Dependente: BIS

Para efetuar a análise de variância do modelo recorreu-se ao teste de F que tem associado o seguinte $p - value$ (sig) de 0,000. De acordo com o seu valor, rejeitamos $H_0: \beta_1 = 0$, pelo que podemos dizer que o modelo é **altamente significativo**.

Tabela 4.7: Coeficientes^a

Modelo		Coeficientes				
		Coeficientes não padronizados		padronizados		
		B	Erro Padrão	Beta	t	Sig.
1	(Constante)	99,687	,387		257,265	0,0E0
	Propofol Ce ((µg/ml))	-10,225	,145	-,736	-70,488	0,0E0

a. Variável Dependente: BIS

b. Preditores: (Constante), Propofol Ce (µg/ml)

Assim, a equação que exprime a relação entre o sinal BIS e a concentração propofol Ce (µg/ml) é a seguinte:

$$BIS = 99,687 - 10,225 \text{ propofol} \quad (4.15)$$

O teste ao coeficiente de regressão β_0 é dado pelo teste $t - student$ ao qual está associado um valor de significância de 0,000 (< 0,05). Concluimos, portanto, que se deve rejeitar a hipótese $H_0: \beta_0 = 0$, o que significa que a reta ajustada não passa pela origem.

Quanto ao teste para β_1 é dado pelo teste $t - student$ ao qual está associado um valor de significância de 0,000 (< 0,05). Logo rejeitar a hipótese $H_0: \beta_1 = 0$, o que significa que a variável propofol influencia significativamente o BIS.

Verificação dos pressupostos do modelo

O modelo definido em (4.15) só será adequado se validados todos os pressupostos. Vamos nesta subsecção fazer uma análise desses pressupostos.

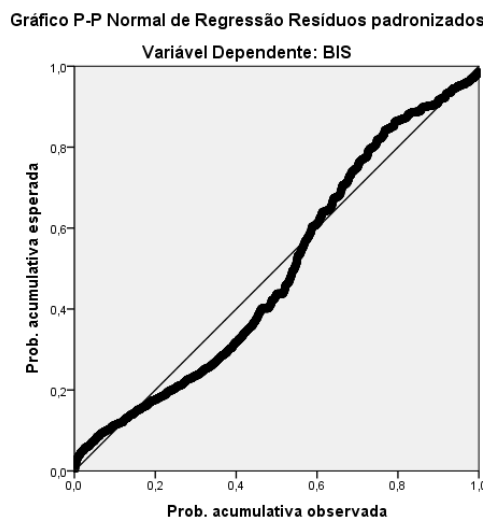


Figura 4.8: Gráficos dos Resíduos versus preditos; resíduos padronizados e da probabilidade normal dos resíduos

A partir da análise da Figura 4.8, podemos concluir que as observações se aproximam da reta sem nenhum afastamento sistemático, pelo que somos levados a concluir que os resíduos não são normalmente distribuídos.

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Estatística	df	Sig.	Estatística	df	Sig.
Propofol Ce (ug/ml)	,080	4214	2,091E-72	,970	4214	1,0643E-28

a. Correlação de Significância de Lilliefors

O p-value (exato) é 2,0881E-20, logo rejeitamos a hipótese de que a variável em estudo segue uma distribuição normal para o nível de significância de $\alpha = 0,05$.

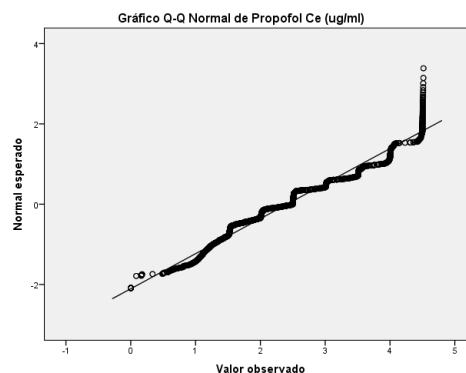


Figura 4.9: Q-Q plot

Usando a variável RES (os resíduos guardados) e fazendo um Q-Q plot (figura 4.9.) e os testes de ajustamento de Kolmogorov-Smirnov e de Shapiro-Wilk podemos concluir que os resíduos não têm uma distribuição normal (o Q-Q plot identifica um ajuste entre os quantis amostrais e os quantis de distribuição normal) e os testes de ajustamentos fornecem os p-values inferiores aos níveis de significância usual de 0,05, por isso dará um teste significativo.

Mas de acordo com o Teorema do Limite Central a dimensão da amostra suficientemente grande e quando a variação da população tem variância finita, considera satisfatório a aproximação da média de X à normal quando $n \geq 30$, que é a situação das amostras em análise onde $n=4214$.

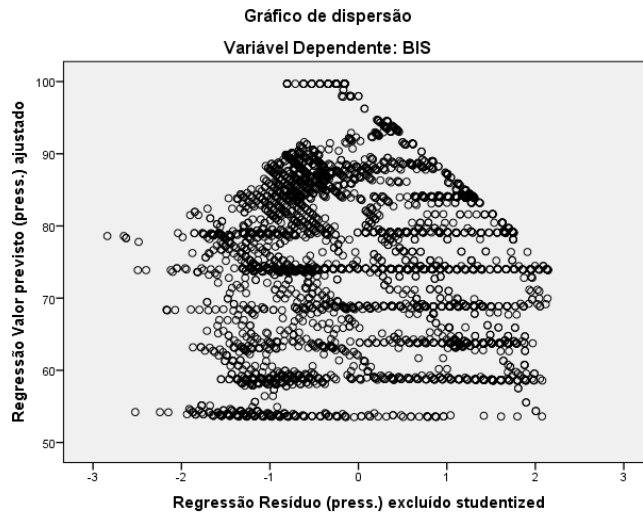


Figura 4.10: Gráficos dos Resíduos press

Pela análise gráfica dos resíduos estandardizados (Figura 4.8) e dos resíduos press (Figura 5.10) podemos concluir que existem Outliers, dado que há resíduos que apresentam valores absolutos superiores a 1,96, como mostra a Figura 5.10 [26].

Tabela 4.9: Estatísticas de resíduos ^a					
	Mínimo	Máximo	Média	Desvio Padrão	N
Valor previsto	53,51	99,69	75,01	11,716	4214
Valor Previsto Padrão	-1,835	2,106	,000	1,000	4214
Erro padrão do valor previsto	,166	,387	,227	,060	4214
Valor previsto ajustado	53,51	99,70	75,01	11,716	4214
Resíduo	-30,591	23,097	,000	10,787	4214
Resíduo Padronizado	-2,836	2,141	,000	1,000	4214
Resíduos Estudantizados	-2,836	2,141	,000	1,000	4214
de Estud.	-30,599	23,103	,000	10,792	4214
Resíduos deletados					
Estudantizados	-2,838	2,142	,000	1,000	4214
Mahal. Distância	,000	4,435	1,000	1,132	4214
Distância de Cook	,000	,003	,000	,000	4214
Valor de ponto alavanca centralizado	,000	,001	,000	,000	4214

a. Variável Dependente: BIS

A confirmação da existência de Outliers pode ser feita, por exemplo, através do valor máximo do *student delete residual* $2,14 > 1,96$.

E também o valor de ponto alavanca centralizado máximo, que é igual a $0,001 > \frac{2(p+1)}{n} = 0,0009$. Interessa verificar se *Outliers* são ou não observações influentes. Olhando ainda para a

tabela 4.9 tudo leva a crer que sim, já que temos como valor máximo da Distância de Cook $0,003 > \frac{4}{n-p-1} = 0,0009$.

Vamos usar mais uma técnica para averiguar a existência de Observações Influentes, recorrendo à análise dos SDFFIT. Visto que as observações da Figura 4.9 SDFFIT $> 2 \sqrt{\frac{p+1}{n-p-1}} = 0,044$. Concluimos que se tratam de observações influentes, devendo ser mantidas no estudo como é sugerido, por exemplo, por [30].

Tabela 4.10: Verificação da multicolinearidade (Variáveis dependente: BIS)

Modelo		Coeficientes				Estatísticas de colinearidade	
		Coeficientes não padronizados		Coeficientes padronizados		Tolerância	VIF
		B	Erro Padrão	Beta	t	Sig.	
1	(Constante)	99,687	,387		257,265	,000	
	Propofol Ce (ug/ml)	-10,225	,145	-,736	-70,488	,000	1,000 1,000

a. Variável Dependente: BIS

Para avaliar a multicolinearidade, o SPSS utiliza a Tolerância de cada variável que é a medida da proporção da variância da variável que não é explicada pelas restantes variáveis independentes e que é calculada aquando da aplicação do método Stepwise, onde se vão selecionar as variáveis que vão entrar na análise. As variáveis do nosso estudo que se encontram nestas condições (Tolerância $> 0,8$) são propofol, assim revelando poder discriminante, pelo que não há a recear a violação do pressuposto de multicolinearidade.

Tabela 4.11: Diagnóstico de colinearidade (Variável Dependente: BIS)

Modelo	Dimensão	Autovalor	Índice de condição	Proporções de variância	
				(Constante)	Propofol Ce (ug/ml)
1	1	1,903	1,000	,05	,05
	2	,097	4,438	,95	,95

a. Variável Dependente: BIS

4.7. Construção do Modelo de Regressão Múltipla

Neste estudo, passamos a considerar o modelo de regressão linear múltipla, que será “estimado” através do método Stepwise. Pretendemos averiguar se a concentração de propofol C_e (ug/ml) e de remifentanil C_e (ng/ml), influenciam o sinal BIS, tendo sido realizada uma análise de regressão linear múltipla.

Tabela 4.12: Variáveis Inseridas/Removidas^a

Modelo	Variáveis inseridas	Variáveis removidas	Método
1	Remifentanil C_e (ng/ml), Propofol C_e (ug/ml) ^b	.	Inserir

a. Variável Dependente: BIS

b. Todas as variáveis solicitadas inseridas.

Tabela 4.13: Resumo do modelo

Modelo	R	R quadrado	R quadrado ajustado	Erro padrão da estimativa
1	,836 ^a	,698	,698	8,748

a. Preditores: (Constante), Remifentanil C_e (ng/ml), Propofol C_e (ug/ml)

A tabela 4.12 apresenta o sumário do modelo, com as variáveis (Variáveis predictoras: Constante), remifentanil C_e (ng/ml), propofol C_e (ug/ml) cuja construção será feita posteriormente.

Neste modelo, encontramos $R^2_{ajust} = 0,698$, donde podemos afirmar que 69,8% da variabilidade da variável dependente BIS é explicada pelas variáveis independentes do modelo ajustado. Podemos considerar que este modelo é um bom ajuste pois o valor de r^2_{ajust} encontra-se próximo de 1. O valor do coeficiente de correlação múltipla é $R=0,836$. Existe pois uma forte relação entre as variáveis.

Tabela 4.14: ANOVA^a

Modelo	Soma dos Quadrados	df	Quadrado Médio	Z	Sig.
1 Regressão	746311,096	2	373155,548	4876,630	,000 ^b

Resíduo	322222,106	4211	76,519
Total	1068533,202	4213	

a. Variável Dependente: BIS

b. Preditores: (Constante), Remifentanil Ce (ng/ml), Propofol Ce (ug/ml)

Para efectuar a análise de variância do modelo recorreu-se ao teste de F que tem associado o seguinte $p - value (sig)$ de 0,000. De acordo com o seu valor, rejeitamos $H_0: \beta_1 = 0$, pelo que podemos dizer que o modelo é **altamente significativo**.

O modelo ajustado (tabela 4.14) é dado por:

Tabela 4.15: Coeficientes (Variáveis dependente: BIS)

Model		Unstandardized		Standardized		Sig.
		Coefficients		Coefficients		
		B	Std. Error	Beta	t	
1	(Constant)	102,190	,319		320,654	,000
	Propofol Ce (ug/ml)	-8,741	,122	-,629	-71,762	,000
	Remifentanil Ce (ng/ml)	-10,464	,223	-,411	-46,858	,000

a. Dependent Variable: BIS

$$BIS = 102,190 - 8,741 \text{ Propofol Ce (ug/ml)} - 10,464 \text{ Remifentanil Ce (ng/ml)} \quad (4.15)$$

O teste ao coeficiente de regressão β_0 é dado pelo teste $t - student$ ao qual está associado um valor de significância de 0,000 ($< 0,05$). Concluimos, portanto, que se deve rejeitar a hipótese $H_0: \beta_0 = 0$, o que significa que a reta ajustada não passa pela origem.

Quanto ao teste para β_1 é dado pelo teste $t - Student$ ao qual está associado um valor de significância de 0,000 ($< 0,05$). Logo rejeitar a hipótese $H_0: \beta_1 = 0$, o que significa que as variáveis propofol e remifentanil influenciam significativamente o BIS e o modelo melhora com a inclusão do remifentanil, sendo aquela que tem maior contribuição individual (-46,858).

Verificação dos pressupostos do modelo

O modelo definido em (4.15) só será adequado se validados todos os pressupostos. Vamos nesta subsecção fazer uma análise desses pressupostos.

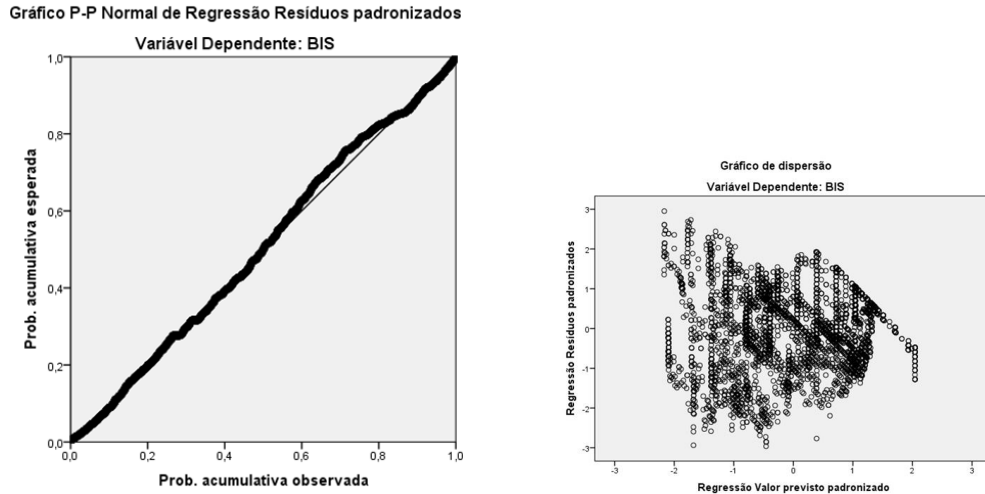


Figura 4.11: Gráficos dos Resíduos versus preditos; resíduos padronizados e da probabilidade normal dos resíduos (RLM)

Procedamos pois à realização de testes exatos da ocorrência de outliers e ao pressuposto da distribuição normal. Assim, para testar a normalidade, optou-se por usar o teste Kolmogorov-Smirnov que é um teste paramétrico tradicional, baseado na distribuição t-Student e é obtido sob a hipótese de que a população tem distribuição normal, e o teste de Shapiro-Wilk, por se tratar de um teste de ajustamento específico para a distribuição normal que tem uma melhor performance que o teste anterior em amostras reduzidas ($n < 30$).

Tabela 4.16: Teste de Kolmogorov-Smirnov de uma amostra

		Studentized Deleted Residual
N		4214
Parâmetros normais ^{a,b}	Média	-,0000095
	Erro Desvio	1,00033285
Diferenças Mais Extremas	Absoluto	,044
	Positivo	,017
	Negativo	-,044
Estatística de teste		,044
Significância Sig. (2 extremidades)		2,0881E-20

- a. A distribuição do teste é Normal.
- b. Calculado dos dados.
- c. Correção de Significância de Lilliefors.

O p-value (exato) é 2,0881E-20, logo rejeitamos a hipótese de que a variável em estudo segue uma distribuição normal para o nível de significância de $\alpha = 0,05$.

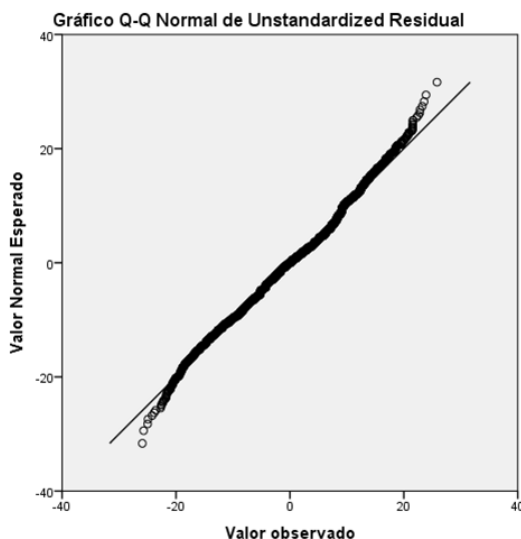


Figura 4.12: Q-Q plot (RLM)

Usando a variável RES (os resíduos guardados) e fazendo um Q–Q plot (figura 4.12) e os testes de ajustamento de Kolmogorov-Smirnov e de Shapiro-Wilk, podemos concluir que os resíduos não têm uma distribuição normal o (Q–Q plot identifica um ajuste entre os quantis amostrais e os quantis de distribuição normal) e os testes de ajustamentos fornecem os p-values inferiores aos níveis de significância usual de 0,05, por isso dará um teste significativo.

Tabela 4.17: Testes de Normalidade

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Estatística	df	Sig.	Estatística	df	Sig.
Unstandardized Residual	,044	4214	1,9099E-20	,995	4214	1,5202E-11

a. Correlação de Significância de Lilliefors

Mas, de acordo com o Teorema do Limite Central, a dimensão da amostra suficientemente grande e quando a variação da população tem variância finita, considera satisfatória a aproximação da média de X à normal quando $n \geq 30$, que é a situação das amostras em análise, onde $n=4214$.

Tabela 4.18: Verificação da multicolinearidade (Variáveis dependente: BIS)

Modelo	Coeficientes não padronizados		Coeficientes padronizados		Estatísticas de colinearidade			
	B	Erro		Beta	t	Sig.	Tolerância	
		Padrão					a	VIF
1 (Constante)	102,190	,319			320,654	,000		
Propofol Ce (ug/ml)	-8,741	,122	-,629	71,762		,000	,932	1,072
Remifentanil Ce (ng/ml)	-10,464	,223	-,411	46,858		,000	,932	1,072

a. Variável Dependente: BIS

Para avaliar a multicolinearidade, o SPSS utiliza a Tolerância de cada variável que é a medida da proporção da variância da variável que não é explicada pelas restantes variáveis independentes e que é calculada aquando da aplicação do método Stepwise, onde se vão seleccionar as variáveis que vão entrar na análise. As variáveis do nosso estudo que se encontram nestas condições (Tolerância > 0,8) são propofol e remifentanil, assim revelando poder discriminante, pelo que não há a recear a violação do pressuposto de multicolinearidade.

Tabela 4.19: Diagnóstico de colinearidade (Variável Dependente: BIS)

Modelo	Dimensão	Autovalor	Índice de condição	Proporções de variância		
				(Constante)	Propofol Ce (ug/ml)	Remifentanil Ce (ng/ml)
1	1	2,526	1,000	,02	,02	,06
	2	,378	2,585	,08	,06	,94
	3	,096	5,121	,90	,92	0,004531

a. Variável Dependente: BIS

Resumo da Análise de Regressão Linear Simples e Múltipla

Da nossa análise, podemos concluir que o teste de significância da equação de Regressão Linear indicou que o modelo construído pode ser considerado significativo para um nível de significância de 5% e, conseqüentemente, o modelo de regressão é válido. Em suma, **o modelo é altamente significativo**. Como o *p-value* encontrado foi inferior a 0,05, podemos assegurar que

o modelo de regressão considerado é melhor que a média para predizer os valores do BIS. Com um nível de confiança de 95%, a variável remifentanil é a mais significativa, sendo aquela que tem maior contribuição individual (-46,858). Os resíduos têm uma distribuição normal e os testes de ajustamentos fornecem os p-values inferiores aos níveis de significância usual de 0,05, o modelo foi declarado ajustado.

O modelo RLS o $R_{ajust}^2 = 0,541$, donde podemos afirmar que 54,1% da variabilidade da variável dependente BIS é explicada pela variável independente do modelo ajustado, sendo a restante variabilidade explicada por fatores não incluídos no modelo.

O modelo RLM o $R_{ajust}^2 = 0,698$, donde podemos afirmar que 69,8% da variabilidade da variável dependente BIS é explicada pelas variáveis independente do modelo ajustado, sendo a restante variabilidade explicada por fatores não incluídos no modelo. As variáveis propofol e remifentanil influenciam significativamente o BIS e o modelo melhora com a inclusão do remifentanil, sendo aquela que tem maior contribuição individual (-46,858).

Para avaliarmos a qualidade do modelo, podemos comparar a variação do sinal BIS que é explicada pelo modelo, com a variação do que não é explicada pelo modelo e o modelo será tanto melhor quanto maior for este quociente (R_{ajust}^2). Assim, podemos concluir que o melhor modelo é o que apresenta um $R_{ajust}^2 = 0,698$ Modelo RLM.

Superfície de Resposta ajustada entre remifentanil , propofol e BIS

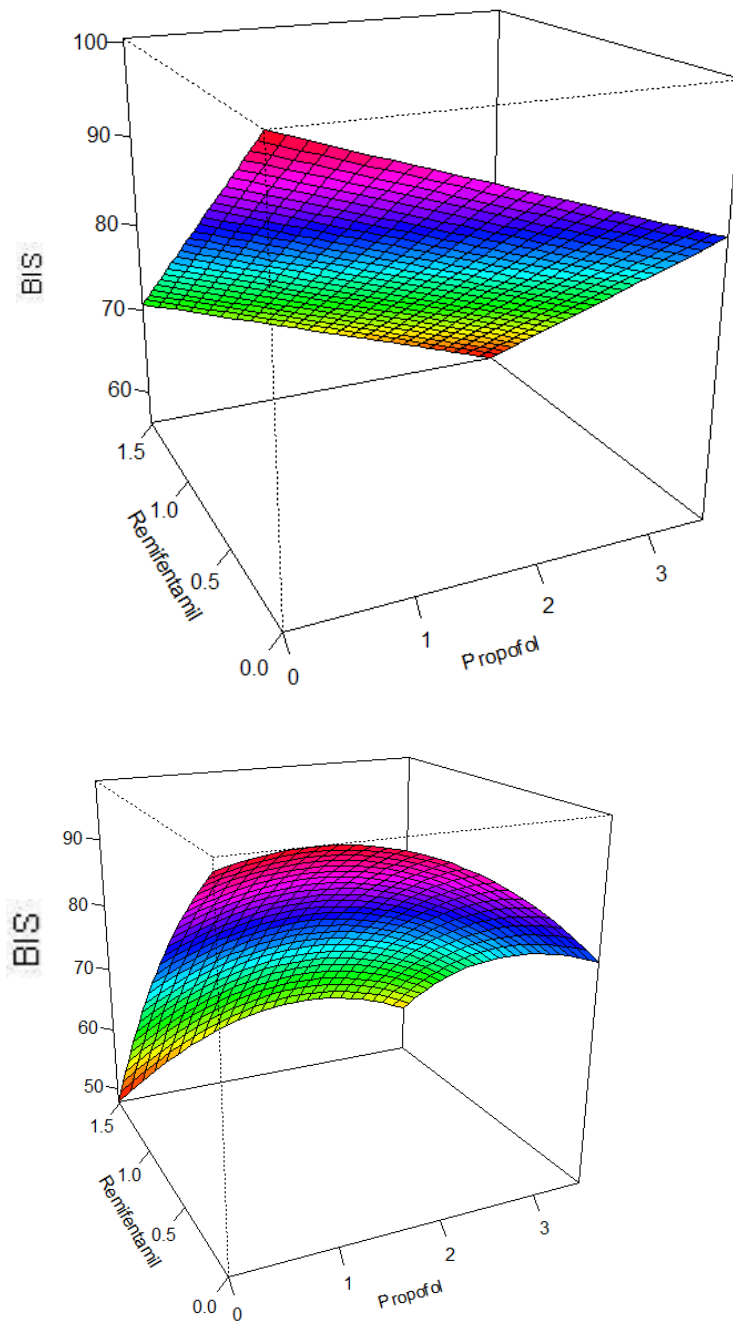


Figura 4.13:Superfície de Resposta ajustada entre remifentanil e propofol e BIS do modelo de 1º Ordem. e 2º Ordem.

Capítulo 5

Análise de Clusters

5.1. Introdução

A análise de *Clusters* é uma técnica de análise multivariada que tem como principal objetivo o agrupamento de elementos [39]. Neste caso específico, o objetivo passa por agrupar as concentrações dos fármacos, com base no efeito no sinal cerebral BIS. Pretendendo-se captar diferentes tipos de interações farmacocinéticas/dinâmicas que possam provocar efeitos semelhantes Esta análise permite a formação de grupos distintos onde a variação dentro do grupo é baixa e entre grupos é alta [51].

Este agrupamento é efetuado de forma que elementos pertencentes ao mesmo grupo tenham características semelhantes e elementos de diferentes grupos tenham características dissemelhantes. Genericamente, parte-se de um conjunto com n observações e pretende-se formar k grupos com um menor número de observações.

Assim, para detetar grupos homogéneos de concentrações de fármacos e o efeito no sinal cerebral BIS, usam-se métodos hierárquicos ou métodos não hierárquicos aos 4214 casos. Entre os métodos hierárquicos, estes podem ser aglomerativos ou divisivos. Devido à natureza intrínseca de cada variável, em que cada uma tem uma escala de medição, temos que ter em conta o interesse ou não da amplitude das mesmas, ou seja, podemos ter um agrupamento por amplitude ou por padrão. A primeira técnica, agrupamento por amplitude, tem em consideração as diferentes escalas de medição, ou seja, as variáveis com maior amplitude terão uma maior contribuição. Esta técnica tende a agrupar as variáveis com maior amplitude entre elas. A segunda técnica, agrupamento por padrão, considera todas as variáveis com a mesma escala, ou seja, o agrupamento é feito em variáveis estandardizadas. Esta técnica tende a agrupar os casos com base na natureza das variáveis. Dado o objectivo deste trabalho, a segunda técnica consiste em caracterizar e classificar casos de acordo com padrões, consideramos a técnica de agrupamento por padrão, tendo para tal estandardizado as variáveis em estudo. A primeira técnica determinar as concentrações de fármacos e o efeito no sinal cerebral BIS.

Na Análise de *Clusters*, os agrupamentos de sujeitos (casos ou itens) ou variáveis é feito a partir de medidas de semelhanças ou de medidas dissemelhança (distância) entre, inicialmente, dois sujeitos e mais tarde entre dois *Clusters* de observações, usando técnicas hierárquicas ou não-hierárquicas de agrupamento de *Clusters* [33].

Genericamente, a análise de *Clusters* compreende cinco etapas [34]:

1. A seleção de indivíduos ou de uma amostra de indivíduos a serem agrupados;
2. A definição de um conjunto de variáveis a partir das quais será obtida a informação necessária ao agrupamento dos indivíduos;
3. A definição de uma medida de semelhança ou distância entre cada dois indivíduos;
4. A escolha de um critério de agregação ou desagregação dos indivíduos, isto é, a definição de um algoritmo de partição / classificação;
5. Por último, a validação dos resultados encontrados.

Os algoritmos de agrupamento operam, geralmente, sobre dois tipos de estrutura de dados:

- (i) uma matriz de dimensão $n \times p$ correspondendo as n linhas aos sujeitos e as p colunas aos seus atributos ou características;
- (ii) quadro de dimensão $n \times n$ cujos elementos medem as proximidades³ entre cada par de indivíduos.

5.2. As variáveis

5.2.1. Seleção das variáveis

A escolha das variáveis adequadas para a definição de grupos pode estar relacionada com conhecimento anteriormente adquirido pelo investigador sobre o tema a estudar, o que permitirá, à partida, rejeitar as variáveis irrelevantes.

Variáveis que assumem praticamente o mesmo valor para todos os sujeitos são pouco discriminatórias e a sua inclusão pouco contribuiria para a determinação da estrutura do agrupamento. Por outro, a inclusão de variáveis com grande poder de discriminação, porém

³ As proximidades poderão ser semelhanças (medem o grau de similitude entre cada par de sujeitos) ou distâncias (medem o grau de afastamento ou diferença)

pouco significativas na abordagem do problema, pode mascarar os grupos e levar a resultados equivocados.

Acontece, com frequência, o número de variáveis medidas ser grande, dificultando a análise. Respeitando o princípio da parcimónia, devemos tentar diminuir o seu número de forma que a seleção considere tanto a sua relevância como o seu poder de discriminação face ao problema em estudo. Em último caso, pode-se ainda tentar utilizar técnicas estatísticas para redução da dimensionalidade da matriz de dados (e.g. a Análise de Componentes Principais).

5.2.3. Escala das variáveis

Quando as variáveis estão definidas em diferentes escalas de medida e se aplica a análise de *Clusters*, qualquer medida de semelhança ou de distância vai refletir sobretudo o peso das variáveis que maiores valores e maior dispersão apresentam. Visando anular este efeito, surgiram várias propostas de estandardização das variáveis [39]. Apresentam-se as mais comuns:

Consideremos as observações originais x_1, \dots, x_n

A transformação mais comum é feita por

$Z_i = \frac{x_i - \bar{x}}{s}$, $i = 1, \dots, n$ onde \bar{x} e s denotam a média e o desvio padrão das observações. Esta transformação faz com que as novas variáveis tenham média nula e variância unitária.

Outra forma de se transformar variáveis é tornar-se os desvios em relação ao menor valor e normalizá-los pela amplitude, ou seja,

$$Z_i = \frac{x_i - x_{(1)}}{x_{(n)} - x_{(1)}}, \quad i = 1, \dots, n \quad (5.1)$$

Onde $x_{(n)}$ e $x_{(1)}$ denotam o mínimo e o máximo da amostra, respetivamente.

Podemos ainda tomar a média como fator normalizador,

$$Z_i = \frac{x_i}{\bar{x}}, \quad i = 1, \dots, n \quad (5.2)$$

É importante ter cuidado no processo de estandardização, pois não deve ser tomado como solução ideal para todos os casos. Este processo reduz as diferenças entre os sujeitos anulando os agrupamentos naturais que possam existir nos dados.

5.3. Medidas de semelhança e medidas de dissemelhança

Na análise de Clusters, a escolha da medida de semelhança ou dissemelhança que melhor se adequa ao tipo de dados recolhidos representa um passo.

Segundo [35, 36] um coeficiente tem de ser considerado no contexto do estudo estatístico, incluindo a natureza dos dados e do tipo de análise pretendido.

Indicam alguns critérios para a escolha das medidas, no entanto, apesar da sua análise sobre o assunto concluíram que não é possível dar uma resposta definitiva.

As medidas de semelhança (ou dissemelhança) podem ser entre sujeitos ou entre variáveis, de acordo com o objetivo do estudo, respetivamente *Clusters* de sujeitos ou *Clusters* de variáveis.

5.3.3. Medidas de semelhança e medidas de dissemelhança entre sujeitos

São obtidas de uma matriz multivariada $X_{n \times p}$ resultante da observação de variáveis em sujeitos, e são escolhidas de acordo com o tipo de variáveis.

Os números d_{ij} (valor de uma medida de dissemelhança entre o sujeito i e sujeito j) ou s_{ij} (valor de uma medida de semelhança entre o sujeito i e sujeito j) são colocados numa matriz $n \times n$, conhecida por matriz de semelhança (ou dissemelhança).

A proximidade de dois sujeitos i e j é tanto maior quanto menor é a dissemelhança ou distância entre eles.

O estudo das relações de semelhança é inspirado em modelos geométricos, os sujeitos são representados por pontos no espaço. Deste modo as dissemelhanças observadas entre os sujeitos são visualizadas como a distância entre os respetivos pontos.

5.3.3.1. Dissemelhanças e distâncias – propriedades

Uma medida de dissemelhança entre um sujeito e um sujeito deverá satisfazer algumas propriedades:

- $d_{ij} \geq 0, \forall i, j = 1:n$;
- $d_{ii} = 0, \forall i, j = 1:n$; (Identidade)
- $d_{ij} = d_{ji}, \forall i, j = 1:n$; (Simetria)
- $d_{ij} \leq d_{ik} + d_{kj}, \forall i, j, k = 1:n$ (Desigualdade triangular).

No caso de as medidas de dissemelhança verificarem além das três primeiras condições a desigualdade triangular fala-se em distância.

5.3.3.2. Semelhanças – propriedades

- $0 \leq s_{ij} \leq 1, \forall i, j$

Quando $s_{ij} = 0$ os objetos não são semelhantes

Quando $s_{ij} = 1$ significa que a semelhança é máxima

- $s_{ij} = s_{ji} \quad \forall i, j = 1:n$ (Simetria)
- $s_{ij} = 1 \quad \forall i = 1:n$ (Identidade)

As medidas de semelhança (ou dissemelhança) dependem em primeiro lugar, do tipo de variáveis que caracterizam os sujeitos.

5.3.3.3. Medidas de dissemelhança e de semelhança para variáveis quantitativas

São várias as medidas que podem ser utilizadas como medidas de distância ou dissemelhança entre cada par de sujeitos. Assim para dois sujeitos i e j , para as variáveis $\vartheta = 1, 2, \dots, p$ [26, 37].

- **Distância Euclidiana**

$$d_{ij} = \left[\sum_{\vartheta=1}^p (X_{i\vartheta} - X_{j\vartheta})^2 \right]^{1/2}$$

- **Distância Euclidiana ao quadrado**

$$d_{ij} = \sum_{\vartheta=1}^p (X_{i\vartheta} - X_{j\vartheta})^2$$

- **Distância de Minkowski**

$$d_{ij} = \left[\sum_{\vartheta=1}^p (X_{i\vartheta} - X_{j\vartheta})^r \right]^{1/r}$$

para $r=1$, d_{ij} é o módulo da distância absoluta entre os sujeitos i e j relativamente às p -variáveis medidas (conhecida por Distância *city-block*);

para $r = 2$, tem-se a distância euclidiana habitual [26, 51].

- **Distância absoluta ou de Manhattan**

$$d_{ij} = \sum_{\vartheta=1}^p |X_{i\vartheta} - X_{j\vartheta}|$$

- **Distância de Chebishev**

$$d_{ij} = \max_{\vartheta} |X_{i\vartheta} - X_{j\vartheta}|$$

- **Distância de Mahalanobis**

$$d_{ij} = (X_i - X_j)' \Sigma^{-1} (X_i - X_j)$$

Apenas a distância de Mahalanobis, também chamada distância generalizada, utiliza a matriz de variância e covariância Σ fazendo implicitamente a estandardização das variáveis [26, 52].

- **Medida de Semelhança do Cosseno** - define-se no intervalo $[-1, 1]$

$$CoSIN(i, j) = \frac{\sum_{k=1}^p x_{ik} x_{jk}}{\sqrt{\sum_{k=1}^p x_{ik}^2 \sum_{k=1}^p x_{jk}^2}} = s_{ij}$$

$CoSIN(i, j)$, cosseno do ângulo formado pelas duas semirretas que unem a origem aos respectivos sujeitos, representados como pontos no espaço.

$s_{ij}=1$ significa que a semelhança é máxima

$s_{ij}=-1$ significa que a semelhança é mínima

Coefficiente de correlação de Pearson

Para dois sujeitos i e j , caracterizados por p atributos, dado por:

$$r_{ij} = \frac{\sum_{k=1}^p (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^p (x_{ik} - \bar{x}_i)^2 \sum_{k=1}^p (x_{jk} - \bar{x}_j)^2}}$$

O seu valor varia entre -1 e 1.

com

p = número total de variáveis

x_{ik} = valor da variável k para o sujeito i ($k = 1, \dots, p$)

x_{jk} = valor da variável k para o sujeito j

\bar{x}_i = média de todas as variáveis para o sujeito i

\bar{x}_j = média de todas as variáveis para o sujeito j

5.3.3.4. Medidas de dissemelhança e de semelhança para variáveis qualitativas

Na procura de elementos semelhantes é frequente o uso de critérios qualitativos, para medir o grau de semelhança entre os sujeitos, segundo variáveis qualitativas.

Estas medidas de semelhança (dissemelhança) geralmente têm valores no intervalo $[0, 1]$.

Se dois sujeitos i e j têm valores iguais para todas as variáveis, então têm coeficiente de semelhança, igual a um, $s_{ij}=1$.

Se dois sujeitos i e j diferem no máximo para todas as variáveis, então têm coeficiente de semelhança, igual a zero, $s_{ij}=0$.

Na literatura muitas são as propostas deste tipo de coeficientes.

5.3.3.4.1. Medidas de semelhança para variáveis nominais binárias. Medidas de associação

Indicadas para definir a semelhança entre os sujeitos de uma amostra multivariada caracterizados por variáveis qualitativas, em especial binárias⁴ (as medidas de distância métrica não são

⁴ Variáveis que apenas podem tomar dois diferentes valores (1/0, Presente/ Ausente, Sim/Não, Masculino/ Feminino, etc.)

aplicáveis). De entre os vários coeficientes de associação (s_{ij}) existentes, destaquem-se: coeficientes de emparelhamento simples (*simple matching*), coeficientes de Jaccard, coeficiente de Russel & Rao, coeficiente de Sorenson e coeficiente de Gower e Legendre [33,35,34,36].

Considere-se dois sujeitos i e j caracterizados por p variáveis nominais dicotômicas onde 1 e 0 significam, respetivamente, a presença e a ausência da característica em questão, as medidas de semelhança entre os dois sujeitos baseiam-se, em geral, nas seguintes quatro quantidades:

- a – o número de variáveis para os quais ambos os sujeitos toma o valor 1 (Presente);
- b - o número de variáveis para os quais o sujeito toma o valor 1 (Presente) e o sujeito toma o valor 0 (Ausente);
- c – o número de variáveis para os quais o sujeito toma o valor 0 (Ausente) e o sujeito toma o valor 1 (Presente);
- d – o número de variáveis em que ambos os sujeitos toma o valor 0 (Ausente).

O resumo do número de presenças e ausências das características das variáveis sob estudo para cada sujeito i e j pode ser representado na tabela de contingência:

Tabela 5. 1 Tabela de contingência[33]

		Sujeito j		Totais
		1	0	
Sujeito i	1	a	b	$a + b$
	0	c	d	$c + d$
Totais		a+b	c+d	$p = a + b + c + d$

De acordo com a informação da tabela os coeficientes são definidos como:

- **Coefficientes de emparelhamento simples (simple matching measures)**

$$s_{ij} = \frac{(a+b)}{(a+b+c+d)} \quad \text{ou} \quad d_{ij} = \frac{(b+c)}{(a+b+c+d)}$$

s_{ij} -Mede a semelhança entre cada dois sujeitos, varia entre 0 e 1. Representa a razão entre o número de características presentes e ausentes simultaneamente nos dois sujeitos e o número de características totais;

d_{ij} - Mede a distância entre os dois sujeitos, varia entre 0 e 1. Representa a razão entre o número de características presentes num sujeito, mas ausentes no outro, e o número total de características.

- **Coefficientes de Jaccard**– Medem a semelhança ou dissemelhança entre dois sujeitos. Não contemplam o número de características ausentes em ambos os sujeitos [26, 50, 53].

$$s_{ij} = \frac{a}{a+b+c}; \quad 0 \leq s_{ij} \leq 1 \quad \text{ou} \quad d_{ij} = \frac{b+c}{a+b+c}; \quad 0 \leq d_{ij} \leq 1$$

- **Coefficiente de Russel & Rao**– Dá-nos a perfeita semelhança ($s_{ij} = 1$) quando $b = c = d = 0$ e a máxima dissemelhança ($s_{ij} = 0$) quando $a = 0$ [52].

$$s_{ij} = \frac{a}{a+b+c+d}; \quad 0 \leq s_{ij} \leq 1$$

- **Coefficiente de Sorenson**– Valoriza a ocorrência simultânea da característica presente nos sujeitos.

$$s_{ij} = \frac{2a}{2a+b+c}; \quad 0 \leq s_{ij} \leq 1$$

- **Coefficiente de Gower e Legendre**– Toma a diferença entre concordâncias e discordâncias, relativamente ao número total de variáveis observadas. Ao contrário dos anteriores coeficientes, pode tomar valores negativos, situação que ocorre caso haja mais discordâncias do que concordâncias nos valores das variáveis para os sujeitos i e j . Toma valores entre -1 e 1. [35,36,26]

$$s_{ij} = \frac{(a+d)-(b+c)}{a+b+c+d}$$

5.3.3.4.2. Medidas de semelhança para variáveis nominais com mais de dois níveis

Quando a variável qualitativa nominal possui mais do que dois níveis, o artifício usual é a transformação em variáveis binárias através da criação de variáveis fictícias (*dummies*).

Supor o vetor de variáveis qualitativas nominais:

$$y' = (y_1, y_2, \dots, y_l)$$

onde a i -ésima componente assume l_i níveis, codificados de modo que

$$y_i = i, \text{ com } i = 1, 2, \dots, l_i$$

Supondo também que $\sum l_i = p$ Cada componente irá dar origem a l_i , variáveis binárias $x_{k(i)}$ tal que

$$x_{k(i)} = \begin{cases} 1 & \text{se } y_i = k \\ 0 & \text{em caso contrário} \end{cases}$$

Assim, o vetor y de dimensão l é transformado no vetor x de dimensão p , formado por componentes binárias. Esquemáticamente tem-se:

$$y_1 y' = (y_1, y_2, \dots, y_l) \rightarrow x' = \left(\underbrace{0, \dots, 1, \dots, 0}_{l_1}; \dots; \underbrace{0, \dots, 1, \dots, 0}_{l_l} \right)$$

Sem perda de generalidade, o vetor x será indicado p por coordenadas binárias x_i , isto é

$$x' = (x_1, x_2, \dots, x_p)$$

e tem-se a situação anterior.

Para corrigir o desequilíbrio causado pelo diferente número de níveis de cada variável, faz-se intervir no cálculo do coeficiente o número de níveis de cada variável. Supondo que há p variáveis, y_1, \dots, y_p com, l_1, \dots, l_p níveis respetivamente, então o coeficiente de semelhança s_{ij} será

$$s_{ij} = \frac{\sum_{k=1}^p l_n l_k I(y_{k(i)}, y_{k(j)})}{\sum_{k=1}^p l_n l_k}$$

Sendo I a função indicatória dos níveis dos sujeitos, i e j , na variável k , isto é,

$$I(y_{k(i)}, y_{k(j)}) = \begin{cases} 1 & \text{se } y_{k(i)} = y_{k(j)} \\ 0 & \text{se } y_{k(i)} \neq y_{k(j)} \end{cases}$$

$y_{k(i)}$ e $y_{k(j)}$ são, respetivamente, os níveis dos sujeitos i e j , na variável k .

5.3.3.4.3. Medidas de semelhança para variáveis ordinais

No caso de variáveis qualitativas do tipo ordinal, uma solução simples é considerá-las simplesmente qualitativas e aplicar qualquer um dos coeficientes definidos anteriormente. Este procedimento deixa de considerar a importante propriedade da ordem.

Podemos utilizar uma extensão do conceito de variáveis fictícias para este tipo de variáveis. Assim, a mesma estratégia usada anteriormente para transformar cada possível realização, numa variável binária, de acordo com a ocorrência do atributo, também pode ser usada nesta situação. Porém, deve ser considerada a questão da ordem.

5.3.3.5. Coeficiente de semelhança para variáveis de diferentes tipos

É frequente a presença de diferentes tipos de variáveis (quantitativas e qualitativas) na procura de sujeitos semelhantes.

A seguir, são apresentadas algumas estratégias para aplicar a uma matriz de dados que contém diferentes tipos de variáveis.

Coefficientes combinados de semelhança

Determina-se os coeficientes de semelhança de mesmo sentido (semelhança ou dissemelhança), s_n, s_o, s_q , e para cada grupo de variáveis (nominais, ordinais e quantitativas), depois constrói-se um único coeficiente ponderado.

Para dois sujeitos A e B, esquematicamente, tem-se:

$$S_{AB} = W_1.S_{nAB} + W_2.S_{oAB} + W_3.S_{qAB}$$

Onde W_p , $p = 1,2,3$ os são os pesos associados. É comum ponderar pelo número de variáveis envolvidas.

Proposta de Gower

Gower [50] propõe uma forma mais elaborada do coeficiente de semelhança combinado.

O coeficiente de entre os sujeitos A e B, segundo as p variáveis, de qualquer tipo, passa a ser:

$$S_{AB} = \sum I_{iAB} S_{iAB} / \sum I_{iAB}$$

Para cada variável x_j é definido um coeficiente de semelhança S_i com valores entre 0 e 1. A variável I_i , assume o valor 1 quando a comparação dos sujeitos é possível segundo o critério i , e assume o valor 0 em caso contrário, isto é se o valor da variável é omissa em pelo menos um dos sujeitos A e B. Este coeficiente será indefinido quando todos $I_{iAB} = 0$, ou seja, a comparação dos dois sujeitos não é válida segundo nenhum critério.

Este coeficiente torna-se idêntico ao coeficiente de Jaccard quando as variáveis são todas binárias.

Proposta de Romesburg

Romesburg [49] sugere esquecer a natureza das variáveis, tratar todas como variáveis quantitativas (todas são codificadas com números) e aplicar a distância euclidiana. A grande desvantagem está na interpretação dos valores dos coeficientes de semelhança, pois estes dependem da codificação das variáveis.

5.3.3.6. Conversão das semelhanças em dissemelhanças

Quando o ponto de partida é uma matriz de semelhanças de um conjunto de sujeitos, pode-se trabalhar diretamente sobre estas medidas de semelhança ou, alternativamente, converter as medidas de semelhança em dissemelhança. É possível estabelecer uma relação entre as semelhanças e dissemelhanças dos sujeitos:

$$d_{ij} = 1 - S_{ij}$$

$$d_{ij} = 1 - S_{ij}^2$$

$$d_{ij} = \sqrt{1 - S_{ij}^2}$$

$$d_{ij} = \sqrt{1 - S_{ij}}$$

A primeira destas regras de conversão, no caso de se utilizar o Coeficiente de emparelhamento simples como medida de semelhança, equivale a tomar o coeficiente de dissemelhança.

5.3.4. Medidas de semelhança entre variáveis

É também possível agrupar variáveis através duma Análise de Clusters. Neste caso, o ponto de partida será uma matriz de semelhanças (ou dissemelhanças) entre variáveis.

As variáveis tomam o lugar dos sujeitos e podemos aplicar as medidas de (dis) semelhança utilizadas na análise de sujeitos.

De um modo geral, o agrupamento de variáveis é baseado em medidas de correlação ou associação.

Independentemente dos critérios de semelhança e/ ou dissemelhança adotados entre variáveis, o procedimento de classificação que se segue será análogo aos procedimentos de classificação de sujeitos anteriormente apresentados.

5.3.4.1. Medidas de semelhança entre variáveis quantitativas

Seja o comportamento das variáveis X e Y , representados pelos vetores $X = (x_1, \dots, x_n)'$ e $Y = (y_1, \dots, y_n)'$. O valor observado da variável no i -ésimo objeto é representado pelo i -ésimo componente de cada vetor, $i = 1, \dots, n$.

Coefficiente de correlação de Pearson – este coeficiente que varia entre $[-1,1]$ mede a intensidade e a direção da associação de tipo linear entre duas variáveis quantitativas e define-se como

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Poderemos também usar a Medida de Semelhança do Cosseno com variáveis quantitativas.

Medida de Semelhança do Cosseno

$$s_{XY} = \text{CoSIN}(X, Y) = \frac{\sum_{i=1}^n X_i Y_i}{\sqrt{\sum_{i=1}^n X_i^2 \sum_{i=1}^p Y_i^2}}$$

$\text{CoSIN}(X, Y)$, cosseno do ângulo formado pelos dois vetores. À medida que o vetor x “aproxima” de y , o cosseno do ângulo cresce [26, 27, 51].

5.3.4.2. Medidas de semelhança entre variáveis nominais binárias

Retomemos à Tabela 1, nesta situação, em que i e j representam a i -ésima e a j -ésima variáveis e tomando os valores 1 e 0 para representar as duas categorias das variáveis, os coeficientes de correlação de Pearson e de medida de semelhança do cosseno são dados por [26, 27, 51]:

$$r_{ij} = \frac{ad - bc}{\sqrt{(a + b)(c + d)(a + c)(b + d)}}$$

Medida de Semelhança do Cosseno

$$\text{CoSIN}(i, j) = \frac{a}{\sqrt{(a + b)(a + c)}}$$

5.3.4.3. Medidas de semelhança entre variáveis nominais com mais de dois níveis

Consideremos as variáveis X e Y com as categorias $1, \dots, p$ e $1, \dots, q$ respectivamente. Suponhamos uma tabela de contingência onde n representa o número total de observações, n_{ij} representa a frequência absoluta do par (X_i, Y_j) , e n_i e n_j as frequências marginais de X e Y , respectivamente. Analogamente, definamos as frequências relativas, f_{ij} , f_i e f_j .

Tabela 5.2: Tabela de Contingência [61].

X	Y						Total
	Y_1	Y_2	...	Y_j	...	Y_q	
X_1	n_{11}	n_{12}	...	n_{1j}	...	n_{1q}	$n_{1.}$
X_2	n_{21}	n_{22}	...	n_{2j}	...	n_{2q}	$n_{2.}$
...

X_i	n_{i1}	n_{i2}	...	n_{ij}		n_{iq}	$n_{i.}$
...			
X_p	n_{p1}	n_{p2}	...	n_{pj}		n_{pq}	$n_{p.}$
Total	$n_{.1}$	$n_{.2}$...	$n_{.j}$		$n_{.q}$	n

Apresenta-se de seguida a medida mais usada *Qui-quadrado* e as variações desta:

Qui-quadrado x^2 de Pearson

$$x^2 = n \sum_{i=1}^p \sum_{j=1}^q (f_{ij} - f_i.f_j)^2 / f_i.f_j$$

Coefficiente de contingência quadrática média

$$\phi^2 = \frac{x^2}{n}$$

Coefficiente de contingência de Pearson

$$P = \left(\frac{\phi^2}{1 + \phi^2} \right)^{1/2}$$

Coefficiente de Tschuprow

$$T = \left[\frac{\phi^2}{(p-1)(q-1)} \right]^{1/2}$$

Coefficiente V de Cramer

$$V = \sqrt{\frac{\phi^2}{\min(p-1, q-1)}}$$

5.3.4.4 Medida de semelhança entre variáveis ordinais

Como já foi referido o coeficiente de Spearman é uma das medidas de associação entre variáveis ordinais mais usada. Toma valores no intervalo $[-1,1]$ e pode obter-se usando a fórmula dos coeficientes de correlação de Pearson, substituindo os valores das observações X e Y pelas respectivas ordens r_1 e r_2 :

$$r_s = \frac{\sum_{i=1}^n (r_{1i} - \bar{r}_1)(r_{2i} - \bar{r}_2)}{\sqrt{\sum_{i=1}^n (r_{1i} - \bar{r}_1)^2 (r_{2i} - \bar{r}_2)^2}}$$

5.4. Métodos hierárquicos

Permitem a obtenção de *Clusters* quer para indivíduos quer para variáveis. Dividem-se em aglomerativos e divisivos; os mais divulgados e mais utilizados são os hierárquicos aglomerativos. No método aglomerativo, o agrupamento em classes procede por etapas, em geral determinando-se a partir de n subgrupos (de um indivíduo cada) sucessivas fusões de subgrupos considerados mais “semelhantes”, até se encontrar apenas um grupo ou *Cluster* que incluirá a totalidade de n indivíduos. No método divisivo o processo é inverso.

O ponto de partida para os métodos hierárquicos é, em geral, uma matriz $n \times n$ cujo elemento genérico (i, j) é uma medida de semelhança (ou dissemelhança) entre o indivíduo i e o indivíduo j [26, 55].

O método aglomerativo desenrola-se de acordo como seguinte algoritmo:

1. Começar com n *Clusters* (um para cada indivíduo ou variável) e calcular a matriz de dissemelhança (ou de semelhança) $D_n \times n$;
2. Encontrar na matriz os pares *Clusters* (indivíduos ou variáveis) mais semelhantes de acordo com uma medida de distância escolhida;
3. Com os *Cluster* i e j encontrados formar um *Cluster* maior, *Cluster* ij e recalculer a distância deste *Cluster* para os restantes *Cluster* originais;
4. Repetir os passos 2 e 3 até sobrar um único *Cluster*.

5.4.1. Métodos de (des) agregação - Características

As medidas de distância entre os indivíduos não representam a única opção a fazer numa Análise de *Clusters*. É necessário também escolher o método de (des) agregação dos indivíduos.

Os métodos hierárquicos de agrupamento diferem no modo como calculam as distâncias entre grupos e os restantes (grupos ou indivíduos). Os métodos mais utilizados são os seguintes [26, 27,55]:

- **Menor distância** (*Single linkage ou Nearest neighbor*) – Consiste em considerar que a distância entre dois grupos é a menor distância entre um elemento dum grupo e um elemento do outro grupo. Seja k a distância e dois grupos (i, j) :

$$d_{(i,j)k} = \min\{d_{ik}; d_{jk}\}$$

Este método tende a produzir classes com indivíduos que podem estar muito distantes entre si, mas pertencendo a uma mesma classe. Este facto resulta da existência de um elemento numa classe “próximo” de um único elemento de outra classe para que estas sejam atraídas.

- **Maior distância** (*Complete linkage ou farthest-neighbor*) – Consiste em considerar que a distância entre dois grupos é a maior distância entre um elemento dum grupo e um elemento de outro grupo. Seja k a distância e dois grupos (i, j) :

$$d_{(i,j)k} = \max\{d_{ik}; d_{jk}\}$$

Caso os sujeitos sejam representáveis em R^p , este método tem tendência a produzir classes onde não há grandes diferenças nas distâncias entre pares de elementos mais distantes, ao longo de várias direções - *classes “esféricas”*.

- **Distância Média entre Clusters** (*Average linkage between groups*) – Consiste em considerar que a distância entre dois grupos é a média de todas as distâncias entre pares de elementos (um de cada grupo). Seja k a distância e dois grupos (i, j) :

$$d_{(i,j)k} = \text{média}\{d_{ik}; d_{jk}\}$$

Assim como o método da menor distância, este método também tem tendência a produzir classes “esféricas”.

- **Distância Média dentro dos Clusters** (*Average linkage within groups*) – Semelhante à “Distância média entre *clusters*” mas neste método os clusters são unidos de modo a que a soma de quadrados dos erros seja mínima.
- **Distância Mediana** (*Median linkage*) – Após formado o primeiro *Cluster*, a distância deste aos restantes é a mediana das distâncias de cada um dos elementos constituintes deste *Cluster* a cada um dos restantes indivíduos ou variáveis. Seja k a distância e dois grupos (i, j) :

$$d_{(i,j)k} = \text{mediana}\{d_{ik}; d_{jk}\}$$

- **Método Centróide** – Toma-se a distância entre dois grupos como sendo a distância entre os centros de gravidade ou outros pontos considerados “representativos” (centróides) dos grupos. O método do centróide calcula a distância entre dois grupos como a diferença entre as suas médias, para todas as variáveis.
Este método também tem tendência a produzir classes “esféricas”.
- **Método Ward** (*Método da Inércia Mínima*) – Neste método os *Clusters* são formados de modo a minimizar a soma dos quadrados dos desvios das observações individuais relativamente às médias dos grupos em que são classificadas.
O método Ward tem tendência a produzir classes com um número aproximado igual de sujeitos.

A escolha do método de (des) agregação condicionará a classificação obtida. A escolha do método deve ser justificada com base na natureza dos dados e no objetivo da análise.

O processo de agrupamento é possível ser visualizado através de uma representação gráfica denominada *dendrograma*. O primeiro mostra todas as fases do processo do agrupamento desde a separação total dos indivíduos até à sua inclusão num grupo apenas. A posição na escala horizontal indica a distância a que os *Clusters* são agrupados. Um corte no dendrograma a qualquer nível de aglomeração produz uma classificação em K grupos ($1 \leq K \leq n$).

Numa classificação, segundo uma dada escolha de distâncias entre indivíduos e classes, a representação em dendrograma não é única, uma vez que a ordem dos indivíduos é arbitrária. Reordenações dos indivíduos podem produzir dendrogramas de aspeto diferente, mas a informação neles contida é idêntica.

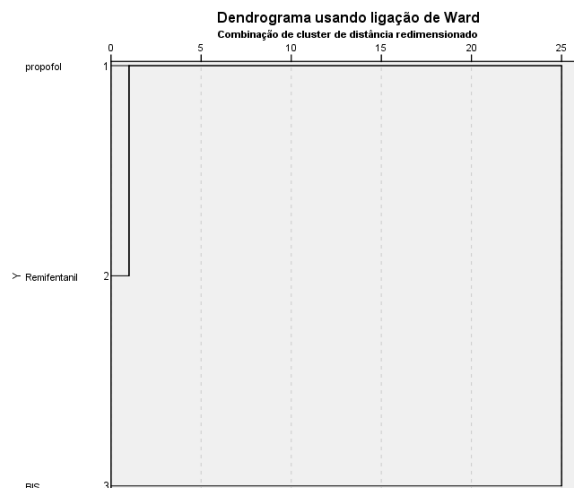


Figura 5.1: Dendrograma Fonte: elaborado pelo autor

A partir do dendrograma anterior, fazendo um corte a uma distância de aproximadamente 3 é possível identificar a existência de dois grupos: (2, 5, 3, 4) e (1, 6, 7).

5.5. Escolha do número de *Clusters*

O objetivo da análise de *Clusters* é formar grupos homogêneos. Assim, coloca-se o problema da escolha do número apropriado de *Clusters*. Existem várias técnicas para se determinar o número adequado de *Clusters* [26,34,51].

5.5.1. Análise do dendrograma

O procedimento básico consiste no exame do dendrograma procurando grandes alterações (saltos) na distância para as sucessivas fusões.

Qualquer método produzirá sempre uma classificação em qualquer número de classes, conforme o nível em que decidamos cortar o dendrograma. Por este facto a análise de *Clusters* produz

classificações mesmo onde elas possam não fazer sentido. Assim, é importante verificar a robustez das classificações obtidas.

Uma boa classificação deverá corresponder a um agrupamento que resulte de cortar o dendrograma numa zona onde as separações entre classes correspondam a grandes distâncias (barras de junção de classes relativamente compridas) - heterogeneidade entre classes. A homogeneidade interna das classes, será tanto maior quanto mais próximo dos indivíduos se fizer o corte.

5.5.2. Coeficiente de fusão

É o valor numérico (distância ou semelhança) para o qual os vários indivíduos se unem para formar um grupo. A comparação gráfica do número de *Clusters* com o valor do coeficiente de fusão permite sugerir a escolha do número de *Clusters*. Quando a divisão de um novo grupo não introduz alterações no coeficiente de fusão, poderá tornar-se essa partição como sendo ótima [34]. A escolha ótima coincidirá com uma marcada horizontalidade da curva.

5.6. Métodos não hierárquicos

Os métodos não hierárquicos são válidos apenas para a obtenção de *Clusters* de indivíduos (e não de variáveis). Estes métodos são enquadrados como partitivos (dividem os n dados existentes em K partições). Fixa-se, à partida, o número de partições que se pretende constituir e (regra geral) faz-se a classificação dos n indivíduos em K *Clusters*, de modo a otimizar algum critério de homogeneidade interna e heterogeneidade externa. Existem vários métodos não-hierárquicos, o desempenho destes métodos depende da primeira agregação dos indivíduos em *Clusters*, e do modo como as novas distâncias entre os centroides⁵ dos *Clusters* e os indivíduos são calculada.

Um dos métodos partitivos mais frequente nos softwares estatísticos é o *K-means* que se desenrola nos seguintes passos [37].

1. Partição inicial dos sujeitos em K *Clusters* definidos, à partida, pelo investigador;

⁵Centroides são os valores médios contidos em cada uma das variáveis do *cluster*.

2. Cálculo dos centróides para cada um dos K *Clusters* e cálculo da distância euclidiana dos centróides a cada sujeito na base de dados;
3. Agrupar os sujeitos aos *Clusters* de cujos centróides se encontram mais próximos, e voltar ao passo 2 até que não ocorra variação significativa na distância mínima de cada sujeito da base de dados a cada um dos centróides dos K *Clusters* (ou até que o número máximo de interações ou o critério de convergência definido pelo analista seja alcançado).

Cada indivíduo é transferido para o *Cluster* que apresenta uma menor distância (por exemplo, distância euclidiana) entre o indivíduo e o centróide do *Cluster*. Assim, é necessário conhecer os centróides de cada *Cluster* ou calculá-los a partir dos dados originais.

5.7. Outros métodos

5.7.1. *TwoStep Cluster* (Análise de *Clusters* em duas fases)

Os tradicionais métodos de agrupamento são eficientes e rigorosos quando aplicados a pequenos conjuntos de dados. O mesmo não se verifica no caso de conjuntos de dados muito grandes. Para aplicar os métodos tradicionais, é necessário previamente reduzir a dimensão da base de dados, ou seja, o agrupamento é realizado em dois passos como BIRCH [38]. O método *TwoStep Cluster* usa este procedimento, permitindo dar resposta a conjuntos de dados de enorme dimensão e a utilização de variáveis contínuas, categóricas ou os dois tipos de variável em simultâneo.

Passo 1: formação de uma série de preclusters. O objetivo deste passo é reduzir o tamanho da matriz das distâncias entre todos os pares de casos possíveis. Nesta primeira etapa, os dados são percorridos um a um e o algoritmo decide se um determinado indivíduo deve migrar para um precluster previamente formado ou iniciar um novo *preCluster*. No fim deste procedimento e todos os indivíduos pertencentes ao mesmo precluster são tratados como uma só entidade. Assim, a matriz de distâncias é menor, pois o seu tamanho passa a depender do número de *preClusters*.

Passo 2: agrupamento dos preclusters. No segundo passo, ocorre o agrupamento hierárquico (dos *preCluster* formados na etapa anterior) de acordo com o número de *Clusters* pretendido.

Medidas de distância: se só existirem variáveis contínuas, é possível usar a distância euclidiana entre o centro de dois *Clusters*. Quando existem misturadas variáveis contínuas e categóricas, é utilizada a função log-verossimilhança, a distância entre dois *Clusters* é expressa pelo decréscimo da função log-verossimilhança. Neste caso, o algoritmo fornece melhores resultados quando se verifica a normalidade das variáveis contínuas, e a distribuição multinomial no caso das variáveis categóricas. Estes pressupostos são difíceis de encontrar em dados reais, no entanto, o algoritmo encontra uma solução razoável mesmo quando os pressupostos são quebrados.

Principais vantagens: utilização de variáveis contínuas e categóricas simultaneamente; agrupamento em duas etapas, aumentando a eficiência do método; o próprio algoritmo encontra um número ótimo de *Clusters*, sendo também possível especificar o número de *Clusters* desejado; de fácil interpretação, são disponibilizadas informações sobre a importância de cada variável na formação de cada *Cluster* e uma medida de significância estatística (Qui-quadrado para variáveis categóricas e t – *test* para variáveis contínuas), permitindo a confirmação dos perfis definidos.

5.7.2. Técnicas de densidade

Com aplicação quando se espera observar grupos naturais. Os agrupamentos são formados através da procura de regiões que contenham uma concentração relativamente densa de pontos (casos). Regiões com muitos pontos próximos no espaço, separadas por áreas com poucos pontos (representando ruídos), sugerem *Clusters*. Em geral, usa-se o método da ligação simples para a obtenção dos *Clusters*. Começa-se por escolher um raio r e o número de pontos P . Uma região densa – pontos de densidade - é uma região onde uma vizinhança (círculo de raio r) de cada ponto contém pelo menos P pontos. Os primeiros *Clusters* são definidos pelos pontos de densidade. Um ponto cuja distância a todos os pontos de densidade seja superior a r , forma o seu próprio *Cluster*. O algoritmo repete-se e só para quando não houver mais possibilidades de se juntar *Clusters* da etapa precedente [26, 27, 52].

5.7.3. Agrupamentos fuzzy (Difusos)

O agrupamento *fuzzy* é uma generalização dos métodos partitivos, permitindo que haja sobreposição dos grupos (*fuzzy Clusters*). É possível observar o grau de associação de cada elemento em cada grupo, que geralmente se verifica em domínios de dados reais, onde um elemento pertence a diferentes grupos, com diferentes graus de associação – vantagem relativamente a outros métodos por partição. No entanto, apresenta a desvantagem do número de coeficientes de associação crescer rapidamente com o aumento do número de elementos e de grupos. Trata-se de uma técnica válida, pois associa graus de incerteza aos elementos nos grupos, situação que se aproxima das características reais dos dados [39].

Um método de agrupamento *fuzzy* frequentemente utilizado é o *fuzzy c-means*. Este método iterativo inicia-se com c valores arbitrários, com base nos quais associa cada elemento ao valor com menor distância, formando c grupos. Depois, determina-se o centro de cada *Cluster* formado e os elementos são reagrupados ao centro mais próximo. Este procedimento termina quando as diferenças entre os centros do passo atual e do anterior sejam mínimas.

5.8. Métodos hierárquicos / métodos não hierárquicos

Enquanto a aplicação dos métodos hierárquicos requer o cálculo de uma matriz de dissimilaridades, os métodos não-hierárquicos aplicam-se diretamente sobre os dados originais, permitindo a sua aplicação a matrizes de dados muito grandes.

Os métodos não-hierárquicos permitem reagrupar os indivíduos num *Cluster* diferente daquele em que foram inicialmente incluídos. Nos métodos hierárquicos, os indivíduos que sejam incluídos num mesmo *Cluster* em qualquer etapa do processo não poderão mais ser separados em etapas posteriores.

A necessária definição, à partida, do número de *Clusters*, sem conhecimento da estrutura dos dados, pode representar uma desvantagem nos métodos não-hierárquicos.

Num problema de análise de *Clusters*, a validade das soluções encontradas aumenta, se o processo de análise começar com um método hierárquico aglomerativo para determinar o número de grupos, e proceder com o *k-means* para refinar (otimizar) a partição encontrada [26, 55].

5.9. Escolha da técnica a utilizar

Qualquer um dos métodos de análise de *Clusters* impõe um certo grau de estrutura nos dados e, para o investigador assegurar que o resultado obtido não é um artefacto da técnica utilizada, é aconselhável usar diferentes critérios de agrupamento e escolher a estrutura resultante da maior parte deles.

Podemos utilizar outra estratégia para averiguar a estabilidade do agrupamento. Esta estratégia consiste em formar, ao acaso, dois subconjuntos do conjunto de observações e aplicar em cada um deles o mesmo critério. A alocação dos sujeitos nas subamostras e na amostra total será semelhante se o agrupamento for estável [26].

5.10. Análise Discriminante

A análise discriminante tem como principal objetivo construir um modelo preditivo através da determinação de uma regra de classificação que permita afetar um elemento a um de dois ou mais grupos especificados *a priori* com base num conjunto de variáveis dependentes [42,43]. Com um objetivo semelhante à regressão logística, a análise discriminante, que se insere na análise de variância múltipla (MANOVA), pode ser vista como um método descritivo e/ou preditivo para um conjunto de variáveis dependentes definidas por uma variável independente categórica (para 2 ou mais grupos).

Na análise discriminante pretende-se assim explorar os dados de forma a maximizar a diferença entre os K grupos existentes. Para isso são criadas funções preditivas, combinações lineares das variáveis preditivas, que discriminam os K grupos, maximizando a diferença entre eles [44], ou seja, dados n elementos, K grupos e p variáveis preditivas (X_1, X_2, \dots, X_p) , pretende-se formar uma combinação linear da seguinte forma:

$$y_j = \sum_{i=1}^p \vartheta_{ji} X_i \quad (5.3)$$

Seja Y o vector das funções lineares y_j e ϑ , a matriz dos coeficientes lineares ϑ_{ji} e X o vector das variáveis preditivas, temos que:

$$Y = \vartheta X \quad (5.4)$$

Para cada função discriminante y_j queremos determinar os coeficientes ϑ_{ji} , componentes do vector ϑ_j , que maximizam as diferenças entre os grupos. Existem duas aproximações:

- o critério discriminante, em que se pretende maximizar a razão entre a variabilidade inter-grupos e a variabilidade intra-grupos;
- as correlações canónicas, em que se pretende maximizar a correlação entre a combinação linear das p variáveis preditivas e a combinação linear das q variáveis *dummy*, ($q = K - 1$).

As funções discriminantes são obtidas após maximização do critério em questão (critério discriminante ou correlação canónica). As soluções obtidas nas duas aproximações são idênticas uma vez que estas se encontram relacionadas [44].

5.10.1. Critério Discriminante

Nesta primeira aproximação, pretende-se maximizar λ , razão entre a soma dos quadrados inter-grupos (SS_b) e a soma dos quadrados intra-grupos (SS_w) das respectivas matrizes de covariância, dado por:

$$\lambda \equiv \frac{SS_b(Y)}{SS_w(Y)} \quad (5.5)$$

A dedução matemática da maximização do critério discriminante (5.5) em [26,44].

Seja W a matriz *SSCP* intra-grupos dada por $W = \sum_{k=1}^K S_k(x)$ e B a matriz *SSCP* inter-grupos dada por $B = (\bar{x} - \bar{\bar{x}})'(\bar{x} - \bar{\bar{x}})$, em que \bar{x} é a matriz compostapelas médias das p variáveis para cada grupo K , dada por (5.6), e $\bar{\bar{x}}$ é a matriz composta pelasmédias das p variáveis dada por (5.7).

$$\bar{X} = \begin{bmatrix} \bar{X}_{11} & \dots & \bar{X}_{p1} \\ \dots & \dots & \dots \\ \bar{X}_{1K} & \dots & \bar{X}_{pK} \end{bmatrix} \quad (5.6)$$

$$\bar{\bar{X}} = \begin{bmatrix} \bar{\bar{X}}_1 & \dots & \bar{\bar{X}}_p \\ \dots & \dots & \dots \\ \bar{\bar{X}}_1 & \dots & \bar{\bar{X}}_p \end{bmatrix} \quad (5.7)$$

Uma vez que $SS_w(Y) = v'wv$ e $SS_b(Y) = v'Bv$, o critério discriminante passa a ser dado por:

$$\lambda = \frac{v'Bv}{v'Wv} \quad (5.8)$$

Querendo determinar os coeficientes das combinações lineares, componentes do vector v , que maximize λ , temos que resolver a seguinte equação $\left(\frac{\partial \lambda}{\partial v} = 0\right)$ ou seja:

$$\frac{\partial \lambda}{\partial v} \left(\frac{v'Bv}{v'Wv} \right) = \quad (5.9)$$

Derivando (5.9) [26,44,59] obtemos a seguinte equação:

$$(W^{-1}B - \lambda I)v = 0 \quad (5.10)$$

Em que equação característica da matriz $W^{-1}B$ é dada por:

$$|W^{-1}B - \lambda I| = 0 \tag{5.11}$$

Determinando os valores próprios λ , obtemos os vetores próprios v_i e, por conseguinte, obtemos os coeficientes das funções discriminantes.

5.10.2. Correlações Canónicas

Nas correlações canónicas queremos maximizar a correlação entre as p variáveis preditivas e as q variáveis *dummy* ($q = K - 1$, em que K é o número de grupos). Estas variáveis *dummy* (Y_j) são utilizadas para atribuir os casos nos K grupos existentes, ou seja, para cada k entre 1 e $K - 1$, todos os elementos do k -ésimo grupo têm valor 1 na variável Y_k e valor 0 nas restantes Y_j ($j \neq k$), enquanto todos os elementos do último grupo K , têm o valor 0 nas $K - 1$ variáveis Y_j , ou seja, a matriz Y das variáveis *dummy* é dada por:

$$Y = \begin{matrix} & Y_1 & Y_2 & Y_3 \dots & Y_{K-1} \\ \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & 0 \\ 0 & 0 & 0 & \dots & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} & \begin{matrix} \text{todos os casos pertencem ao grupo 1} \\ \text{todos os casos pertencem ao grupo 2} \\ \text{todos os casos pertencem ao grupo 3} \\ \dots \\ \text{todos os casos pertencem ao grupo } K - 1 \\ \text{todos os casos pertencem ao grupo } K \end{matrix} \end{matrix}$$

Temos, desta forma, dois conjuntos de variáveis Y_i e Y_j ($i = 1, p$ $j = 1, q$) a partir dos quais se constroem as seguintes combinações lineares:

$$Z_l = \sum_{i=1}^p \mu_i X_i \text{ (função discriminante)} \tag{5.12}$$

$$W_m = \sum_{j=1}^q \vartheta_j Y_j \text{ (função de critério)} \tag{5.13}$$

Pretende-se, portanto, determinar os dois conjuntos de coeficientes $u' = [\mu_1, \mu_2, \dots, \mu_p]$ e $v' = [v_1, v_2, \dots, v_p]$ tal que a correlação entre as funções Z_l e W_m seja máxima.

A correlação entre estas duas funções é dada por:

$$r_{z_l w_m} = \frac{\sum z_l w_m}{\sqrt{(\sum z_l^2)(\sum w_m^2)}} \quad (5.14)$$

A dedução matemática da maximização do coeficiente de correlação (5.14) encontra-se em [26,44].

Temos que $\sum Z_l^2 = u' S_{xx} u$, $\sum w_m^2 = v' S_{yy} v$ e $\sum Z_l W_m = u' S_{xy} v$. Considerando $u' S_{xx} u = v' S_{yy} v = 1$ [44] o coeficiente de correlação reduz-se a:

$$r_{z_l w_m} = u' S_{xy} v \quad (5.15)$$

A maximização de $r_{z_l w_m}$ em função de u e v , dadas as considerações feitas, reduz-se à seguinte equação [44]:

$$(S_{xx}^{-1} S_{xy} S_{yy}^{-1} S_{xy} - \mu^2 I) \mu = 0 \quad (5.16)$$

Em que S_{xx} é a matriz *SSCP* das variáveis preditivas dada por X , S_{yy} é a matriz *SSCP* das variáveis *dummy* dada por y (equação 5.16), S_{xy} e S_{yx} são as matrizes da soma dos produtos entre X e Y e Y e X , respetivamente.

Os valores próprios e os vetores próprios da equação característica do produto das matrizes *SSCP* são obtidos através da resolução da seguinte equação:

$$|S_{xx}^{-1} S_{xy} S_{yy}^{-1} S_{yx} - \mu^2 I| = 0$$

O maior valor próprio, $(\mu^2)_1$, corresponde o valor máximo de $r_{z_l w_m}$ (primeiro coeficiente canónico) e os elementos do vetor próprio associado, μ_1 , dão os coeficientes da primeira combinação linear das variáveis X_i , isto é, da primeira função discriminante [27, 43,44,].

5.10.3 Testes de Significância

Uma vez determinadas as funções discriminantes, é preciso calcular a significância das mesmas em discriminar os K grupos. Para testar a existência, ou não, de diferenças estatisticamente significativas entre grupos, recorre-se aos testes estatísticos de distribuição F .

Os testes estatísticos, que seguem uma distribuição F , testam a hipótese nula de que as variâncias dos grupos não diferem. O cálculo de F , na sua forma multivariada (para mais do que 2 grupos), baseia-se na matriz $SSCP$, ou seja, baseia-se não só na soma dos quadrados inter- e intra-grupos (variância), como acontece nos testes ANOVA, mas também na soma dos produtos cruzados (covariâncias).

Passa-se a descrever resumidamente alguns desses testes cuja dedução matemática pode ser encontrada em [33, 37, 60].

Seja W a matriz $SSCP$ intra-grupo, B a matriz $SSCP$ inter-grupo e T a matriz $SSCP$ total.

Teste Lambda de Wilks

Este teste baseia-se nos valores próprios da matriz $w^{-1}T$, α_i , isto é:

$$\Lambda = \frac{1}{\prod_{i=1}^p (1 + \alpha_i)} \quad (5.17)$$

Quanto menor for Λ , menor será a proporção de variância intra-grupos e maior será a diferença entre grupos [33, 37, 60].

Teste Hotelling-Lawley Trace

Este teste baseia-se nos valores próprios da matriz $w^{-1}B$, α_i , isto é:

$$\tau = \sum_{i=1}^p \alpha_i \quad (5.18)$$

Quanto maior for τ , maior será a diferença entre os grupos [33, 37, 60].

Teste Pillai-Bartlett Trace

Este teste baseia-se nos valores próprios da matriz $T^{-1}B$, α_i , isto é:

$$V = \sum_{i=1}^p \left(\frac{\alpha_i}{1 + \alpha_i} \right) \quad (5.19)$$

Quanto maior for V , maior será a diferença entre os grupos [33, 37, 60].

Teste Roy's Largest Root

Este teste baseia-se no primeiro valor próprio da matriz $T^{-1}B$, α_1 , isto é:

$$\theta = \frac{\alpha_1}{1 + \alpha_1} \quad (5.20)$$

Quanto maior for θ , maior será a diferença entre os grupos, explicada pela primeira função discriminante [33, 37, 60].

De uma forma geral, para uma amostra grande, todos eles produzem resultados semelhantes; no entanto, não existe nenhum consenso quanto ao teste mais adequado, pois estudos como os realizados por [57] e por [56] encontram resultados contraditórios, relativos à classificação dos 4 testes [44] O primeiro teste, Lambda de Wilks, é no entanto o teste mais usado por ser o mais comum e o mais tradicional quando temos mais do que 2 grupos. O segundo teste, de Hotelling-Lawley Trace, é usado como caso particular do teste anterior para 2 grupos. O teste de Pillai-Bartlett Trace tem demonstrado ser um teste mais robusto, enquanto o teste de Roy's Largest Root, aparece como sendo o teste menos robusto e mais sensível às violações de normalidade multivariada dos dados.

5.11. Identificação de grupos homogêneos

Com o intuito de explorar os dados e de confirmar a existência de grupos homogêneos na amostra, foram utilizadas as técnicas de agrupamento hierárquico. Consideramos o método aglomerativo com minimização da variância entre grupos (método de Ward), e as seguintes variáveis estandardizadas (z-scores):

- Remifentanil Ce ([ng/ml]) - concentração de fármaco opióide;
- Propofol Ce ($\mu\text{g/ml}$) - concentração de fármaco hipnótico;
- BIS- sinal cerebral;

Uma vez que os métodos aglomerativos iniciam com tantos *Clusters* quantos casos, os casos/ *Clusters* mais similares são agrupados nas primeiras iterações e os casos/ *Clusters* mais dissimilares são agrupados nas últimas. Em cada iteração é calculado o coeficiente de aglomeração que reflete a distância entre os casos/*Clusters* a serem agrupados, ou seja, quanto mais dissimilares forem os casos/ *Clusters* a agrupar, maior será o valor do coeficiente de aglomeração, sendo que para casos/ *Clusters* similares este coeficiente tenderá para zero.

Podemos, portanto, da análise dos coeficientes de aglomeração em função das iterações, identificar a região a partir da qual o coeficiente começa a agrupar *Clusters* cada vez mais dissimilares, isto é, a região a partir da qual o coeficiente de aglomeração começa a aumentar de forma significativa.

A região de interesse, para a amostra em estudo, situa-se entre as iterações 4208 e 4213 em que, nestas 5 iterações, o coeficiente de aglomeração aumenta 1035319,827 unidades, ou seja, passa 40390,258 para 1075710,085.

Analisando cada um dos passos de aglomeração temos que:

- da iteração 4208 para a iteração 4209 (o que corresponde a aglomeração de 5 *Clusters* para 4) o coeficiente aumenta 12050,176 unidades (de 40390,258 para 52440,434);
- da iteração 4210 para a iteração 4211 (o que corresponde a aglomeração de 4 *Clusters* para 3) o coeficiente aumenta 112778,848 unidades (de 72506,327 para 185285,175);

- da iteração 4211 para a iteração 42122 (o que corresponde a aglomeração de 3 *Clusters* para 2) o coeficiente aumenta 163141,154 unidades (de 185285,175 para 348426,329);

- da iteração 4212 para a iteração 4213 (o que corresponde a aglomeração de 2 *Clusters* para 1) o coeficiente aumenta 1075710,085 unidades (348426,329 para 1075710).

Assim, pressupõe-se, portanto existir na amostra de estudo entre 5 a 3 *Clusters* (figura 5.2).

Passa-se a descrever as 3 soluções obtidas com o agrupamento hierárquico com método de Ward na amostra em estudo, isto é, para 3, 4 e 5 *Clusters*.

Escala das variáveis				
	Propofol Ce (ug/ml)	Remifentanil Ce (ng/ml)	BIS	z-scores
moderados	2,4130	,5814	75,0138	1,5
baixos	,00	,00	27,00	-1,5
altos	4,52	2,02	98,00	1,5

1. Solução com 3 *Clusters*

Considerando a solução com 3 *Clusters*, obtemos a seguinte caracterização (figuras 5.2):

- o *Cluster Ward's 1*, constituído por 54,2% da amostra (2286) casos com valores de BIS baixo e os valores de propofol Ce ($\mu\text{g/ml}$) altos, sendo os valores remifentanil Ce ($\eta\text{g/ml}$) moderados;
- o *Cluster Ward's 2*, constituído por 15% da amostra (630casos), com valores nulos de propofol Ce ($\mu\text{g/ml}$), baixos de remifentanil Ce ($\eta\text{g/ml}$) e BIS valores altos;
- o *Cluster Ward's 3*, constituído por 30,8% da amostra (1293 casos), com valores altos de remifentanil Ce ($\eta\text{g/ml}$), valores baixos de propofol Ce ($\mu\text{g/ml}$) e BIS com valores baixos.

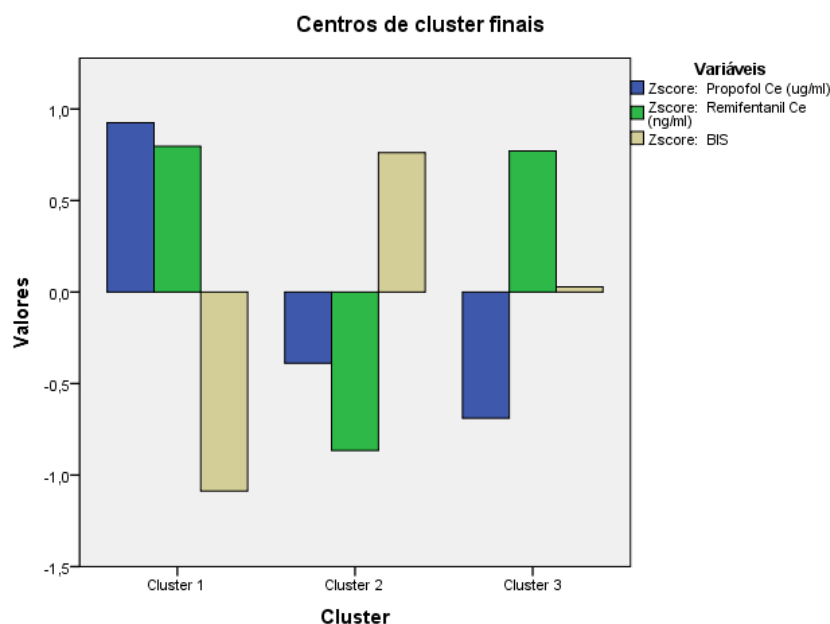


Figura 5.2: Valor estandardizado (z-scores) da mediana para pelos 3 *Clusters*. O propofol Ce ($\mu\text{g/ml}$), remifentanil Ce (ng/ml) e BIS

Com este método de agrupamento hierárquico e para uma solução com 3 *Clusters*, temos uma reprodutibilidade do modelo para 97,5% (casos corretamente classificados pelas funções discriminantes, tabela 5.3).

Duas funções discriminantes (Zscore- padronizados) (figura 5.3 e tabelas 5.4 e 5.5) com nível de significância na discriminação dos 3 *Clusters* de $p < 0,001$ (para as 2 funções, teste Lambda de Wilks) são criadas, em que a primeira função se correlaciona negativamente com os valores de BIS, enquanto que a segunda função se correlaciona positivamente com os valores de remifentanil Ce (ng/ml) e BIS, sendo que a função 1 explica 74,6% da variância entre os clusters e a função 2 explica 25,4%.

Tabela 5.3: Reprodutibilidade do modelo de agrupamento para 3 clusters em que a previsão para o grupo é a classificação prevista pelas funções discriminantes.						
Método Ward			Associação ao grupo prevista			Total
			1	2	3	
Original	Contagem	1	2251	2	33	2286
		2	18	564	48	630
		3	2	1	1295	1298
	%	1	98,5	,1	1,4	100,0
		2	2,9	89,5	7,6	100,0
		3	,2	,1	99,8	100,0

a. 97,5% de casos agrupados originais classificados corretamente.

Tabela 5.4: Valores próprios e coeficientes canônicos das 2 funções discriminantes.				
Função	Autovalor	% de variância	% cumulativa	Correlação canônica
1	4,104 ^a	74,6	74,6	,897
2	1,400 ^a	25,4	100,0	,764

a. As primeiras 2 funções discriminantes canônicas foram usadas na análise.

As funções discriminantes (Zscore- padronizados) são dadas por (tabela 5.6):

$$Função_1 = +0,104 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,667 \text{ Remifentanil Ce } (\eta\text{g/ml}) - 0,534 \text{ BIS}$$

$$Função_2 = -0,047 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,736 \text{ Remifentanil Ce } (\eta\text{g/ml}) + 0,827 \text{ BIS}$$

Tabela 5.5: Matriz de estrutura: coeficientes de correlação entre cada variável e as 2 funções discriminantes.		
	Função	
	1	2
Zscore: Remifentanil Ce (ηg/ml)	,819	,571
Zscore: BIS	-,769	,632
Zscore: Propofol Ce (ug/ml)	,335	-,474

Correlações entre grupos no conjunto entre variáveis discriminantes e funções discriminantes canônicas padronizadas Variáveis ordenadas por tamanho absoluto de correlação na função.

*. Maior correlação absoluta entre cada variável e qualquer função discriminante

Tabela 5.6: Coeficientes de funções discriminantes canônicas padronizados		
	Função	
	1	2
Zscore: Propofol Ce (ug/ml)	,104	-,047
Zscore: Remifentanil Ce(ηg/ml)	,677	,796
Zscore: BIS	-,534	,827

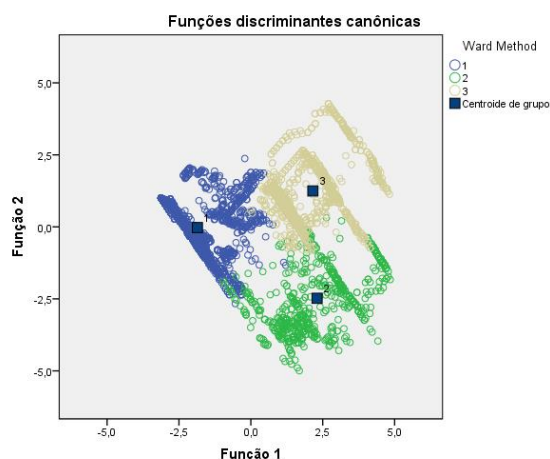


Figura 5.3: Funções discriminantes para o agrupamento hierárquico com método de Ward's param 3 clusters (z-scores).

1. Solução com 4 Clusters

Considerando a solução com 4 Clusters, obtemos a seguinte caracterização (figuras 5.4 e 5.5):

- o *Cluster Ward's 1*, constituído por 33,2% da amostra (1401 casos), com valores de BIS (27) baixos e os valores de propofol Ce (4,51 $\mu\text{g/ml}$) e elevados sendo os valores remifentanil Ce (1,012 $\eta\text{g/ml}$) moderados;
- o *Cluster Ward's 2*, constituído por 21% da amostra (885 casos), com valores baixos de propofol Ce (0 $\mu\text{g/ml}$) e de remifentanil Ce (0 $\eta\text{g/ml}$) e BIS(98) valores altos;
- o *Cluster Ward's 3*, constituído por 14,9% da amostra (630 casos), com valores altos de propofol Ce (3,584 $\mu\text{g/ml}$), valores de remifentanil Ce (0,269 $\eta\text{g/ml}$) baixos e BIS (58) com valores moderados.
- o *Cluster Ward's 4*, constituído por 30,8% da amostra (1298 casos), com valores moderados de remifentanil Ce (0,264 $\eta\text{g/ml}$) e valores moderados de propofol Ce (3,563 $\mu\text{g/ml}$) e BIS (76) com valores altos.

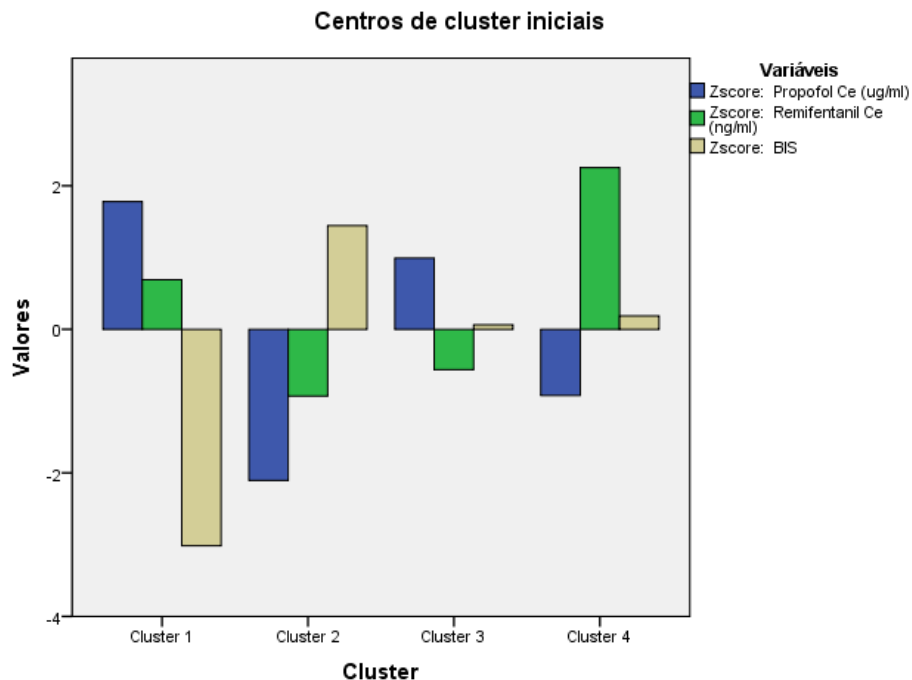


Figura 5.4: Valor estandardizado (z-scores) da mediana para pelos 4 *Clusters*. o propofol Ce ($\mu\text{g/ml}$), remifentanil Ce($\eta\text{g/ml}$) e BIS

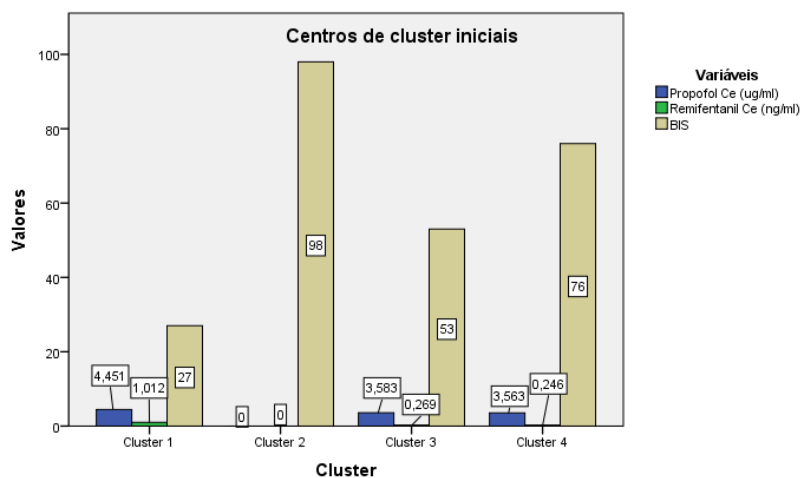


Figura 5.5: Valor da mediana para pelos 4 *Clusters*. o propofol Ce ($\mu\text{g/ml}$), remifentanil Ce (ng/ml) e BIS

Ao considerarmos mais um *Cluster* na amostra, diminuimos a reprodutibilidade do modelo para 96,7% (casos corretamente classificados pelas funções discriminantes, tabela 5.7).

Tabela 5.7: Resultados da classificação ^a							
Método Ward			Associação ao grupo prevista				Total
			1	2	3	4	
Original	Contagem	1	1313	58	0	30	1401
		2	4	879	2	0	885
		3	0	21	593	16	630
		4	0	2	5	1291	1298
	%	1	93,7	4,1	,0	2,1	100,0
		2	,5	99,3	,2	,0	100,0
		3	,0	3,3	94,1	2,5	100,0
		4	,0	,2	,4	99,5	100,0

a. 96,7% de casos agrupados originais classificados correctamente.

Três funções discriminantes (Zscore- padronizados) (figura 5.6 e tabelas 5.8 e 5.9) com nível de significância na discriminação dos 4 *Clusters* de $p < 0,001$ (para as 3 funções, teste Lambda de Wilks) são criadas em que:

- a função 1 correlaciona-se negativamente com os valores de: propofol Ce ($\mu\text{g/ml}$) e remifentanil Ce (ng/ml) positivamente com os valores de BIS, explicando 66,6% da variância entre os *Clusters*;

- a função 2 correlaciona-se negativamente com os valores de: propofol Ce ($\mu\text{g/ml}$) e positivamente com os valores de remifentanil Ce ($\eta\text{g/ml}$) e BIS, explicando 32,1% da variância entre os *Clusters*;
- a função 3 correlaciona-se negativamente com os valores de : propofol Ce ($\mu\text{g/ml}$) e positivamente com os valores de remifentanil ($\eta\text{g/ml}$) e BIS, explicando 1,3% da variância entre os *Clusters*;

Tabela 5.8: Valores próprios				
Função	Autovalor	% de variância	% cumulativa	Correlação canónica
1	5,350 ^a	66,6	66,6	,918
2	2,579 ^a	32,1	98,7	,849
3	,105 ^a	1,3	100,0	,308

a. As primeiras 3 funções discriminantes canónicas foram usadas na análise.

As funções discriminantes (Zscore- padronizados) são dadas por (tabela 5.9):

$$\text{Função}_1 = -0,387 \text{ Propofol Ce } (\mu\text{g/ml}) - 0,175 \text{ Remifentanil Ce } (\eta\text{g/ml}) + 0,757 \text{ BIS}$$

$$\text{Função}_2 = -0,278 \text{ Propofol Ce } (\mu\text{g/ml}) + 1,043 \text{ Remifentanil Ce } (\eta\text{g/ml}) + 0,407 \text{ BIS}$$

$$\text{Função}_3 = -0,908 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,281 \text{ Remifentanil Ce } (\eta\text{g/ml}) + 0,719 \text{ BIS}$$

Tabela 5.9: Matriz de estruturas			
	Função		
	1	2	3
Zscore: BIS	,914 [*]	,045	,403
Zscore: Remifentanil Ce (ng/ml)	-,508	,860 [*]	,047
Zscore: Propofol Ce (ug/ml)	-,566	-,302	,767 [*]

Tabela 5.10 : Coeficientes de funções discriminantes canónicas padronizados			
	Função		
	1	2	3
Zscore: Propofol Ce (ug/ml)	-,387	-,278	,908
Zscore: Remifentanil Ce ($\eta\text{g/ml}$)	-,175	1,043	,281
Zscore: BIS	,757	,407	,719

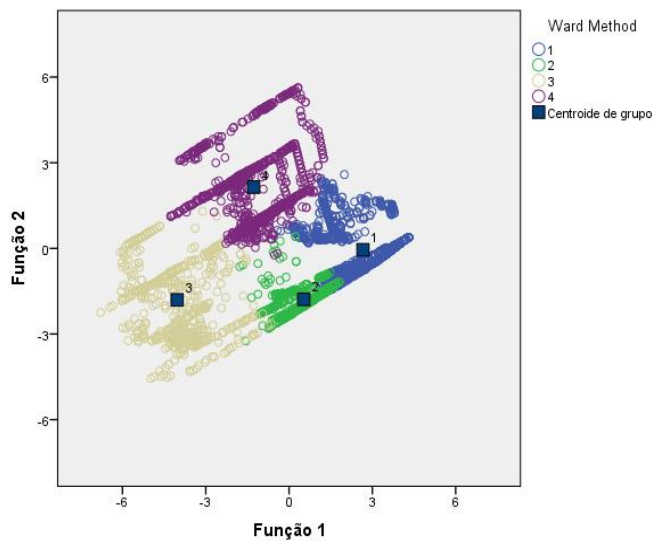


Figura 5.6: Representação das 2 principais funções discriminantes para o agrupamento hierárquico com método de Ward's para 4 Clusters

5. Solução com 5 Clusters

Considerando a solução com 5 Clusters, obtemos a seguinte caracterização (figuras 5.8):

- o Cluster Ward's 1, constituído por 33,2% da amostra (1401 casos) com valores de BIS (46) baixos, os valores de propofol Ce (3,665 $\mu g/ml$) altos e remifentanil Ce (0,844 $\eta g/ml$) moderadamente;
- o Cluster Ward's 2, constituído por 21% da amostra (885 casos), com valores nulos de propofol Ce (0 $\mu g/ml$) e de remifentanil Ce (0 $\eta g/ml$) e BIS (98) valores altos;
- o Cluster Ward's 3, constituído por 14,9% da amostra (630 casos), com valores baixos de propofol Ce (1,399 $\mu g/ml$) e valores baixos de remifentanil Ce (0,091 $\eta g/ml$) e BIS (81) com valores altos.
- o Cluster Ward's 4, constituído por 23% da amostra (971 casos), com valores altos de propofol Ce (4,440 $\mu g/ml$) e valores moderados de remifentanil Ce (0,51 $\eta g/ml$) e valores baixos de BIS (65);
- o Cluster Ward's 5, constituído por 7,8% da amostra (327 casos), com valores altos de remifentanil Ce (1,012 $\eta g/ml$) e valores altos de propofol Ce (4,451 $\mu g/ml$) e BIS (48) com valores baixos.

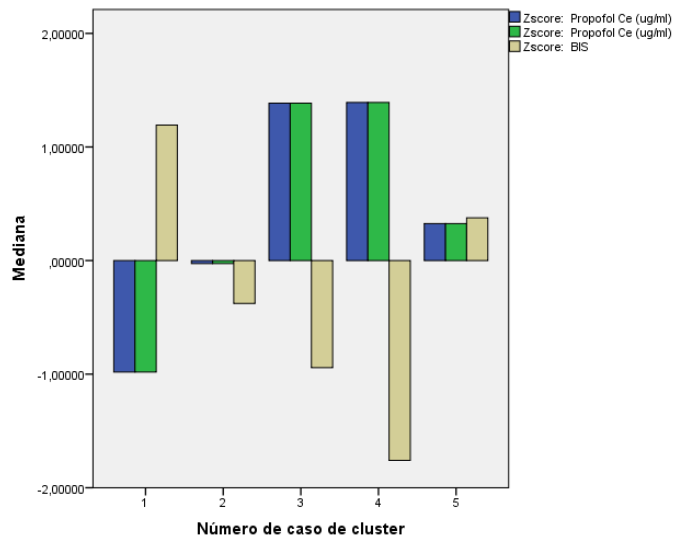


Figura 5.7: Valor estandardizado (z-scores) da mediana para pelos 5 Clusters. O propofol Ce ($\mu\text{g/ml}$), remifentaniil Ce (ng/ml) e BIS

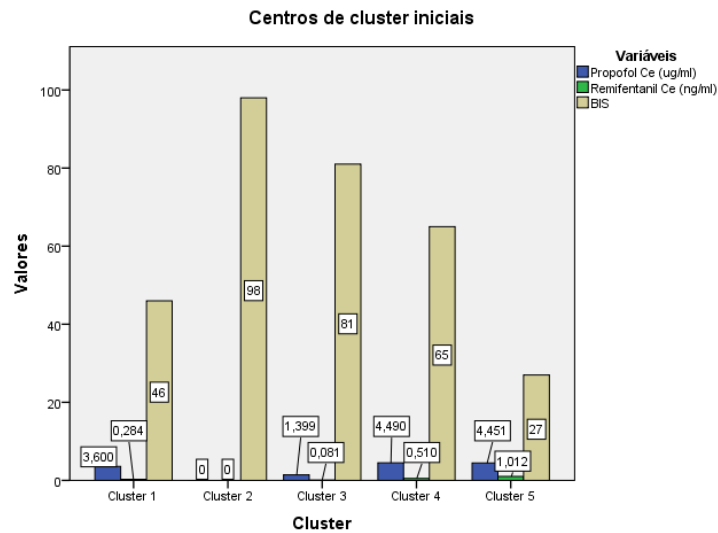


Figura 5.8: Valor real da mediana para pelos 5 Clusters. O propofol Ce (ug/ml), remifentaniil Ce (ng/ml) e BIS

Ao considerarmos mais um *Cluster* na amostra, diminuimos a reprodutibilidade do modelo para 95,2% 6 (casos corretamente classificados pelas funções discriminantes, tabela 5.11).

Tabela 5.11: Resultados da classificação ^a								
Método Ward			Associação ao grupo prevista					Total
			1	2	3	4	5	
Original	Contagem	1	1283	58	0	16	44	1401
		2	1	878	0	6	0	885

		3	0	17	586	27	0	630
		4	0	2	0	955	14	971
		5	0	0	0	17	310	327
	%	1	91,6	4,1	,0	1,1	3,1	100,0
		2	,1	99,2	,0	,7	,0	100,0
		3	,0	2,7	93,0	4,3	,0	100,0
		4	,0	,2	,0	98,4	1,4	100,0
		5	,0	,0	,0	5,2	94,8	100,0

a. 95,2% de casos agrupados originais classificados corretamente.

Quatro funções discriminantes (Zscore- padronizados) (figura 5.9 e tabelas 5.12 e 5.13) com nível de significância na discriminação dos 5 *Clusters* de $p < 0,001$ (para as 3 funções, teste Lambda de Wilks) são criadas, em que:

- a função 1 correlaciona-se negativamente com os valores de: propofol Ce ($\mu\text{g/ml}$) e de BIS e positivamente com os valores remifentanil Ce ($\eta\text{g/ml}$), explicando 69,3% da variância entre os *Clusters*;
- a função 2 correlaciona-se negativamente com os valores de: propofol Ce ($\mu\text{g/ml}$) e BIS e positivamente com os valores de *remifentanil* Ce ($\eta\text{g/ml}$), explicando 29,5% da variância entre os *Clusters*;
- a função 3 correlaciona-se positivamente com os valores de: propofol Ce ($\mu\text{g/ml}$) e Ce ($\eta\text{g/ml}$) e BIS, explicando 1,2% da variância entre os *Clusters*.

Tabela 5.12: Valores próprios				
Função	Autovalor	% de variância	% cumulativa	Correlação canônica
1	7,136 ^a	69,3	69,3	,937
2	3,034 ^a	29,5	98,8	,867
3	,123 ^a	1,2	100,0	,331

a. As primeiras 3 funções discriminantes canônicas foram usadas na análise.

As funções discriminantes (Zscore- padronizados) são dadas por (tabela 5.13):

$$Função_1 = -0,554 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,201 \text{ Remifentanil Ce } (\eta\text{g/ml}) - 0,8727 \text{ BIS}$$

$$Função_2 = -0,201 \text{ Propofol Ce } (\mu\text{g/ml}) + 1,081 \text{ Remifentanil Ce } (\eta\text{g/ml}) - 0,202 \text{ BIS}$$

$$Função_4 = 0,812 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,301 \text{ Remifentanil Ce } (\eta\text{g/ml}) + 0,699 \text{ BIS}$$

Tabela 5.13: Matriz de estruturas			
	Função		
	1	2	3
Zscore: BIS	,826	-,289	,484
Zscore: Remifentanil Ce (ng/ml)	-,282	,957*	,075
Zscore: Propofol Ce (ug/ml)	-,607	-,106	,788*

Tabela 5.14: Coeficientes de funções discriminantes canônicas padronizados			
	Função		
	1	2	3
Zscore: Propofol Ce (ug/ml)	-,554	-,227	,812
Zscore: Remifentanil Ce (ng/ml)	,201	1,081	,301
Zscore: BIS	,872	,202	,699

Correlações entre grupos no conjunto entre variáveis discriminantes e funções discriminantes canônicas padronizadas Variáveis ordenadas por tamanho absoluto de correlação na função.

*. Maior correlação absoluta entre cada variável e qualquer função discriminante

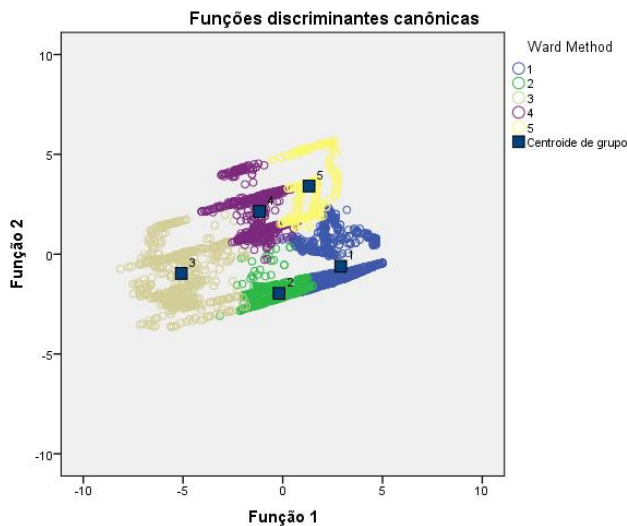


Figura 5.9: Representação das 2 principais funções discriminantes para o agrupamento hierárquico com método de Ward's para 5 Clusters

5.12. Clusters k-means variáveis estandardizadas

Comparando as 3 soluções obtidas através das técnicas de agrupamento hierárquico com método de Ward's, a solução que apresenta uma melhor reprodutibilidade é a solução com 3 *Clusters* (sendo a reprodutibilidade desta solução de 97,6%, respetivamente). Como tal, o número de *Clusters* que melhor descreve a amostra em estudo, (4214 casos), é de 3 *Clusters*.

No intuito de melhorar o agrupamento com 3 *Clusters*, obtido anteriormente, um novo modelo de agrupamento de dados por partição k-means foi criado baseado nas mesmas variáveis estandardizadas.

A solução obtida distingue os 3 *Clusters* em 3 iterações (tabela 5.15) sendo a distância entre os centros superior a 2 desvios padrões (tabela 5.16).

Neste agrupamento, o *Clusters* 1 e 2 são os *Clusters* mais próximos, uma vez que a distância entre os respectivos centros é de 3,036, e os *Clusters* 1 e 3 são os *Clusters* mais afastados, dado que a distância entre os respetivos centros é de 1,995.

Tabela 5.15: Número de iterações necessárias para a convergência da solução			
Iteração	Alteração em centros de cluster		
	1	2	3
1	1,519	1,520	1,300
2	,112	,058	,038
3	,097	,043	,004
4	,120	,048	,002
5	,088	,031	,005
6	,044	,015	,003
7	,034	,012	,003
8	,026	,009	,002
9	,016	,004	,003
10	,012	,004	,001

As iterações foram interrompidas porque o número máximo de iterações foi atingido.

As iterações não convergiram. A mudança de coordenada absoluta máxima para qualquer centro é ,010. A iteração atual é 10. A distância mínima entre os centros iniciais é 3,943

Tabela 5.16: Distância Euclidiana entre os centros dos 3 Clusters.

Cluster	1	2	3
1		3,036	1,995
2	3,036		2,367
3	1,995	2,367	

A tabela da ANOVA: todas as variáveis referidas apresentam elevados valores de Z, contribuindo fortemente para a definição dos grupos. (valor de Z, tabela 5.17).

Tabela 5.17: Tabela de ANOVA

	Cluster		Erro		Z	Sig.
	Quadrado Médio	df	Quadrado Médio	df		
	Zscore: Propofol Ce (ug/ml)	1086,504	2	,484		
Zscore: Remifentanil Ce (ng/ml)	1575,918	2	,252	4211	6253,690	,000
Zscore: BIS	1464,409	2	,305	4211	4801,986	,000

Analisando e caracterizando os 3 *Clusters* temos (figura 5.11):

- o *Clusters* 1, é portanto, caracterizado por 18,2% da amostra (769 casos), com valores de BIS altos e os valores de propofol Ce ($\mu\text{g/ml}$) baixos, sendo os valores remifentanil Ce ($\eta\text{g/ml}$) moderadamente baixos;
- o *Clusters* 2, constituído por 51,7% da amostra (2128 casos), com valores altos de propofol Ce ($\mu\text{g/ml}$) e de remifentanil Ce ($\eta\text{g/ml}$) e BIS valores baixos;
- o *Clusters* 3, constituído por 30,1% da amostra (1267 casos), com valores altos de remifentanil Ce ($\eta\text{g/ml}$) e valores moderadamente altos de propofol Ce ($\mu\text{g/ml}$) e BIS com valores baixos.

Resumindo, o *Cluster* 1 é, portanto, caracterizado por ter elevados valores de propofol Ce ($\mu\text{g/ml}$) e elevados valores de BIS; o *Cluster* 2 é caracterizado por ter valores altos de propofol Ce ($\mu\text{g/ml}$) e de remifentanil Ce ($\eta\text{g/ml}$) e valores altos de BIS e o *Cluster* 3 é caracterizado por ter um elevado valor de remifentanil Ce ($\eta\text{g/ml}$).

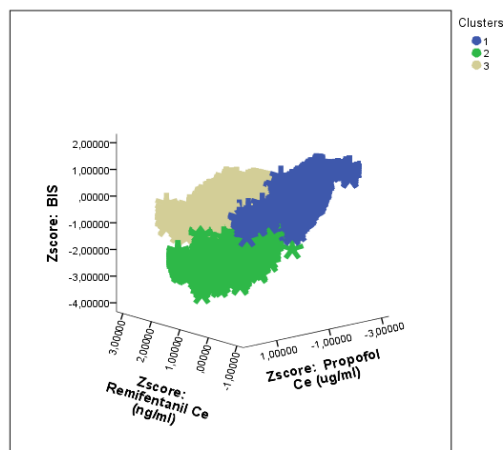


Figura 5.10: Representação 3D dos 3 Clusters Ward's em função dos valores (z-scores) de propofol Ce ($\mu\text{g/ml}$), remifentanil Ce (ng/ml) e BIS

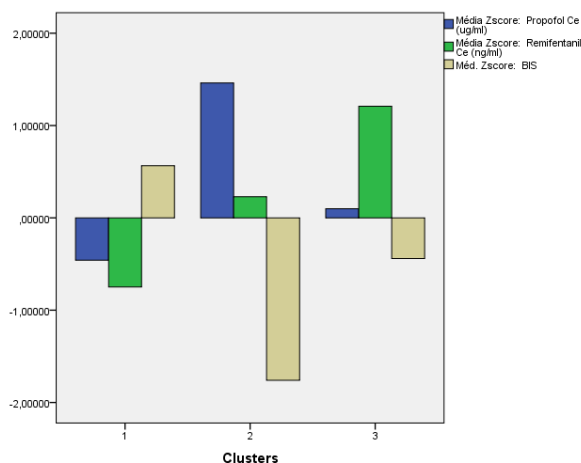


Figura 5.11: Valor estandardizado (z-scores) da mediana para pelos 3 Clusters. O propofol Ce ($\mu\text{g/ml}$), remifentanil Ce (ng/ml) e BIS.

O modelo tem uma reprodutibilidade de 95,6% (casos corretamente classificados pelas funções discriminantes, (tabela 5.18).

Duas funções discriminantes (figura 5.12 e tabelas 5.18 e 5.19) com nível de significância na discriminação dos 3 *Clusters* de $p < 0,001$ (para as 2 funções, teste Lambda de Wilks) são criadas, em que a primeira função se correlaciona negativamente com os valores de BIS, enquanto a segunda função se correlaciona positivamente com os valores de Remifentanil Ce (ng/ml), e BIS, sendo que a função 1 explica 73,5% da variância entre os clusters e a função 2 explica 26,5%.

Tabela 5.18: Reprodutibilidade do modelo de agrupamento clusters k-means em que a previsão para o grupo é a classificação prevista pelas funções discriminantes.

Número de caso de cluster			Associação ao grupo prevista			Total
			1	2	3	
Original	Contagem	1	621	83	65	769
		2	0	2154	24	2178
		3	2	13	1252	1267
	%	1	80,8	10,8	8,5	100,0
		2	,0	98,9	1,1	100,0
		3	,2	1,0	98,8	100,0

a. 95,6% de casos agrupados originais classificados corretamente.

Tabela 5.19: Valores próprios e coeficientes canônicos das 2 funções discriminantes.

Função	Autovalor	% de variância	% cumulativa	Correlação canônica
1	3,571 ^a	73,5	73,5	,884
2	1,286 ^a	26,5	100,0	,750

a. As primeiras 2 funções discriminantes canônicas foram usadas na análise.

As funções discriminantes são dadas por (tabela 5.21):

$$Função_1 = 0,203 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,739 \text{ Remifentanil Ce } (\eta\text{g/ml}) - 0,405 \text{ BIS}$$

$$Função_2 = -0,336 \text{ Propofol Ce } (\mu\text{g/ml}) + 0,706 \text{ Remifentanil Ce } (\eta\text{g/ml}) ,657 + \text{ BIS}$$

Tabela 5.20: Matriz de estrutura: coeficientes de correlação entre cada variável e as 2 funções discriminantes.

	Função	
	1	2
Zscore: Remifentanil Ce ($\eta\text{g/ml}$)	,859 [*]	,513
Zscore: BIS	-,711	,609
Zscore: Propofol Ce ($\mu\text{g/ml}$)	,382	-,651 [*]

Correlações entre grupos no conjunto entre variáveis discriminantes e funções discriminantes canônicas padronizadas. Variáveis ordenadas por tamanho absoluto de correlação na função.

*. Maior correlação absoluta entre cada variável e qualquer função discriminante

Tabela 5.21: Coeficientes canônicos das 2 funções discriminantes.

	Função	
	1	2
Zscore: Propofol Ce ($\mu\text{g/ml}$)	,203	-,366
Zscore: Remifentanil Ce ($\eta\text{g/ml}$)	,739	,706
Zscore: BIS	-,405	,657

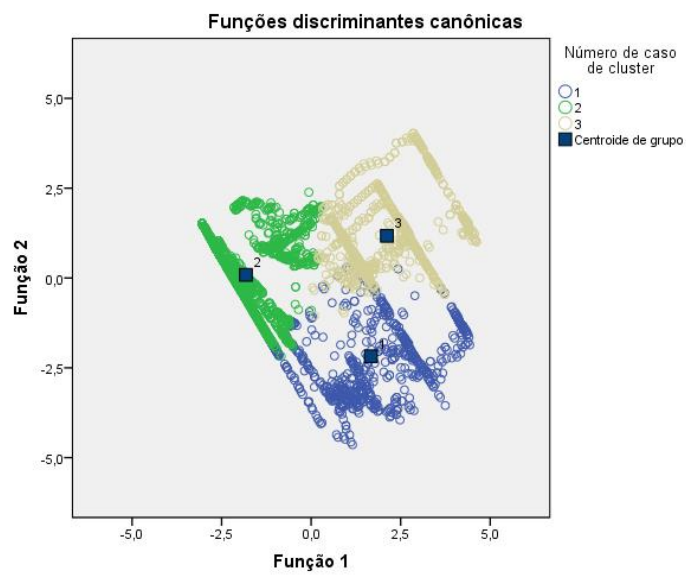


Figura 5.12: Funções discriminantes para o agrupamento por clusters kmeans

5.13. Clusters k-means por Amplitude

No intuito de melhorar o agrupamento com 3 *Clusters*, obtido anteriormente, um novo modelo de agrupamento de dados por partição k-means por amplitude foi criado baseado nas mesmas variáveis com maior amplitude.

A solução obtida distingue os 3 *Clusters* em 9 iterações (tabela 5.23) sendo a distância entre os centros superior a 2 desvios padrões (tabela 5.16).

Neste agrupamento, os *Clusters* 1 e 2 são os *Clusters* mais próximos, uma vez que a distância entre os respectivos centros é de 16,860, e os *Clusters* 2 e 3 são os *Clusters* mais afastados, uma vez que a distância entre os respectivos centros é de 39,341.

Tabela 5.22: Distância entre centros de cluster finais			
Cluster	1	2	3
1		16,860	22,483
2	16,860		39,341
3	22,483	39,341	

Tabela 5.23: Histórico de iteração ^a			
Iteração	Alteração em centros de cluster		
	1	2	3
1	8,101	8,180	25,258
2	5,667	4,878	4,341
3	,780	,582	,784
4	,000	,001	,001
5	2,403E-7	6,907E-7	3,496E-7
6	1,334E-10	7,525E-10	2,335E-10
7	7,142E-14	8,286E-13	1,493E-13
8	,000	2,220E-16	,000
9	,000	,000	,000

a. Convergência alcançada devido a nenhuma ou pequena alteração em centros de cluster. A mudança de coordenada absoluta máxima para qualquer centro é ,000. A iteração atual é 9. A distância mínima entre os centros iniciais é 33,308.

A tabela da ANOVA: todas as variáveis referidas apresentam elevados valores de Z, contribuindo fortemente para a definição dos grupos, (valor de Z, tabela 5.23).

Tabela 5.23: Tabela de ANOVA						
	Cluster		Erro		Z	Sig.
	Quadrado Médio	df	Quadrado Médio	df		
Propofol Ce (ug/ml)	1139,839	2	,772	4211	1476,22 1	,000
Remifentanil Ce (ng/ml)	284,167	2	,256	4211	1110,64 2	,000
BIS	448952,79 1	2	40,520	4211	11079,9 19	,000

Os testes F devem ser usados apenas para finalidades descritivas porque os cluster foram escolhidos para maximizar as diferenças entre os casos em clusters diferentes. Os níveis de significância observados não estão corrigidos para isso e, dessa forma, não podem ser interpretados como testes da hipótese de que as médias de cluster são iguais.

Analisando e caracterizando os 3 *Cluster* temos (figura 5.13):

- o *Clusters* 1, constituído por 14% da amostra (539 casos), com valores de BIS de 65 e os valores de propofol Ce (4,49 µg/ml) altos sendo os valores remifentanil Ce (0,51 ng/ml) moderados;
- o *Clusters* 2, constituído por 21% da amostra (885 casos), com valores de propofol Ce (0µg/ml) e de remifentanil Ce (0 ng/ml) e BIS (98);
- o *Clusters* 3, constituído por 12% da amostra (490 casos), com valores de remifentanil Ce (1,02 ng/ml) e valores altos de propofol Ce (4,551 µg/ml) e BIS com valores baixos de 27.

Resumindo, o *Cluster* 1 é, portanto, caracterizado por ter elevados valores de propofol Ce (4,49 µg/ml) e valores de BIS (65) moderados; o *Cluster* 2 é caracterizado por ter valores nulos de propofol Ce (0 µg/ml) e de remifentanil Ce (0 ng/ml) e valores máximos de BIS (98) e o *Cluster* 3 é caracterizado por ter elevados valores de propofol Ce (4,551 µg/ml) e de remifentanil Ce (1,02 ng/ml) e um BIS muito baixo.

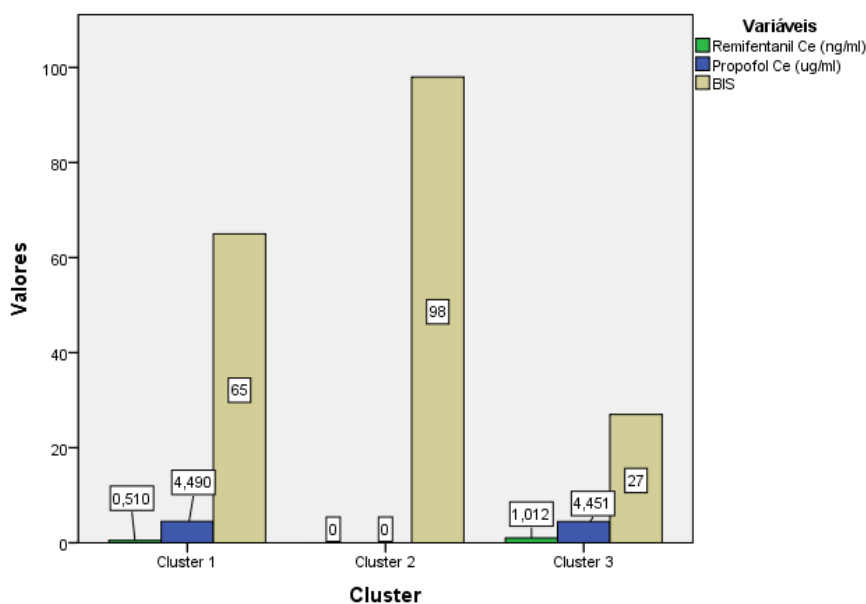


Figura 5.13: Valor por amplitude para propofol Ce (µg/ml), remifentanil Ce (ng/ml) e BIS pelos 3 Clusters.

Tabela 5.18: Tabela de Resultados da classificação ^a						
		Distância do caso de seu centro de cluster de classificação	Associação ao grupo prevista			Total
			1,0000	2,00000	3,00000	
Original	Contagem	1,00000	196	280	113	589
		2,00000	217	500	168	885
		3,00000	159	86	245	490
		Casos não agrupados	302	298	1650	2250
	%	1,00000	33,3	47,5	19,2	100,0
		2,00000	24,5	56,5	19,0	100,0
		3,00000	32,4	17,6	50,0	100,0
		Casos não agrupados	13,4	13,2	73,3	100,0

a. 47,9% de casos agrupados originais classificados corretamente.

O modelo tem uma reprodutibilidade de 47,9% (casos corretamente classificados pelas funções discriminantes, tabela 5.18).

Duas funções discriminantes (tabelas 5.19 e 5.20) com nível de significância na discriminação dos 3 clusters de $p < 0,001$ (para as 2 funções, teste Lambda de Wilks) são criadas, em que a primeira função se correlaciona negativamente com os valores de BIS, enquanto que a segunda

função se correlaciona positivamente com os valores de remifentanil Ce (ng/ml), e BIS, sendo que a função 1 explica 93,1% da variância entre os clusters e a função 2 explica 6,9%.

Tabela 5.19: Valores próprios e coeficientes canônicos das 2 funções discriminantes.				
Função	Autovalor	% de variância	% cumulativa	Correlação canônica
1	,153 ^a	93,1	93,1	,364
2	,011 ^a	6,9	100,0	,106

a. As primeiras 2 funções discriminantes canônicas foram usadas na análise.

As funções discriminantes são dadas por (tabela 5.21):

$$Função_1 = -0,283 \text{ Propofol Ce } (\mu\text{g/ml}) - 0,240 \text{ Remifentanil Ce } (\text{ng/ml}) + 0,666 \text{ BIS}$$

$$Função_2 = 0,803 \text{ Propofol Ce } (\mu\text{g/ml}) - 0,646 \text{ Remifentanil Ce } (\text{ng/ml}) + 0,119 \text{ BIS}$$

Tabela 5.20: Matriz de estrutura: coeficientes de correlação entre cada variável e as 2 funções discriminantes.		
	Função	
	1	2
BIS	,974*	,054
Propofol Ce (ug/ml)	-,677	,698*
Remifentanil Ce (ng/ml)	-,663	-,671*

Correlações entre grupos no conjunto entre variáveis discriminantes e funções discriminantes canônicas padronizadas

Variáveis ordenadas por tamanho absoluto de correlação na função.

*. Maior correlação absoluta entre cada variável e qualquer função discriminante

Tabela 5.21: Coeficientes de funções discriminantes canônicas padronizados		
	Função	
	1	2
Propofol Ce (ug/ml)	-,283	,803
Remifentanil Ce (ng/ml)	-,240	-,646
BIS	,666	,119

5.14. Caracterização dos Clusters

A diferença nos dois modelos de agrupamento por partição k-means reside portanto no facto de um considerar as variáveis com maior amplitude entre elas e outro considerar as variáveis estandardizadas. A vantagem de utilizar o desvio padrão absoluto em vez do desvio padrão normal (função do quadrado das diferenças) reside no facto deste ser menos sensível a outliers (valores extremos), tornando-se desta forma mais robusto.

Concluimos que, o *Cluster 1* é, portanto, caracterizado por ter elevados valores de propofol Ce (4,49 µg/ml) e de Remifentanil Ce (0,51 ng/ml) baixos, e valores de BIS (68) moderados; o *Cluster 2* é caracterizado por ter valores baixos de propofol Ce (0 µg/ml) e de remifentanil Ce (0 ng/ml) e valores máximos de BIS (98) e o *Cluster 3* é caracterizado por ter elevados valores de propofol Ce (4,551 µg/ml) e de remifentanil Ce (1,02 ng/ml) e um BIS (27) muito baixo.

Estes resultados demonstram claramente a interação farmacodinâmica dos dois fármacos, quando analisamos o Cluster 1 e o Cluster 3. Para concentrações semelhantes de propofol o efeito no BIS é claramente diferente dependendo da grandeza da concentração de remifentanil.

A reprodutibilidade no modelo com as variáveis estandardizadas é de 95,6% e no modelo com as variáveis com maior amplitude é 49,7%.

Capítulo 6

Considerações finais

6.1. Conclusões

Neste trabalho apresentaram-se diferentes contribuições na área da estatística. Uma destas contribuições foi o desenvolvimento dos dois modelos lineares para modelar o efeito dos fármacos no sinal cerebral BIS. Da nossa análise, podemos concluir que o teste de significância da equação de Regressão Linear indicou que os modelos construídos podem ser considerado significativo para um nível de significância de 5% e, conseqüentemente, o modelo de regressão é válido. Em suma, **o modelo é altamente significativo**. Como o *p-value* encontrado foi inferior a 0,05, podemos assegurar que o modelo de regressão considerado é melhor que a média para prever os valores do BIS. Com um nível de confiança de 95%, a variável remifentanil é a mais significativa, sendo aquela que tem maior contribuição individual (-46,858). Os resíduos têm uma distribuição normal e os testes de ajustamentos fornecem os *p-values* inferiores aos níveis de significância usual de 0,05, o modelo foi declarado ajustado. Relativamente ao RLS o $r^2=0,541$, podemos afirmar que 54,1% da variabilidade da variável dependente BIS é explicada pela variável independente do modelo ajustado. O valor do coeficiente de correlação é $r=0,736$. No que concerne ao modelo RLM o $r^2=0,698$, podemos afirmar que 69,8% da variabilidade da variável dependente BIS é explicada pelas variáveis independentes do modelo ajustado. O valor do coeficiente de correlação é $r=0,836$.

Os dados utilizados neste trabalho, foram recolhidos em voluntários e durante sedação. Nessa perspectiva, não existiu muita variabilidade individual entre estes e as condições foram controladas. Permitindo assim, a obtenção de dados fiáveis e limpos de artefactos. Os modelos encontrados mostram claramente a relação entre os dois fármacos e o seu efeito cerebral (profundidade de anestesia/sedação). Um fármaco analgésico (para alívio da dor) opiáceo como o remifentanil não tem propriedades de sedação e por si só não deprime o sistema nervoso central (EEG e BIS). Este trabalho demonstra claramente, que quando o remifentanil é administrado com o propofol (um hipnótico) o efeito deste último é potenciado, levando o sinal BIS a valores bastante baixos.

Este trabalho demonstra ainda ser possível construir um modelo que permite agrupar as concentrações dos fármacos, com base no efeito no sinal cerebral BIS, com o apoio de técnicas matemáticas (de agrupamento e discriminantes), numa amostra de 4214 casos de 8 pacientes.

Estes resultados demonstram claramente a interação farmacodinâmica dos dois fármacos, quando analisamos o Cluster 1 e o Cluster 3. Para concentrações semelhantes de propofol o efeito no BIS é claramente diferente dependendo da grandeza da concentração de remifentanil.

6.2. Trabalho futuro

No conjunto dos trabalhos desenvolvidos e apresentados nesta dissertação muitas foram as portas que ficaram abertas para futuras contribuições. Assim ficam assinaladas algumas das ideias para trabalho futuro.

Com o modelo de regressão linear, iremos futuramente construir uma superfície de resposta que possa ser utilizada como um modelo de navegação/previsão para os clínicos num ambiente de sedação com propofol e remifentanil.

Outra das contribuições deste trabalho foi estimar o efeito que as concentrações de fármacos têm na depressão do electroencefalograma através de um modelo de regressão linear. Os resultados obtidos até agora podem ser considerados como um passo para a definição de um modelo simples para estimar a dosagem individualizada de propofol e remifentanil com base nos atributos dos pacientes. Podem ser destacadas várias vantagens inerentes à estrutura e às propriedades associadas a modelos lineares. Em primeiro lugar, permite facilmente a inclusão de outros parâmetros clínicos a fim de obter uma representação adequada das características individuais do paciente e da forma de pensar do clínico. Permite, também, a utilização direta de outras medidas de DoA, em vez da concentração de efeito pretendida, sempre que elas possam ser consideradas como função dessa concentração. Finalmente, o modelo linear é um modelo estatístico, permitindo assim outros procedimentos inferenciais, tais como a construção de intervalos de confiança e a realização de testes de hipóteses.

A análise de *Clusters* futuramente ainda deverá ser mais explorada com a inclusão de mais voluntários ou doentes, para tentar perceber se a intersecção entre os dois fármacos pode ser afectada também por características individuais de cada pessoa.

Perante os resultados obtidos e tendo em vista a sua aplicação na área da medicina, deixam-se algumas considerações e sugestões:

- É decisivo o ajuste da escolha de variáveis e métodos em função do contexto da Anestesia, percebendo e analisando os fenômenos alvos do estudo estatístico de modo que todas as opções a fazer sejam devidamente fundamentadas tanto na vertente estatística como na vertente Anestésica.
- É relevante alertar os utilizadores da Estatística na área da Saúde para a importância da correta utilização dos métodos, não só validando pressupostos mas também na escolha dos modelos.
- A concretização deste estudo viabilizou uma experiência enriquecedora, apesar de haver requerido grande empenho e esforço. Atende-se que este estudo proporcione, de algum modo, possa contribuir, ainda que de forma modesta, para o progresso de mudanças a nível das estratégias adotadas e da clareza do conhecimento veiculado.

Bibliografia

- [1] M. M. Silva, T. Wigren, and T. Mendonça, "Nonlinear identification of a minimal neuromuscular blockade model in anesthesia," *IEEE Transactions on Control Systems Technology*, vol. 20, pp. 181-188, January 2012.
- [2] W. M. Haddad, K. Y. Volyansky, and J. M. Bailey, "Neuroadaptive output feedback control for automated anesthesia with noisy eeg measurements," in *2008 American Control Conference*, pp. 813-818, June 2008.
- [3] D. V. Caiado, J. M. Lemos, and B. A. Costa, "Robust pole-placement control of neuromuscular blockade," in *Program of the 10th Portuguese Conference on Automatic Control*, (Funchal, Portugal), University of Madeira, July 2012.
- [4] M. Martins da Silva, T. Mendonça, and T. Wigren, "Online nonlinear identification of the effect of drugs in anaesthesia using a minimal parameterization and BIS measurements," *Tech. Rep. 2010-008*, Department of Information Technology, Uppsala University, Março 2010.
- [5] J. M. Ehrenfeld, R. D. Urman, and S. Segal, *Anesthesia Student Survival Guide: A Case-based Approach*. New York: Springer, 2010.
- [6] A. R. Absalom and M. M. R. F. Struys, *Overview of target Controlled Infusion and Total Intravenous Anaesthesia*. Academia Press, 2 ed., 1997.
- [7] Cherink DA, Gillings D, Liane H, et al. Validity and reliability of the observer's assessment of alertness/sedation scale: study with intravenous midazolam. *Journal of Clinical Psychopharmacology* 1994; 10: 244-51.
- [8] R. Meier, J. Nieuwland, S. Hacisalihzade, D. Steck, and A. Zbinden, "Fuzzy control of blood pressure during anesthesia with isourane," in *IEEE International Conference on Fuzzy Systems*, pp. 981-987, March 1992.
- [9] M. Agarwal and R. Grieths, "Monitoring the depth of anaesthesia," *Anaesthesia and Intensive Care Medicine*, vol. 5 (10), pp. 343-344, October 2004.
- [10] S. A. Abdulla and P. Wen, "Depth of anaesthesia patient models and control," in *2011 IEEE/ICME International Conference on Complex Medical Engineering*, (Harbin, China), pp. 37-41, May 2011.
- [11] H. Labbaf, M. Aliyari, and M. Teshnehlab, "A new approach in drug delivery control in anesthesia," in *Proceedings of the IEEE International Conference on Systems, Man and Cybernetics*, Istanbul, Turkey, 10-13 October 2010, pp. 2064-2068, IEEE, 2010.

- [12] O. Simanski, A. Schubert, R. Kaehler, M. Janda, and J. B. and, "Automatic drug delivery in anesthesia: From the beginning until now," in 2007 Mediterranean Conference on Control and Automation, (Athens-Greece), pp. 1-6, July 2007.
- [13] T. Mendonça, J. M. Lemos, H. Magalhães, P. Rocha, and S. Esteves, "Drug delivery for neuromuscular blockade with supervised multimodel adaptive control," IEEE Transactions on Control Systems Technology, vol. 17 (6), pp. 1237-1244, November 2009.
- [14] B. Appadu and A. Vaidya, "Monitoring techniques: neuromuscular blockade and depth of anaesthesia," Anaesthesia and Intensive Care Medicine, vol. 9 (6), pp. 247- 250, June 2008.
- [15] A. T. Hindle, M. O. Columb, and M. V. Shah, "Drug interactions and anaesthesia," Current Anaesthesia and Critical Care, vol. 6 (2), pp. 103-112, 1995.
- [16] J. G. Bovill, "Adverse drug interactions in anesthesia," Journal of Clinical Anesthesia, vol. 9 (6), pp. 3S-13S, 1997.
- [17] I. Kissin, "Anesthetic interactions following bolus injections," Journal of Clinical Anesthesia, vol. 9 (6), pp. 14S-17S, 1997.
- [18] N. C. D. Castro, "Predictive control of depth of anaesthesia," Master's thesis, Grenoble INP -ESISAR and INESC-ID, 2008.
- [19] H. Alonso, J. M. Lemos, and T. F. Mendonça, "A target control infusion method for neuromuscular blockade based on hybrid parameter estimation," in 30th Annual International Conference of the IEEE. Engineering in Medicine and Biology Society, (Vancouver, British Columbia, Canada), pp. 707-710, August 2008.
- [20] G. Kenny and N. Sutcliffe, "Target-controlled infusions: Stress free anesthesia?," Journal of Clinical Anesthesia, vol. 8, pp. S15-S20, May 1996.
- [21] S. E. Milne and G. N. C. Kenny, "Target controlled infusions," Current Anaesthesia and Critical Care, vol. 9, no. 4, pp. 174-179, 1998.
- [22] M. M. Silva, C. Sousa, R. Sebastião, J. Gama, T. Mendonça, P. Rocha, and S. Esteves, "Total mass tci driven by parametric estimation," in 17th Mediterranean Conference on Control and Automation, (Makedonia Palace, Thessaloniki, Greece), pp. 24-26, June 2009.
- [23] P. F. White, "Intravenous anesthesia and analgesia: what is the role of target-controlled infusion," Journal of Clinical Anesthesia, vol. 8, no. 3, pp. 26-28, 1998.
- [24] Murteira B., Ribeiro C., Silva J., Pimenta C. – Introdução à estatística. McGraw Hill (ISBN: 972-773-116-3), 2002.
- [25] Bowers D. – Medical statistics from scratch. John Wiley & Sons (ISBN: 0-470-84474-4), 2002.
- [26]. Maroco, J.; Análise Estatística – Com utilização do SPSS, 2ª edição; Edições Sílabo; 2003.
- [27]. Hair, JF, et al. Multivariate Data Analysis, 7th Edition, Pearson Education Limited 2014

- [28]. C.S. Nunes, F. Lobo, and P. Amorim, "Remifentanil adjusted Propofol Response Surface Model for BIS during Sedation in Volunteers," in Proceedings of the 2010 Annual Meeting of the American Society Anesthesiologists, 2010, p. A586.
- [29]. Branco, J.A. e Pires, A.M.; Introdução aos Métodos Estatísticos Robustos; Edições SPE; 2007.
- [30]. Pestana, M. H. e Gageiro, J.N.; Análise de Dados para Ciências Sociais: A Complementaridade do SPSS. 4ª ed., Lisboa: Sílabo; 2005a.
- [31] O. Simanski, A. Schubert, R. Kaehler, M. Janda, and J. B. and, "Automatic drug delivery in anesthesia: From the beginning until now," in 2007 Mediterranean Conference on Control and Automation, (Athens-Greece), pp. 1{6, July 2007.
- [32] C. F. Minto, T. W. Schnider, T. D. Egan, E. Youngs, H. J. Lemmens, P. L. Gambus, V. Billard, J. F. Hoke, K. H. Moore, D. J. Hermann, K.T. Muir, M. J. W., and S. L. Shafer, Influence of age and gender on the pharmacokinetics and pharmacodynamics of remifentanil. I. Model development, *Anesthesiology*, vol. 86, 1997, pp 10-23.
- [33] Maroco, João; Análise Estatística com utilização do SPSS. Lisboa, Edições Sílabo 2010.
- [34] Reis, Elizabeth; Estatística Multivariada Aplicada. 2ª ed., Edições Sílabo 2001.
- [35] Gower, J.C.; "Measures of Similarity, Dissimilarity, and Distance," in Encyclopedia of Statistical Sciences, Vol 5, Eds. S. Kotz, N.L. Johnson and C.B. Read, New York: John Wiley and Sons, 1985, 397-405.
- [36] Legendre, P., Dallot, S., and Legendre, L. ;"Succession of Species Within a Community: Chronological Clustering, with Applications to Marine and Freshwater Zooplankton," *American Naturalist*, 125, 1985, 257-288
- [37] Jonhson, R.A., Wichen D.W ; Applied Multivariate Statistical Methods. Prentice Hall, 2002.
- [38] Zhang, T., R. Ramkrishishnon, M. Livny; Birch: An efficient data clustering method for very large databases 1996
- Disponível em: <http://www.cs.sfu.ca/CourseCentral/459/han/papers/zhang96.pdf>.
- [39] Kaufman, L.; Rousseeuw, P. J.); Finding Groups in Data. Wiley Interscience, 1990.
- [40] T. W. Schnider, C. F. Minto, S. L. Schafer, P. L. Gambus, C. Andersen, D. B. Goodale, e E. J. Youngs. The influence of age on propofol pharmacodynamics. *Anesthesiology*, 90: 1502–1516, 1999. 16.
- [41] C. F. Minto, T. W. Schnider, e S. L. Shafer. Pharmacokinetics and pharmacodynamics of remifentanil: Ii. model application. *Anesthesiology*, 86: 24–33, 1997. 16.

- [42] Amado C., Pires A. – Uso de técnicas de reamostragem para selecção de variáveis em análise discriminante. Curia, Portugal: V Congresso Anual da Sociedade Portuguesa de Estatística, 1997.
- [43] PIRES A., BRANCO J., TURKMAN M. – Comparação de métodos de análise discriminante no diagnóstico da doença coronária. Luso, Portugal: II Congresso Anual da Sociedade Portuguesa de Estatística, 1994.
- [44] Tatsuoka M. – Multivariate analysis. Techniques for educational and psychological research. Publishing Company (ISBN: 0-02-419120-5), 1988. [45] T. W. Schnider, C. F. Minto, S. L. Schafer, P. L. Gambus, C. Andersen, D. B. Goodale, e E. J. Youngs. The influence of age on propofol pharmacodynamics. *Anesthesiology*, 90:1502–1516, 1999. 16
- [46] Montgomery. D. C.: (2001): Design and Analysis of Experiments, 5th Ed, John Wiley & Sons.
- [47] Morrison (1984): Multivariate Statistical Methods. 2nd Edition, International Student Edition.
- [48] Oliveira, T. A. (2004): Estatística Aplicada, Universidade Aberta.
- [49] Romesburg, H. (1984). Cluster analysis for researchers. USA: Lifetime Learning Publications
- [50] Gower, J. C. A general coefficient of similarity and some of its properties. *Biometrics*, v. 27, p. 857-874, 1971
- [51] Johnson, R.A.; Wichern, D.W. Applied multivariate statistical analysis. New Jersey: Prentice Hall, 1998. 816p.
- [52] Anderberg, M. R. (1973), Cluster Analysis for Applications, New York: Academic Press.
- [53] Russel, P. F., Rao, T. R. (1940) On habitat and association of species of anopheline larvae in south-eastern Madras. *J. Malaria Inst. India* 3: 154-178.
- [54] Jaccard, P. (1908) Nouvelles recherches sur la distribution florale. *Bull Soc Vaud Sci Nat*, 44, 223-270.
- [55] ED 9164 Design of Hydraulics ... Richard A. Johnson and Dean W. Wichern, Applied Multivariate Statistical . Analysis, Pearson Education, Asia, 5th Edition, 2002.
- [56] Olson, C.L. 1976. On choosing a test statistic multivariate analysis of variance. *Psychological Bulletin*, 83, 579-586.
- [57] Schatzoff, M., 1966. Sensitivity comparisons among test of the general linear hypothesis. *Journal of the American Statistical Association*, 61, 415-435
- [58] Santos, C. M. A.; Estatística Descritiva – Manual de auto-aprendizagem; Edições Sílabo; 2007.

[59] SHARMA, S. Applied multivariate techniques. New York: John Wiley & Sons, 1996.

[60] SPSS – Classification trees. SPSS Inc. (ISBN: 1-56827-354-1), 2004b.

[61] Pearson, K. (1978) The History of Statistics in the 17 th and 18 th Centuries, Edited by E.S. Pearson. New York: MacMillan.

Anexos A

Tabela das correlações

		Propofol Ce (ug/ml)	Remifentanil Ce (ng/ml)	Idade Voluntários	Altura	Peso	BIS
Propofol Ce (ug/ml)	Correlação de Pearson	1	,260**	-,088**	,259**	,372**	-,736**
	Sig. (2 extremidades)		,000	,000	,000	,000	,000
	N	4214	4214	4214	4214	4214	4214
Remifentanil Ce (ng/ml)	Correlação de Pearson	,260**	1	-,058**	-,191**	-,249**	-,574**
	Sig. (2 extremidades)	,000		,000	,000	,000	,000
	N	4214	4214	4214	4214	4214	4214
Idade Voluntários	Correlação de Pearson	-,088**	-,058**	1	,078**	-,115**	,044**
	Sig. (2 extremidades)	,000	,000		,000	,000	,004
	N	4214	4214	4214	4214	4214	4214
Altura	Correlação de Pearson	,259**	-,191**	,078**	1	,834**	-,104**
	Sig. (2 extremidades)	,000	,000	,000		,000	,000
	N	4214	4214	4214	4214	4214	4214
Peso	Correlação de Pearson	,372**	-,249**	-,115**	,834**	1	-,219**
	Sig. (2 extremidades)	,000	,000	,000	,000		,000
	N	4214	4214	4214	4214	4214	4214
BIS	Correlação de Pearson	-,736**	-,574**	,044**	-,104**	-,219**	1
	Sig. (2 extremidades)	,000	,000	,004	,000	,000	
	N	4214	4214	4214	4214	4214	4214

** . A correlação é significativa no nível 0,01 (2 extremidades).

Anexos B

Saídas do SPSS nos modelos de regressão (Stepwise Forward,Enter)

```
REGRESSION
/MISSING LISTWISE
/STATISTICS COEFF OUTS R ANOVA COLLIN TOL
/CRITERIA=PIN(.05) POUT(.10)
/NOORIGIN
/DEPENDENT BIS
/METHOD=ENTER propofol Remifentanil
/SCATTERPLOT=(*ZRESID ,*ZPRED) (*SRESID ,*ADJPRED)
/RESIDUALS DURBIN NORMPROB(ZRESID)
/SAVE ZPRED LEVER RESID SDRESID SDBETA SDFIT.
```

Variáveis Inseridas/Removidas^a

Modelo	Variáveis inseridas	Variáveis removidas	Método
1	Remifentanil Ce (ng/ml), Propofol Ce (ug/ml) ^b		Inserir

a. Variável Dependente: BIS

b. Todas as variáveis solicitadas inseridas.

Resumo do modelo^b

Modelo	R	R quadrado	R quadrado ajustado	Erro padrão da estimativa	Durbin-Watson
1	,836 ^a	,698	,698	8,748	,055

a. Preditores: (Constante), Remifentanil Ce (ng/ml), Propofol Ce (ug/ml)

b. Variável Dependente: BIS

ANOVA^a

Modelo		Soma dos Quadrados	df	Quadrado Médio	Z	Sig.
1	Regressão	746311,096	2	373155,548	4876,630	,000 ^b
	Resíduo	322222,106	4211	76,519		
	Total	1068533,202	4213			

a. Variável Dependente: BIS

b. Preditores: (Constante), Remifentanil Ce (ng/ml), Propofol Ce (ug/ml)

Coefficientes^a

Modelo		Coeficientes não padronizados		Coeficientes padronizados	t	Sig.	Estatísticas de colinearidade	
		B	Erro Padrão	Beta			Tolerância	VIF
1	(Constante)	102,190	,319		320,654	,000		
	Propofol Ce (ug/ml)	-8,741	,122	-,629	-71,762	,000	,932	1,072
	Remifentanil Ce (ng/ml)	-10,464	,223	-,411	-46,858	,000	,932	1,072

a. Variável Dependente: BIS

Diagnóstico de colinearidade^a

Modelo	Dimensão	Autovalor	Índice de condição	Proporções de variância		
				(Constante)	Propofol Ce (ug/ml)	Remifentanil Ce (ng/ml)
1	1	2,526	1,000	,02	,02	,06
	2	,378	2,585	,08	,06	,94
	3	,096	5,121	,90	,92	,00

a. Variável Dependente: BIS

Estatísticas de resíduos^a

	Mínimo	Máximo	Média	Desvio Padrão	N
Valor previsto	46,11	102,19	75,01	13,310	4214
Valor Previsto Padrão	-2,172	2,042	,000	1,000	4214
Erro padrão do valor previsto	,136	,394	,228	,049	4214
Valor previsto ajustado	46,08	102,20	75,01	13,310	4214
Resíduo	-25,911	25,827	,000	8,745	4214
Resíduo Padronizado	-2,962	2,953	,000	1,000	4214
Resíduos Estudantizados de Estud.	-2,963	2,955	,000	1,000	4214
	-25,929	25,870	,000	8,752	4214
Resíduos deletados Estudantizados	-2,966	2,958	,000	1,000	4214
Mahal. Distância	,019	7,559	2,000	1,380	4214
Distância de Cook	,000	,005	,000	,000	4214
Valor de ponto alavanca centralizado	,000	,002	,000	,000	4214

a. Variável Dependente: BIS

Teste de Kolmogorov-Smirnov de uma amostra

		Studentized Deleted Residual
N		4214
Parâmetros normais ^{a,b}	Média	-,0000095
	Erro Desvio	1,00033285
Diferenças Mais Extremas	Absoluto	,044
	Positivo	,017
	Negativo	-,044
Estatística de teste		,044
Significância Sig. (2 extremidades)		,000 ^c

a. A distribuição do teste é Normal.

b. Calculado dos dados.

c. Correção de Significância de Lilliefors.

Descrição do modelo

Nome do modelo	MOD_1
Série ou sequência	1
Transformação	Unstandardized Residual
Diferenciação não sazonal	Nenhum
Diferenciação sazonal	0
Comprimento do Período Sazonal	0
Padronização	Nenhuma periodicidade
Distribuição	Não aplicado
Tipo	Normal
Localização	estimado
Escala	estimado
Método de Estimação de Posição Fracionária	de Blom
Classificações atribuídas a empates	A classificação média de valores empatados

Aplicando as especificações de modelo de MOD_1

Resumo de processamento do caso

		Unstandardized Residual
Comprimento da série ou sequência		4214
Número de Valores Ausentes na	Usuário ausente	0
Parcela	Sistema ausente	0

Os casos não são ponderados.

Parâmetros de Distribuição Estimados

		Unstandardized Residual
Distribuição normal	Localização	,0000000
	Escala	8,74544547

Os casos não são ponderados.

Resumo de processamento do caso

	Casos					
	Válido		Ausente		Total	
	N	Porcentagem	N	Porcentagem	N	Porcentagem
Unstandardized Residual	4214	100,0%	0	0,0%	4214	100,0%

Descritivos

		Estatística	Erro Padrão	
Unstandardized Residual	Média	,0000000	,13472081	
	95% Intervalo de Confiança para Média	Limite inferior	-,2641238	
		Limite superior	,2641238	
	5% da média aparada	,0845946		
	Mediana	-,0013141		
	Variância	76,483		
	Desvio Padrão	8,74544547		
	Mínimo	-25,91119		
	Máximo	25,82729		
	Intervalo	51,73848		
	Intervalo interquartil	12,66041		
	Assimetria	-,137	,038	
	Curtose	-,386	,075	

Testes de Normalidade

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Estatística	df	Sig.	Estatística	df	Sig.
Unstandardized Residual	,044	4214	,000	,995	4214	,000

a. Correlação de Significância de Lilliefors

```
REGRESSION
/MISSING LISTWISE
/STATISTICS COEFF OUTS R ANOVA COLLIN TOL
/CRITERIA=PIN(.05) POUT(.10)
/NOORIGIN
/DEPENDENT BIS
/METHOD=STEPWISE propofol Remifentanil
/SCATTERPLOT=(*ZRESID ,*ZPRED) (*SDRESID ,*ADJPRED)
/RESIDUALS DURBIN
/SAVE ZPRED LEVER RESID SDRESID SDBETA SDFIT.
```

Variáveis Inseridas/Removidas^a

Modelo	Variáveis inseridas	Variáveis removidas	Método
1	Propofol Ce (ug/ml)		Em Etapas (Critérios: Probabilidade de F a ser inserido <= ,050, Probabilidade de F a ser removido >= ,100).
2	Remifentanil Ce (ng/ml)		Em Etapas (Critérios: Probabilidade de F a ser inserido <= ,050, Probabilidade de F a ser removido >= ,100).

a. Variável Dependente: BIS

Resumo do modelo^c

Modelo	R	R quadrado	R quadrado ajustado	Erro padrão da estimativa	Durbin-Watson
1	,736 ^a	,541	,541	10,78843	
2	,836 ^b	,698	,698	8,74752	,055

a. Preditores: (Constante), Propofol Ce (ug/ml)

b. Preditores: (Constante), Propofol Ce (ug/ml), Remifentanil Ce (ng/ml)

c. Variável Dependente: BIS

ANOVA^a

Modelo		Soma dos Quadrados	df	Quadrado Médio	Z	Sig.
1	Regressão	578297,807	1	578297,807	4968,614	,000 ^b
	Resíduo	490235,395	4212	116,390		
	Total	1068533,202	4213			
2	Regressão	746311,096	2	373155,548	4876,630	,000 ^c
	Resíduo	322222,106	4211	76,519		
	Total	1068533,202	4213			

a. Variável Dependente: BIS

b. Preditores: (Constante), Propofol Ce (ug/ml)

c. Preditores: (Constante), Propofol Ce (ug/ml), Remifentanil Ce (ng/ml)

Coefficientes^a

Modelo		Coeficientes não padronizados		Coeficientes padronizados	t	Sig.	Estatísticas de colinearidade	
		B	Erro Padrão	Beta			Tolerância	VIF
1	(Constante)	99,687	,387		257,265	,000		
	Propofol Ce (ug/ml)	-10,225	,145	-,736	-70,488	,000	1,000	1,000
2	(Constante)	102,190	,319		320,654	,000		
	Propofol Ce (ug/ml)	-8,741	,122	-,629	-71,762	,000	,932	1,072
	Remifentanil Ce (ng/ml)	-10,464	,223	-,411	-46,858	,000	,932	1,072

a. Variável Dependente: BIS

Variáveis excluídas^a

Modelo		Beta In	t	Sig.	Correlação parcial	Estatísticas de colinearidade		
						1,072	,932	Tolerância mínima
1	Remifentanil Ce (ng/ml)	-,411 ^b	-46,858	,000	a. Variável Dependente: BIS	,932		

b. Preditores no Modelo: (Constante), Propofol Ce (ug/ml)

Diagnóstico de colinearidade^a

Modelo	Dimensão	Autovalor	Índice de condição	Proporções de variância		
				(Constante)	Propofol Ce (ug/ml)	Remifentanil Ce (ng/ml)
1	1	1,903	1,000	,05	,05	
	2	,097	4,438	,95	,95	
2	1	2,526	1,000	,02	,02	,06
	2	,378	2,585	,08	,06	,94
	3	,096	5,121	,90	,92	,00

a. Variável Dependente: BIS

Estatísticas de resíduos^a

	Mínimo	Máximo	Média	Desvio Padrão	N

Valor previsto	46,1099	102,1899	75,0138	13,30958	4214
Valor Previsto Padrão	-2,172	2,042	,000	1,000	4214
Erro padrão do valor previsto	,136	,394	,228	,049	4214
Valor previsto ajustado	46,0837	102,2048	75,0136	13,31009	4214
Resíduo	-25,91120	25,82729	,00000	8,74545	4214
Resíduo Padronizado	-2,962	2,953	,000	1,000	4214
Resíduos Estudantizados	-2,963	2,955	,000	1,000	4214
de Estud.	-25,92872	25,86970	,00012	8,75207	4214
Resíduos deletados					
Estudantizados	-2,966	2,958	,000	1,000	4214
Mahal. Distância	,019	7,559	2,000	1,380	4214
Distância de Cook	,000	,005	,000	,000	4214
Valor de ponto alavanca					
centralizado	,000	,002	,000	,000	4214

a. Variável Dependente: BIS

Anexos C

Análise Clusters

Estatísticas descritivas

	N	Variância
Zscore: Propofol Ce (ug/ml)	4214	1,000
Zscore: Remifentanil Ce (ng/ml)	4214	1,000
Zscore: BIS	4214	1,000
N válido (de lista)	4214	

Centros de cluster iniciais

	Cluster		
	1	2	3
Zscore: Propofol Ce (ug/ml)	1,83046	-2,10598	-,68978
Zscore: Remifentanil Ce (ng/ml)	-,70330	-,93016	2,26819
Zscore: BIS	-2,26136	1,44334	-,37761

Histórico de iteração^a

Iteração	Alteração em centros de cluster		
	1	2	3
1	1,519	1,520	1,300
2	,112	,058	,038
3	,097	,043	,004
4	,120	,048	,002
5	,088	,031	,005
6	,044	,015	,003
7	,034	,012	,003
8	,026	,009	,002
9	,016	,004	,003
10	,012	,004	,001

a. As iterações foram interrompidas porque o número máximo de iterações foi atingido. As iterações não convergiram. A mudança de coordenada absoluta máxima para qualquer centro é ,010. A iteração atual é 10. A distância mínima entre os centros iniciais é 3,943.

Centros de cluster finais

	Cluster		
	1	2	3
Zscore: Propofol Ce (ug/ml)	1,43133	-,52326	,03076
Zscore: Remifentanil Ce (ng/ml)	,15948	-,76236	1,21372
Zscore: BIS	-1,38865	,74387	-,43590

Distâncias entre centros de cluster finais

Cluster	1	2	3
1		3,036	1,995
2	3,036		2,367
3	1,995	2,367	

ANOVA

	Cluster		Erro		Z	Sig.
	Quadrado Médio	df	Quadrado Médio	df		
Zscore: Propofol Ce (ug/ml)	1086,504	2	,484	4211	2242,787	,000
Zscore: Remifentanil Ce (ng/ml)	1575,918	2	,252	4211	6253,690	,000
Zscore: BIS	1464,409	2	,305	4211	4801,986	,000

Os testes F devem ser usados apenas para finalidades descritivas porque os cluster foram escolhidos para maximizar as diferenças entre os casos em clusters diferentes. Os níveis de significância observados não estão corrigidos para isso e, dessa forma, não podem ser interpretados como testes da hipótese de que as médias de cluster são iguais.

Número de casos em cada cluster

Cluster	1	769,000
	2	2178,000
	3	1267,000
Válido		4214,000
Ausente		,000

Resumo de processamento de caso de análise

Casos não ponderados		N	Porcentagem
	Válido	4214	100,0
Excluídos	Códigos de grupo ausentes ou fora do intervalo	0	,0
	Pelo menos uma variável discriminante ausente	0	,0
	Códigos de grupo ausentes ou fora do intervalo e pelo menos uma variável discriminadora ausentemos	0	,0
	Total	0	,0
	Total	4214	100,0

Estatísticas de grupo

Número de caso de cluster		Média	Desvio Padrão	N válido (de lista)	
				Não ponderado	Ponderado
1	Zscore: Score: Propofol Ce (ug/ml)	1,4313321	,33880442	769	769,000
	Zscore: Zscore: Remifentanil Ce (ng/ml)	,1594792	,72977481	769	769,000
	Zscore: Zscore: BIS	-1,3886490	,75799680	769	769,000
2	Zscore: Zscore: Propofol Ce (ug/ml)	-,5232647	,81679124	2178	2178,000
	Zscore: Zscore: Remifentanil Ce (ng/ml)	-,7623609	,33662125	2178	2178,000
	Zscore: Zscore: BIS	,7438709	,52062468	2178	2178,000
3	Zscore: Zscore: Propofol Ce (ug/ml)	,0307626	,62810286	1267	1267,000
	Zscore: Zscore: Remifentanil Ce (ng/ml)	1,2137195	,56592620	1267	1267,000
	Zscore: Zscore: BIS	-,4358957	,44689993	1267	1267,000
Total	Zscore: Zscore: Propofol Ce (ug/ml)	,0000000	1,00000000	4214	4214,000
	Zscore: Zscore: Remifentanil Ce (ng/ml)	,0000000	1,00000000	4214	4214,000
	Zscore: Zscore: BIS	,0000000	1,00000000	4214	4214,000

Testes de igualdade de médias de grupo

Lambda de Wilks	Z	df1	df2	Sig.
,484	2242,787	2	4211	,000
,252	6253,690	2	4211	,000
,305	4801,986	2	4211	,000

Matrizes dentro de grupos em pool^a

	Zscore: Zscore: Propofol Ce (ug/ml)	Zscore: Zscore: Remifentanil Ce (ng/ml)
Covariância		
Zscore: Zscore: Propofol Ce (ug/ml)	,484	,001
Zscore: Zscore: Remifentanil Ce (ng/ml)	,001	,252
Zscore: Zscore: BIS	-,168	-,081
Correlação		
Zscore: Zscore: Propofol Ce (ug/ml)	1,000	,002

Zscore: Zscore: Remifentanil Ce (ng/ml)	,002	1,000
Zscore: Zscore: BIS	-,436	-,294

a. A matriz de covariâncias possui 4211 graus de liberdade.

Matrizes de covariâncias^a

Número de caso de cluster	Zscore: Zscore: Propofol Ce (ug/ml)	Zscore: Zscore: Remifentanil Ce (ng/ml)	Zscore: Zscore: BIS
1	Zscore: Zscore: Propofol Ce (ug/ml)	,115	,120
	Zscore: Zscore: Remifentanil Ce (ng/ml)	,120	,533
	Zscore: Zscore: BIS	-,112	-,223
2	Zscore: Zscore: Propofol Ce (ug/ml)	,667	-,078
	Zscore: Zscore: Remifentanil Ce (ng/ml)	-,078	,113
	Zscore: Zscore: BIS	-,224	-,047
3	Zscore: Zscore: Propofol Ce (ug/ml)	,395	,064
	Zscore: Zscore: Remifentanil Ce (ng/ml)	,064	,320
	Zscore: Zscore: BIS	-,104	-,055
Total	Zscore: Zscore: Propofol Ce (ug/ml)	1,000	,260
	Zscore: Zscore: Remifentanil Ce (ng/ml)	,260	1,000
	Zscore: Zscore: BIS	-,736	-,574

a. A matriz de covariâncias total possui 4213 graus de liberdade.

Determinantes de log

Número de caso de cluster	Posição	Determinante de log
1	3	-3,891
2	3	-4,602
3	3	-3,889
dentro de grupos em pool	3	-3,614

As posições e os logaritmos naturais de determinantes impressos são aqueles das matrizes de covariâncias de grupo.

Resultados do teste

M de Box		2712,416
Z	Aprox.	225,751
	df1	12
	df2	28525920,722
	Sig.	,000

Testa hipótese nula de matrizes de covariâncias de população igual.

Variáveis Inseridas/Removidas^{a,b,c,d}

Etapa	Inseridas	Lambda de Wilks							
		Estatística	df1	df2	df3	F exato			
						Estatística	df1	df2	Sig.
1	Zscore: Zscore: Remifentanil Ce (ng/ml)	,252	1	2	4211,000	6253,690	2	4211,000	,000
2	Zscore: Zscore: BIS	,105	2	2	4211,000	4402,500	4	8420,000	,000
3	Zscore: Zscore: Propofol Ce (ug/ml)	,096	3	2	4211,000	3132,048	6	8418,000	,000

Em cada etapa, a variável que minimiza o Lambda de Wilks geral é inserida.

- O número máximo de etapas é 6.
- A significância máxima do F a ser inserido é .05.
- A significância mínima do F a ser removido é .10.
- Nível f, tolerância ou VIN insuficiente para cálculos adicionais.

Variáveis na análise

Etapa		Tolerância	Sig. de F a ser removida	Lambda de Wilks
1	Zscore: Zscore: Remifentanil Ce (ng/ml)	1,000	,000	
2	Zscore: Zscore: Remifentanil Ce (ng/ml)	,914	,000	,305
	Zscore: Zscore: BIS	,914	,000	,252
3	Zscore: Zscore: Remifentanil Ce (ng/ml)	,894	,000	,260
	Zscore: Zscore: BIS	,724	,000	,131
	Zscore: Zscore: Propofol Ce (ug/ml)	,792	,000	,105

Variáveis não presentes na análise

Etapa		Tolerância	Mín. Tolerância	Sig. de F a ser inserida	Lambda de Wilks
0	Zscore: Zscore: Propofol Ce (ug/ml)	1,000	1,000	,000	,484
	Zscore: Zscore: Remifentanil Ce (ng/ml)	1,000	1,000	,000	,252
	Zscore: Zscore: BIS	1,000	1,000	,000	,305
1	Zscore: Zscore: Propofol Ce (ug/ml)	1,000	1,000	,000	,131
	Zscore: Zscore: BIS	,914	,914	,000	,105
2	Zscore: Zscore: Propofol Ce (ug/ml)	,792	,724	,000	,096

Lambda de Wilks

Etapa	Número de variáveis	Lambda	df1	df2	df3	F exato			
						Estatística	df1	df2	Sig.
1	1	,252	1	2	4211	6253,690	2	4211,000	,000
2	2	,105	2	2	4211	4402,500	4	8420,000	,000
3	3	,096	3	2	4211	3132,048	6	8418,000	,000

Valores próprios

Função	Autovalor	% de variância	% cumulativa	Correlação canônica
1	3,571 ^a	73,5	73,5	,884
2	1,286 ^a	26,5	100,0	,750

a. As primeiras 2 funções discriminantes canônicas foram usadas na análise.

Lambda de Wilks

Teste de funções	Lambda de Wilks	Qui-quadrado	df	Sig.
1 até 2	,096	9878,537	6	,000
2	,437	3480,554	2	,000

Coefficientes de funções discriminantes canônicas

padronizados

	Função	
	1	2
Zscore: Zscore: Propofol Ce (ug/ml)	,203	-,366

Zscore: Zscore: Remifentanil Ce (ng/ml)	,739	,706
Zscore: Zscore: BIS	-,405	,657

Matriz de estruturas

	Função	
	1	2
Zscore: Zscore: Remifentanil Ce (ng/ml)	,859*	,513
Zscore: Zscore: BIS	-,711*	,609
Zscore: Zscore: Propofol Ce (ug/ml)	,382	-,651*

Funções em centroides de grupo

Número de caso de cluster	Função	
	1	2
1	1,671	-2,180
2	-1,821	,087
3	2,115	1,174

Funções discriminantes canônicas não padronizadas
avaliadas em médias de grupo

Resumo de processamento de classificação

Processado	4214
Excluídos	
Códigos de grupo ausentes ou fora do intervalo	0
Pelo menos uma variável discriminante ausente	0
Usado em saída	4214

Probabilidades a priori para grupos

Número de caso de cluster	A priori	Casos utilizados na análise	
		Não ponderado	Ponderado
1	,182	769	769,000
2	,517	2178	2178,000
3	,301	1267	1267,000
Total	1,000	4214	4214,000

Coeficientes de função de classificação

	Número de caso de cluster		
	1	2	3
Zscore: Zscore: Propofol Ce (ug/ml)	1,634	-,578	,001
Zscore: Zscore: Remifentanil Ce (ng/ml)	-,607	-2,558	4,766
Zscore: Zscore: BIS	-3,817	1,438	-,155
(Constante)	-5,473	-2,321	-4,128

Funções discriminantes lineares de Fisher

Resultados da classificação^a

		Número de caso de cluster	Associação ao grupo prevista			Total
			1	2	3	
Original	Contagem	1	621	83	65	769
		2	0	2154	24	2178
		3	2	13	1252	1267
%		1	80,8	10,8	8,5	100,0
		2	,0	98,9	1,1	100,0
		3	,2	1,0	98,8	100,0

a. 95,6% de casos agrupados originais classificados corretamente.

