



**Tiago André dos  
Santos Silva**

## **Relatório de Estágio em Bioestatística**

### ***Internship Report in Biostatistics***

Relatório apresentado à Universidade de Aveiro para cumprimento dos requisitos necessários à obtenção do grau de Mestre em Biomedicina Farmacêutica, realizado sob a orientação científica do Professor Doutor Pedro Manuel Vargues de Aguiar, Consultor de Epidemiologia e Estatística da Eurotrials, Consultores Científicos e Professor Auxiliar na Escola Nacional de Saúde Pública, da Universidade Nova de Lisboa e da Professora Doutora Alexandra Isabel Cardador de Queirós, Professora Coordenadora da Escola Superior de Saúde da Universidade de Aveiro.



## **jury**

presiding juror

**Prof. Dr. Bruno Miguel Alves Fernandes do Gago**

invited assistant Professor at the Health Sciences Department at University of Aveiro

**Prof. Dr. Vera Mónica Almeida Afreixo**

auxiliary Professor at the Mathematics Department at University of Aveiro

**Prof. Dr. Alexandra Isabel Cardador de Queirós**

coordinator Professor at Higher School of Health of the University of Aveiro

**Prof. Dr. Pedro Manuel Vargues de Aguiar**

assistant Professor at the Institute of Hygiene and Tropical Medicine at New University of Lisbon



## **acknowledgements**

I would like to appreciate the support of my Supervisor, Professor Pedro Aguiar, as well as my Co-Supervisor, Professor Alexandra Queirós, for being fully available and offering valuable suggestions, criticism and advice that were absolutely essential for the development of this report.

I also would like to thank Professor Luís Almeida and Professor Bruno Gago, not only for suggesting me this challenging activity and arranging the terms of my internship, but also for being available and for supporting me throughout this process.

My thanks also go to all Eurotrials collaborators and Administration Board, notably to Dr. Maria João Queiroz and Dr. Inês Costa, members of the Administration Board, for trusting and providing me with the opportunity of working in this wonderful Company, as well as to Dr. Maria João Salgado, for her tireless support and guidance throughout my internship, trusting me with many different activities in various departments of the Company.

Many thanks to the Biostatistics department, namely to Catarina Silva, Filipa Negreiro and Vera Vicente, whose support, availability, trust and friendship were extremely motivating.

I also thank Ana Filipa Bernardo, Catarina Alves, Catarina Ferreira, Isabel Pinto, Luís Veloso and Pedro Noronha for trusting and sharing with me several activities in their respective departments, always providing valuable guidance and criticism throughout this process.

Finally, I would like to thank my course colleagues, with whom I had the pleasure of working during my academic course and also to my parents, who always fully supported me and are my main source of motivation.



**palavras-chave**

Eurotrials, Contract Research Organisation, bioestatística, análise estatística, investigação clínica.

**resumo**

Este relatório descreve a minha experiência de 9 meses enquanto estagiário na Eurotrials, Consultores Científicos, uma Empresa especializada em investigação clínica e consultoria científica.

Este estágio desenrolou-se em duas vertentes: formação multidisciplinar e monodisciplinar. A formação multidisciplinar envolveu alguma forma de participação activa em diferentes departamentos desta Empresa, com o objectivo de obter uma perspectiva alargada do processo multidisciplinar inerente ao desenvolvimento clínico de produtos de saúde.

A formação monodisciplinar concentrou-se na área de estatística médica, sendo realizada no departamento de Bioestatística da Empresa, com o objectivo de obter conhecimentos práticos de aplicação da estatística à investigação em saúde, implicando também a interiorização de conceitos estatísticos fundamentais.

Este estágio permitiu-me compreender de forma mais aprofundada o trabalho multidisciplinar necessário para a realização adequada de um projecto de investigação clínica. Permitiu-me também não só adquirir conhecimentos importantes de análise estatística, mas também compreender, de forma mais clara, o papel da estatística na investigação clínica, como ferramenta essencial no planeamento do estudo, análise e interpretação dos dados obtidos.



**keywords**

Eurotrials, Contract Research Organisation, biostatistics, statistical analysis, clinical research.

**abstract**

This report describes my experience of 9 months as an intern at Eurotrials, Scientific Consultants, a company devoted to clinical research and scientific consulting.

This internship developed in two aspects: multidisciplinary and monodisciplinary training. Multidisciplinary training involved active participation in different departments of this Company, with the objective of obtaining a broad perspective on the multidisciplinary process of the clinical development of medical products.

Monodisciplinary training was focused in medical statistics, being carried out in the Biostatistics department of the Company. The objective was to obtain practical knowledge for the application of statistics in health sciences. This implied the learning of fundamental statistical concepts.

This internship allowed me to understand, in depth, the multidisciplinary work necessary for an adequate performance of a clinical research project. It also allowed me to acquire valuable knowledge in statistical analysis, as well as to clearly understand the role of statistics in clinical research, as an essential tool in study planning, analysis and interpretation of data obtained.



## Table of Contents

|  |    |
|--|----|
| Introduction .....   | 3  |
| 1. Overview of the host Company .....                                      | 5  |
| 1.1. Overview of Contract Research Organisations (CRO) .....               | 5  |
| 1.2. Overview of the Host Company.....                                     | 5  |
| 1.3. Overview of the Biostatistics department .....                        | 8  |
| 2. Multidisciplinary activity.....   | 11 |
| 2.1. Mandatory actions.....  | 11 |
| 2.2. Epidemiology & Late Phase Research (Medical writing activities) ..... | 12 |
| 2.3. Research & Development .....  | 16 |
| 2.4. Data Management .....   | 18 |
| 2.5. Teaching & Training.....  | 20 |
| 2.6. Clinical Trials .....   | 22 |
| 3. Monodisciplinary activities within the Biostatistics department .....   | 25 |
| 3.1. Use of statistics in health research .....                            | 25 |
| 3.2. Overview of the statistical methods performed.....                    | 28 |
| 3.3. Work developed at Eurotrials .....                                    | 44 |
| 4. Discussion and conclusions .....  | 51 |
| 4.1. Multidisciplinary Activity.....                                       | 51 |
| 4.2. Monodisciplinary activity: Biostatistics department.....              | 53 |
| References .....   | 55 |
| Appendix 1 .....   | 57 |

## Tables and Figures

|  |    |
|--|----|
| Table 1. Global Evaluation of the training action (n=8).....                     | 21 |
| Table 2. Association between smoking habits and diagnosis of lung cancer .....   | 34 |
| Table 3. Values and respective ranks for two groups A and B.....                 | 36 |
| Table 4. Symptoms before and after treatment.....                                | 38 |
| Table 5. Distribution of values and respective ranks for groups A, B and C ..... | 40 |
| Table A 1. P-values for the Chi-Square distribution .....                        | 57 |
| Table A 2. Student's t distribution.....   | 57 |
| Table A 3. Critical Values for U.....  | 58 |
| Table A 4. Critical values for the Wilcoxon signed rank test.....                | 58 |
| Table A 5. Critical values for the sign test .....                               | 59 |

|   |    |
|---|----|
| Figure 1. Bar Chart.....  | 29 |
| Figure 2. Pie Chart .....   | 29 |
| Figure 3. Formula for variance.....   | 29 |
| Figure 4. Histogram .....   | 30 |
| Figure 5. Box Plot.....   | 30 |
| Figure 6. Normal Distribution.....  | 31 |
| Figure 7. Formula for z-score.....  | 31 |
| Figure 8. formula for confidence interval .....                                       | 32 |
| Figure 9. Formula for standard error of the difference between two means .....        | 33 |
| Figure 10. Formula for standard error of the difference between two proportions ..... | 33 |
| Figure 11. Formula for expected frequencies.....                                      | 34 |
| Figure 12. Formula for Chi-Square statistic.....                                      | 35 |
| Figure 13. Calculation of two sample t-test statistic .....                           | 36 |
| Figure 14 Example of calculation of U index .....                                     | 37 |
| Figure 15. Calculation of Paired t-test .....   | 37 |
| Figure 16. Calculation of McNemar Statistic.....                                      | 39 |
| Figure 17. Calculation of Kruskal-Wallis statistic (H).....                           | 39 |
| Figure 18. Example of calculation of Kruskal-Wallis statistic.....                    | 40 |
| Figure 19. Formula for Pearson's correlation coefficient .....                        | 40 |
| Figure 20. Formula for Odds.....  | 41 |
| Figure 21. Model of logistic multiple regression .....                                | 41 |
| Figure 22. Calculation of LOGIT Function .....  | 42 |
| Figure 23. Example of application of logistic multiple regression model.....          | 42 |
| Figure 24. Example of calculation of Odds Ratios .....                                | 42 |
| Figure 25. Cross multiplication method.....   | 49 |
| Figure 26. Timeline of internship activities .....                                    | 51 |
| Figure 27. Statistical tests performed.....   | 53 |

## Introduction

From September 2010 to June 2011, I enrolled in an internship within the scope of the Master's degree in Pharmaceutical Biomedicine. This internship occurred at Eurotrials, Scientific Consultants, a Company dedicated to clinical research and scientific consulting services. This internship had two main objectives:

- To gather basic knowledge in various areas relevant to clinical research, enrolling in activities from different departments, in order to understand the multidisciplinary framework needed to properly plan, conduct, manage and report a clinical study.
- To understand in depth the application of statistics in health sciences, working mostly within the Biostatistics department.

This internship report describes my working experience during these 9 months, characterising Eurotrials and defining its role inside the clinical research environment and also reporting all activities performed and lessons learned from such activities. For this purpose, it is divided in 4 chapters, defined as follows:

- Overview of the host Company: this chapter describes Eurotrials, defining where it fits in the clinical research framework, its purpose, organisation and work developed. The Biostatistics department is described in more detail, as it was the main work area during this internship.
- Multidisciplinary activity: this chapter reports the activities developed in several departments of Eurotrials (with the exception of the Biostatistics department) carried out to understand the objectives and comprehend the type of work performed in each of these departments, as well as how they fit in the Company framework.
- Monodisciplinary activity within the Biostatistics department: this chapter describes in depth the role of statistics in clinical research and justifies its importance in this context. It also gives a theoretical overview of the statistical methods I carried out during my internship, as well as all the activities in which I participated.
- Discussion and conclusions: this chapter gives an overview of my internship experience, discussing what was learned during this period for each department where I actively participated and also summarising the importance of each of these departments to the successful planning, conducting, management and reporting of clinical study projects. For the Biostatistics department, these points are discussed in more depth.



## **1. Overview of the host Company**

This chapter describes the host Company, its purpose, activities and where it belongs in the clinical research framework.

### **1.1. Overview of Contract Research Organisations (CRO)**

Pharmaceutical Industry has been experiencing significant structural changes (1). Over the last 30 years, a growing trend towards outsourcing of several services has been observed, including development of medical products (1). CROs are increasingly assuming responsibilities in this area, particularly in the Biotechnology sector, where outsourcing has increased dramatically (2). These companies are scientific organisations (commercial or academic), to which a sponsor may transfer responsibility for some of its tasks or obligations (3).

There is a substantial growth in this business. Since 1994, the Pharmaceutical Industry has eliminated more than 40 000 jobs, many in Research & Development (R&D) (1). It is estimated that more than 60% of all clinical studies now involve significant outsourcing (1). Also, spending on CRO services, as a share of total global development spending, rose steadily from 13,7% in 2001 to 14,8% in 2004 (4).

CROs can complete drug development tasks faster than the sponsors, without compromising data quality, even with large trials, involving multiple study sites (2). This is a significant advantage, financially speaking, as taking a month off development time may result in an additional \$40 million income (approximately €27,6 million) (2).

Many services can be performed by CROs, including: investigator recruiting and training, study monitoring, data management, statistical analysis, auditing activities, adverse events reporting, medical writing or regulatory services (5).

### **1.2. Overview of the Host Company**

My internship occurred at a Portuguese CRO, called Eurotrials, Scientific Consultants. This is a private Company founded in Lisbon, in 1995, by members from different backgrounds in the Academia, Medical Community and Pharmaceutical Industry (6). Being a CRO, it provides outsourced pharmaceutical research services for Pharmaceutical and Biopharmaceutical Industries, as well as consulting and training services, in the field of clinical research (6).

With expertise in clinical research and scientific consulting in Health Sciences, it operates in Europe and Latin America, with projects developing in Africa (6). Its partners include Pharmaceutical and Biotechnology Industries, CROs, Regulatory Agencies, Food Industry, Academia and Clinical Research Centres (6).

Eurotrials offers several services related to clinical research, by means of a contract with a sponsor. Services include (6):

## Internship Report in Biostatistics

### ***a) Research & Development***

This department is responsible for the early planning of the clinical study. Tasks involved include meetings with Academia representatives, pharmaceutical companies and other CROs, in order to develop business connections, new research projects and search for potential sponsorships, as well as strategic and regulatory planning of clinical studies, along with the clinical trials and Regulatory Affairs departments.

### ***b) Clinical Trials***

The biggest sector in this Company, the Clinical Trials department is concerned with clinical trials design, development of Case Report Forms (CRF) (along with the Data Management department), site selection for the trial, study monitoring activities (such as adequate training of the research team or assurance of Good Clinical Practice (GCP) and regulatory compliance), project management activities (such as time and resources management for the research project) and clinical report development. Medical writing tasks are also performed within the scope of clinical trials. The medical writer can assist in writing clinical protocols, clinical reports and help prepare the clinical investigational brochure. Scientific publications may also be prepared and/or submitted by the medical writer.

### ***c) Epidemiology & Late Phase Research***

This department works in planning, designing, implementing, monitoring and managing clinical observational studies, which are clinical studies where the investigator does not interfere (7), as well as other post-marketing studies. Medical writing is also an important part of this department, assisting in the writing of observational study protocols and reports, as well as scientific publications.

### ***d) Data Management***

The Data Management department works to assure proper data collection and preparation for statistical analysis. Some of its tasks include development of the CRF, database development and validation, data validation and quality control, data collection and export to the desired format and development of data management reports.

### ***e) Biostatistics***

The Biostatistics department collaborates in study design and definition of statistical methodology, definition of sample size, development of randomisation envelopes, statistical analysis (including interim analysis), development of statistical reports, results presentation in meetings or conferences/symposia, article submission and statistical consulting and/or training for other companies/professionals.

### *f) Regulatory Affairs*

This department follows the medical product throughout its entire life cycle, from a regulatory perspective. Some tasks include regulatory consulting for medical products, marketing authorisation applications submission, contact with Regulatory Authorities, validation of drug related information, readability studies for patient leaflets, price and reimbursement requests for medical products, clinical trial authorisation requests and import/export of experimental products and study material.

### *g) Pharmacovigilance*

The Pharmacovigilance department works in the risk/benefit assessment of medical products, and it is part of a pharmacovigilance and adverse events notification network, being responsible for notification of any adverse events reported to the Regulatory Authorities. Some tasks include consulting services, receipt, review, validation and notification of adverse events to the sponsor and Authorities, bibliography research for safety information, development and submission of Periodic Safety Update Reports and periodic internal training in pharmacovigilance, required for all collaborators.

### *h) Pharmacoeconomics*

Pharmacoeconomics is defined as the field of study that assesses the behaviour/welfare of individuals, companies and markets relevant to the use of pharmaceutical medicines, services or programs, focusing on the cost and consequences of such use (8). This department follows the latest regulatory requirements for economical evaluation studies of medical products. Therefore it is responsible for design, implementation and analysis of pharmacoeconomic studies, but also acts in health economic studies, such as the assessment of the economic impact of a certain disease.

### *i) Quality*

The Quality department assures that all departments operate in accordance with the ISO 9001:2008 Quality Standard, as well as the GCP Standard. It serves internally as a quality consulting centre, as well as an important aid to the development of Standard Operation Procedures (SOP) for all departments. It also helps preparing the Company for audits performed by a sponsor, as well as inspections that might be carried out.

For external services, the Quality department assists in the quality control of study documentation and data management. For quality assurance, it helps in the preparation and performance of audits inside the Company (internal audits), as well as to study sites (external audits).

### *j) Teaching & Training*

This department is responsible for the development, management and disclosure of training courses, whether internal or external. Eurotrials conducts many training sessions for different audiences, such as health professionals in hospital or industry settings, as well as students, professors and investigators working in the Academia.

More information about the Company and its work is available on its website ([www.eurotrials.com](http://www.eurotrials.com)) (6).

### **1.3. Overview of the Biostatistics department**

My internship occurred mainly in the Biostatistics department. Statistics is essential in the R&D of any medical product, being present in all stages of its development (9). It helps in the understanding of a possible causal relation between the product and a certain outcome, as well as its strength, driving the interest of the researcher to the most relevant significant results (7).

At Eurotrials, statistics acts in different phases of the clinical study. The Biostatistics department gives its input early, in the design phase of the study, developing the statistical methodology, as well as giving input in endpoint definition, key study variables definition and validation of the CRF. It is also responsible for the development of randomization lists and envelopes, when applicable. Development of a statistical analysis plan (SAP) for each study, stating in detail all the statistical activities planned for this study, is also ensured. Statistical analysis must follow this plan and deviations must be justified. The SAP is written based on the planned study design.

After the analysis is concluded, the statistician develops the statistical analysis report, where the results are stated and significant results are highlighted. This department also works in scientific articles writing and submission, working together with a medical writer. The results of the study may be presented in meetings/conferences/symposia or other events. The preparation of these presentations may also be done by the statistician.

The Biostatistics department also plays a strong role in training activities. It is responsible for the preparation and lecturing of several training activities intended for health professionals working in the hospital or industry environment. Several courses are now regularly planned, varying in complexity of the course subject(s). This way, courses are adapted to different audiences, with different backgrounds and objectives.

Some examples of courses lectured (10):

- Statistical interpretation of scientific publications.
- Biostatistics in SPSS®.
- Multiple logistic regression analysis.
- Survival analysis.

## Internship Report in Biostatistics

The contents of these courses are based in the experience acquired within the Company, but always in accordance with proper scientific bibliography within statistical analysis applied to healthcare research.

Private statistical consulting sessions (*Consultório de Estatística*) are also carried out, where the statistician helps the investigator planning the statistical methodology for his/her project, as a consultant. In these sessions, the investigator states the objective of the study and the planned design. The statistician role may vary, from helping solve some questions, to assisting the elaboration of the entire statistical analysis methodology. Many people from different backgrounds use this service, from investigators designing their own study (investigator-driven studies) to academic researchers and PhD students working on research theses.

This department is also responsible for the development of two Eurotrials periodic publications:

- Bulletin *Saúde em Mapas e Números* (Health in Maps and Numbers). Each issue focuses on one subject (usually a disease) and respective epidemiological data is presented separately for Portugal, Europe and World. This may serve as an informative tool, giving updated information on disease trends all over the World, as well as other epidemiologic relevant factors, such as burden of disease, or, occasionally, economic trends towards the disease. During my internship at Eurotrials, I assisted in the writing and releasing of several bulletins. One example is the No 32, concerning sleep disorders (11).
- Publication *Cartas do Amigo Gauss* (Letters from Buddy Gauss). Each issue focus on one particular subject in statistics for clinical research. The objective is to explain the statistical fundamentals in a clear way, understandable for people not specialised in this field, being a valuable tool for anyone working in Clinical Research with no background in statistics, or in need to review/update their knowledge.

These publications may be useful for informational purposes, but are also important marketing tools, as they reveal the expertise of the Company in the epidemiology/statistics fields of clinical research. They are free for download at the Eurotrials website (6).

Eurotrials also promoted the writing of a statistical analysis book entitled, *Guia Prático Climepsi de Estatística em Investigação Epidemiológica: SPSS* (Climepsi Practical Guide on Statistics in Epidemiologic Research), written by Eurotrials Epidemiology and Statistics consultant, Prof. Pedro Aguiar. This book serves as an introductory text to the Epidemiology rationale and statistical application in healthcare research. Statistical techniques are explained in an accessible way, without resorting to complex mathematics, therefore being accessible for people with no statistics/mathematics background. It also contains a guide to the statistical analysis software tool IBM SPSS Statistics® (SPSS®), describing a suggested pathway to perform the analysis, for each analysis included (7).



## **2. Multidisciplinary activity**

One of the objectives of my internship was to acquire multidisciplinary knowledge and experience. This called for short visits and/or participation in research projects in other sections of this Company.

The purpose of this activity was to get a broad perspective on the multidisciplinary process of clinical development of medical products and health research activities, as well as to understand how different areas interact within the Company.

During my internship, I participated on various activities outside the Biostatistics department. Some of these activities are mandatory for all employees, such as internal training sessions. Others were performed solely on the scope of this internship. I performed several activities within the following sectors:

- Epidemiology & Late Phase Research (performing medical writing activities).
- Research & Development.
- Data Management.
- Teaching & Training.
- Clinical Trials.

During this chapter, each of these areas will be discussed in depth, separately. An overview of all mandatory training actions will also be performed.

### **2.1. Mandatory actions**

At the starting point of my internship, I needed to acknowledge the basic procedures of the Company, before enrolling in any project related activity. This was done through self-reading of the Company Quality Manual, internal SOPs and internal training sessions organised by the Quality department. These procedures included rules of conduct, proper handling of various documents and forms, introduction to information technology services, among others.

I also attended an internal pharmacovigilance training session, mandatory for all collaborators, on how to properly handle safety information. Aspects addressed included definition of safety information, proper handling of relevant safety information and procedures to adequately report this information to the pharmacovigilance manager. It also included internal procedures for forwarding or reporting of e-mails, letters or phone calls concerning or containing any kind of safety information that must be notified to the pharmacovigilance manager of the Company.

## **2.2. Epidemiology & Late Phase Research (Medical writing activities)**

Medical writing involves the development of scientific documentation by specialised professionals, called medical writers. This term has a broad definition, as medical writers can work in a large variety of areas as well as many different sponsors, including medical doctors, academic institutions, health science researchers and Pharmaceutical Industry. Medical writers working for the Pharmaceutical Industry may write regulatory documents required to obtain product marketing approval (12). These documents include the following:

- Investigator's brochure: Clinical and nonclinical data of the investigational product relevant to the study of the product in human subjects. Investigators should be familiar with this document to understand the rationale of the study they are involved in. The investigator's brochure also provides information on the clinical management of the subjects during the course of the study (13).
- Synopsis: A brief summary of the study protocol, containing the study rationale, main objectives, endpoints, overall study design and main statistical considerations.
- Protocol: This document describes the objective(s), design, methodology, statistical considerations and organisation of the study. It may also provide the background and rationale of the study (13).
- Report: Document incorporating a description of the clinical and statistical analysis results, incorporating tables and figures. In clinical trials, study reports require the presence of several elements, such as: sample CRF, investigator related information, information related to the experimental product(s), technical statistical documentation, related publications, patient data listings, and technical statistical details such as derivations, computations, analyses, and computer outputs required for data traceability purposes (13) (14).

The medical writer interacts with several departments, such as Biostatistics, Clinical Trials and Epidemiology & Late Phase Research. Medical writing activities performed in Eurotrials include not only writing of synopses, protocols and reports of health related studies but also conference presentations and scientific manuscripts, as well as submitting to scientific journals for publication.

Eurotrials follows internal SOPs for writing study synopses, protocols and reports. These procedures state the essential and optative contents that each document should possess, in accordance with The International Conference on Harmonisation of Technical Requirements for Registration of Pharmaceuticals for Human Use (ICH) guidelines (e.g. ICH-E6 for protocols and ICH-E3 for reports). The sponsor may require incorporation of additional information or even provide a template for the development of these documents. For submission of scientific manuscripts, the medical writer must be sure that the manuscript is submitted according to the guidelines stated in the journal.

Submission of clinical study reports to specialised databases (e.g. Clinicaltrials.gov) is also a task performed by the medical writer.

## Internship Report in Biostatistics

During my internship at Eurotrials, I performed several medical writing activities for the Epidemiology & Late Phase Research department. I worked in several observational studies, assisting in the development of various documents. Each of these activities is described below.

### ***a) Writing of a clinical study report in Rheumatology***

Period: November 2010 – January 2011.

This is a national, observational, cross-sectional, investigator-driven study in Rheumatology. The purpose of this study is to assess the compliance of current medical practice in the prescription of the study product with the applicable guidelines, in the management of rheumatoid arthritis.

My task was to assist in the writing of the study report, particularly the results section. Interaction with the Biostatistics department was very important for this task, in order to write clear, understandable and accurate information. With this project I realised the importance of the interaction between the medical writer and the Biostatistics department. The statistician is often asked to assist in the writing and/or to review of the results section of a study report or manuscript. The study report was finished and sent to the investigator in January 2011.

### ***b) Registration and Submission of a study report in Infectology to Clinicaltrials.gov***

Period: November 2010 – April 2011.

This is a national, observational prospective study in Infectology. The objective is to observe a specific therapeutic regimen in a defined population of the Human Immunodeficiency Virus (HIV) infected patients.

I only started to participate in this study after the report was submitted. The sponsor required the protocol and the report to be submitted in the Clinicaltrials.gov database. This is a federal database, supported by the US National Institutes of Health and holds the registry and results of clinical trials and observational studies conducted around the World (15). I was in charge of submitting all results related information. This included (16):

- Results Point of Contact: point of contact for scientific information.
- Agreements: section to certify the existence of agreements between the sponsor and investigator(s) that may restrict the investigator rights to discuss or publish results.
- Participant flow: summary of participants starting and completing each period.
- Baseline characteristics: demographic characteristics.
- Outcome measures: data summary, respective statistical analyses and respective results.
- Adverse events: summary of adverse event information.

A template with all fields to be filled is provided by Clinicaltrials.gov. The system automatically detects a broad set of errors and incoherencies that may be inserted, immediately warning the person that is inserting

## Internship Report in Biostatistics

the data that there is something wrong and must be corrected. The submission option is locked as long as these errors are unsolved.

After introduction of all information in the database, the results were submitted, at December 2010. At the first submission, the section was reviewed by Clinicaltrials.gov, which detected some inaccuracies. Therefore, a set of queries were sent to be clarified. These queries were solved in April 2011 and at the end of the internship the study was still being reviewed by Clinicaltrials.gov.

### ***c) Writing of the results and discussion sections of a scientific manuscript in Pulmonology***

Period: March 2011 – April 2011.

This is a national, observational, cross-sectional, study in Pulmonology. It aims to obtain prevalence data for Chronic Obstructive Pulmonary Disease.

The study sponsor conducted this study and performed the statistical analysis. Eurotrials was responsible for writing the scientific manuscript in order to be submitted.

My task was to write the results and discussion sections of this document. Information to be included was pre-selected by the sponsor and all relevant statistical results were provided. The manuscript was aimed for publication in a specific scientific journal, also previously stated by the sponsor. The document was written in accordance with the rules stated in the website of the journal.

The study was submitted in April 2011 and at the end of the internship was still being reviewed by the sponsor.

### ***d) Writing of a study protocol in Infectology***

Period: May 2011.

This is a national, observational, prospective study in Infectology. The study was made to assess the effectiveness, as well as safety, of a specific therapeutic change in patients with HIV infection.

My task was to assist the development of the study protocol. A template for the protocol was provided by the sponsor, as well as a description of the study. The study synopsis had already been written by the medical writer, also according to a template given by the sponsor. This synopsis was also very helpful in the development of this protocol. This task also included literature research, in order to write the study introduction. For this section, I searched for information regarding common therapeutic changes in antiretroviral therapy, as well as for the conditions and issues associated this particular change. The synopsis already provided a study rationale. However, I also carried on further research on this section, by request of the medical writer.

The protocol was written and sent to the sponsor for approval.

***e) Writing of a study synopsis in Pulmonology***

Period: May 2011.

This is a national, observational, cross-sectional study in Paediatric Pulmonology. The objective of this study is to obtain prevalence data on Respiratory Syncytial Virus (RSV) infection, in hospital setting. It also has a prospective component for clinical assessment of infected individuals.

I was asked to write the protocol synopsis for this study. Support bibliography was provided by the sponsor: The study outline was provided after a meeting between the sponsor and the Eurotrials Representative for this study. The synopsis template was also provided by the sponsor. For the introduction, I searched for information characterising the biology and epidemiology of this infection. The writing of the rationale involved searching for information that justified the performance of this study.

The synopsis was written and sent to the sponsor for approval.

***f) Manuscript Submission to a scientific journal of a study in Oncology***

Period: May 2011 – June 2011.

This is a national, observational national study in Oncology. The objective of this study was to evaluate the effectiveness and safety of a specific chemotherapy regimen in the treatment of breast cancer.

The sponsor selected a journal for submission of this manuscript. My task was to consult the journal's guidelines for submission, assist in adapting the manuscript in order for it to comply with these guidelines and proceed to the submission.

Guidelines for submission of manuscripts generally include topics such as:

- Maximum allowed number of words, for the manuscript and for the abstract.
- Required language(s).
- Rules for font size or type, as well as the use of commas or points for decimals.
- Reference style.
- Submission of tables and figures.

A set of recommendations for the contents of the manuscript were sent to the sponsor, in order to properly adapt the document in accordance with the journal's guidelines.

### 2.3. Research & Development

The purpose of the R&D department is to assist in the creation and implementation of clinical development plans, establishment of partnerships with pharmaceutical companies, as well as basic research institutions. This department is very important for the strategic planning of the Company, as it helps the development of new projects, programs and work routes.

The R&D department gains special relevance with the concept of translational medicine. This concept is defined as the efficient and effective translation of basic scientific findings relevant to human disease into knowledge that benefits patients (17). It aims to accelerate the rational transfer of new insights and knowledge into clinical practice, improving patients' outcomes and Public Health (18). This is a multidisciplinary effort, involving coordinated efforts between Academia, Regulatory Authorities and Industry (17).

Eurotrials sees great potential in translational medicine as a pathway for the development of novel and effective treatments and works in the establishment of partnerships with pharmaceutical companies, as well as basic research institutes and Academia, with the objective of accelerating knowledge and technology transfer, from basic to applied science. The efforts of Eurotrials towards translational medicine aim to promote the clinical applicability of basic science and knowledge, as well as to streamline clinical research and development.

Being a department with a strong role in the strategic profile of the Company, the information I am allowed to transmit in this report is very limited.

#### *a) Development of the regulatory basis for a clinical development plan for an advanced therapy product*

My first experience with the R&D department occurred in December 2010. I was invited to assist in the planning of a clinical development plan for an advanced therapy product. My role in this project was to search and resume regulatory information relevant for the development of these products.

Several regulatory and guidance documents were found relevant for this project. Among them, the following documents specific to advanced therapy products can be found:

#### European Commission Law

- Regulation (EC) No 1394/2007 of 13 November 2007: this regulation lays down the specific rules concerning the authorisation, supervision and pharmacovigilance of advanced therapy medicinal products (19).
- Regulation (EC) No 668/2009 of July 2009: this regulation implements Regulation (EC) 1394/2007 of the European Parliament and of the Council with regard to the evaluation and certification of quality and non-clinical data relating to advanced therapy medicinal products developed by micro, small and medium-sized enterprises (20).

European Medicines Agency (EMA) guidelines:

- Guideline on safety and efficacy follow-up – risk management of advanced therapy medicinal products: this guideline describes specific aspects of pharmacovigilance, risk management planning, safety and efficacy follow-up, authorised for advanced therapy medicinal products, as well as aspects relevant to the clinical follow-up of patients treated with such products (21).
- Guideline on the minimum quality and non-clinical data for certification of advanced therapy medicinal products: this guideline describes the minimum quality and non-clinical set of data that small and medium-sized enterprises developing advanced therapy medical products should submit for scientific evaluation when seeking EMA certification of quality (22).
  - The purpose of this certification system is to promote the development of such products by small-medium enterprises. It facilitates the evaluation of applications for a clinical trial authorisation or marketing authorisation application (23).

### ***b) Partnership opportunities with research teams in a basic research institute***

In the scope of translational medicine, in February 2011, I was invited to participate in the search for partnerships between Eurotrials and a specific basic research institute. In this institute, several research projects are ongoing either *in silico*, *in vitro* or *in vivo* in animal models. The steps of this activity are the following:

- Identification of research teams that are developing projects with potential for application in clinical setting.
- Contacting with these teams, preferably through a meeting.
- Presenting of the Company and the services it provides that can be useful in this context. The research team then describes the research project and states potential applications and benefits adjacent to a possible application in clinical setting.
- Discussion of possible approaches for the planning of a clinical development, if the project seems feasible.
- If the project is feasible to be applied in a clinical study, an outline describing the product and possible study is developed. After this a sponsor/financing source for clinical development (usually a pharmaceutical company) is sought.

The clinical development plan may be performed by Eurotrials, in collaboration with the sponsor.

For this activity, I attended the meetings, together with the Eurotrials R&D representative, with all scheduled research teams and participated in the development of the meeting minutes, as well as a specific bibliographic research.

Still within the scope of this activity, I also assisted in the development of an outline for an observational investigator-driven study in cardiology. The outline defined: study rationale, main objectives, overall study design and procedures, planned sample size, planned study duration and a timeline.

## 2.4. Data Management

At Eurotrials, the Data Management department follows the study project since its beginning, providing the means for data collection. It is also accountable for the preparation of the study database for data analysis. The head of data management is responsible for all the activities held inside the department.

### *a) Case Report Form development*

The CRF is designed to record all information required by the approved protocol and contains all data to be analysed for the trial, in order to answer the study question and objectives (24). Also important is the annotated CRF. This document identifies the variables found in the CRF, associates each variable with the corresponding filling blank(s) and is a valuable tool for defining database structure. It also helps assisting data entry and statistical analysis. At Eurotrials, the Biostatistics department and the medical writer are responsible for the CRF review, in order to assess if the data is in accordance with the study protocol.

### *b) Database preparation*

#### *1. Database building*

For the construction of each study database, the head of data management assigns the study data manager. The database is built in accordance with variables in the CRF. The database is created only with the variable fields (no data). After this task, it is tested for quality control, where data is inserted for testing purposes.

#### *2. Data entry*

Data received from the study site are verified by the database operator. Errors are notified by the data manager and database operator to the clinical monitoring team. After the database is approved, data is inserted by the data entry operator. During data entry, all deviations are reported to the sponsor.

#### *3. Data validation and cleaning*

During data entry, the data entry operator must check for inconsistencies in the data inserted by the investigators. These are notified to the study data manager, who sends a query to the site. After the query is solved, data is again entered in the database and re-verified. Before closing the database, unsolved queries are reported to the sponsor. After data entry is complete, the study data manager proceeds with data cleaning. This consists in correcting differences between the two entries. This is applicable for double data entry, where data is inserted twice by two operators, unlike single data entry, where data is only inserted once.

Data validation ensures consistency and accuracy of the data and is responsibility of the study data manager. Inaccuracies found are sent as a query to the investigator. After resolution of all queries, the study data manager validates the database. Unsolved queries are noted by the data manager, in the Data Management Report. The data manager also meets with the sponsor to analyse and classify protocol deviations in minor or major.

#### 4. *Database quality control*

A sample of the subjects in the database is selected and the data inserted is compared with the subject study data. Errors are corrected by the quality control staff. If the error rates exceed a limit previously defined, a new sample is drawn, with the same size and without repeating any subject from the previous sample.

#### 5. *Database lock*

After all data is received, inserted and validated, all queries are solved (unsolved queries are listed in the Data Management Report), quality control is complete and protocol deviations are classified, the database is locked by the study data manager, after authorisation from the sponsor, and sent to the statistician.

#### 6. *Randomisation code break*

Opening of randomisation codes may be done only after the database is locked (except for a medical emergency), and only after the approval of the sponsor.

#### **c) *Electronic CRF (e-CRF)***

E-CRFs are the electronic representation of paper CRFs and can be used, for example, by application of data entry screens (25). Here, data are automatically transferred from the e-CRF to the database (no data entry necessary). The data manager inserts the queries into the system, being directly delivered to the investigator.

#### **d) *Activities developed inside the Data Management department***

With the help of the Teaching & Training department, I planned an internal training session directed, in part, to the Data Management department, where some standards of the Clinical Data Interchange Standards Consortium (CDISC) standards were discussed. This training session is further discussed in the section for the Teaching & Training department.

I also attended a Webinar entitled “Using the CDISC Standards End-to-End in Clinical Trials” (28), in April 2011, where following CDISC standards were discussed in the context of clinical trials (28):

- Clinical Data Acquisition Standards Harmonization (CDASH): standards for data capture in clinical trials. This standard identifies and standardises typical data fields found in CRFs (29).
- Study Data Tabulation Model (SDTM): standard structure for tabulation data to be submitted. Tabulation datasets are those in which each registry is one single observation for one subject (30).
- Analysis Data Model (ADaM): standard for creation of analytical datasets. Analytical datasets support statistical analysis results and may contain raw and/or derived data (31).

In March 2011, a day of my internship was dedicated to data management, where the department was presented by the head of data management. Objectives, procedure flowchart and main documents were presented and discussed. Most of the information that I gathered for this chapter came from this presentation.

## 2.5. Teaching & Training

### a) *Training Activities*

The Teaching & Training (T&T) department works in collaboration with all other Eurotrials departments. It helps each sector to schedule, prepare, conduct, manage and evaluate (when applicable) any training activity that is being developed. Training programs may be internal (given only to personnel working at Eurotrials) or external (open to people outside the Company).

My collaboration with this department concerned the development and conduction of an internal training session entitled *Bases de Dados e Análise Estatística: Recomendações para Ensaios Clínicos* (Databases and Statistical Analysis: Recommendations for Clinical Trials) (26), directed to the Biostatistics and Data Management departments. Its scope was to recall basic concepts and give an introduction to Food and Drug Administration (FDA) and CDISC recommendations for data management and statistical analysis for clinical trials. CDISC is an organisation that establishes standards to support the acquisition, exchange, submission and archiving of clinical research data (27).

This training session occurred in January 2011.

This presentation was divided in 3 sections:

- 1) ICH and Statistical analysis: Overview of important ICH guidelines for statistical analysis, namely ICH-E3: Structure and Content of Clinical Study Reports (14) and ICH-E9: Statistical Principles for Clinical Trials (32).
- 2) FDA and Statistical analysis: Introduction to the FDA document Study Data Specifications (v. 3.1.2.). These specifications are required for the submission of animal and human study datasets in electronic format to the FDA (33).
- 3) CDISC Standards for collection, management and analysis of study data. Three CDISC standards were discussed in this session: SDTM, CDASH and ADaM. Each of these standards is discussed in the Data Management department section.

For this session, I performed several activities along with the T&T department assistant. These activities included:

- Invitation of potential participants: an e-mail was sent to all personnel working in the Biostatistics and Data Management departments, as well as other collaborators that could be interested in participating, inviting them to attend this session.
- Scheduling of the training session: a proposal for a date and time free for all was sent by e-mail to all attendees.
- Preparation of the documentation: a presence sheet was required for all the attendees to fill. Also, a feedback evaluation sheet was given to each attendee.
- Room reservation: a room was previously reserved for this session.

## Internship Report in Biostatistics

- Room preparation: in order to properly conduct this session, an image projector was installed.
- Availability of the presentation: a few hours before the beginning of the session, the presentation was placed in a folder assessed by all the attendees within the Eurotrials Internal Server.

After the training session was over, I was responsible for returning all support equipment, as well as to collect all evaluation sheets. The presence sheet was retrieved by the T&T department. The results of the attendees evaluation is presented in **Table 1**. “No Answer” values were considered missing and did not enter in the relative frequencies.

**Table 1. Global Evaluation of the training action (n=8)**

|  |   |       |
|--|---|-------|
| <b>Scientific contents, n (%)</b>              |   |       |
| Very good                                      | 4 | 57,1% |
| Good   | 3 | 42,2% |
| No Answer                                      | 1 |       |
| <b>Clear and objective information, n (%)</b>  |   |       |
| Very good                                      | 6 | 85,7% |
| Good   | 1 | 14,3% |
| No Answer                                      | 1 |       |
| <b>Time spent, n (%)</b>                       |   |       |
| Enough   | 7 | 100%  |
| No Answer                                      | 1 |       |
| <b>Training objectives accomplished, n (%)</b> |   |       |
| Yes  | 7 | 100%  |
| No Answer                                      | 1 |       |

### ***b) Information Management***

Eurotrials has an implemented Information Management System.

Information management consists in maintaining the Eurotrials Digital Library, as well as releasing of a monthly newsletter with the relevant scientific information, guidelines and regulations.

On November 2010 I was invited to participate in the management of the Eurotrials Digital Library, collecting relevant and updated bibliography within the scope of statistics applied to health research, a task that was maintained until the end of the internship. Information retrieved included regulatory and guidance information, as well as scientific articles, on the subjects of statistical planning, analysis and reporting for health research.

## 2.6. Clinical Trials

According to ICH-E6 – GCP a clinical trial is defined as follows:

*Any investigation in human subjects intended to discover or verify the clinical, pharmacological and/or other pharmacodynamic effects of an investigational product(s), and/or to identify any adverse reactions to an investigational product(s), and/or to study absorption, distribution, metabolism, and excretion of an investigational product(s) with the object of ascertaining its safety and/or efficacy (34).*

New medicines must prove to be safe and effective before being authorised for marketing. The clinical trial is the fundamental tool of therapeutic evaluation, being essential for this purpose (35). All clinical trials should be performed according to sound scientific principles and considering all ethical issues involved, in order to achieve the trial objectives (35).

Due to their large variety, clinical trials are generally divided temporally in the clinical development process, in four phases:

- Phase I: assessment of tolerability, preliminary safety, pharmacokinetics and pharmacodynamics, when applicable (35). The drug is administered to small number of healthy volunteers, or patients, in cases such as drugs with high toxicity (eg. cytotoxic drugs) (36).
- Phase II: these trials usually initiate exploration of the therapeutic effect on patients (36). Safety and efficacy is tested in patients and an optimal dose is sought (1).
- Phase III: randomised, controlled trials on hundreds to thousands of patients to determine efficacy and safety on a substantial scale (35). These studies provide an adequate basis for marketing authorisation (36).
- Phase IV: post-licensing studies on large samples of the target population. Studies characterised by broad inclusion criteria (35). Objectives are mostly focused for optimisation of drug use and gathering of additional information (36).

Due to their sensitivity relatively to ethical and regulatory issues, the planning, conduction and submission of a clinical trial are highly regulated activities. Some of the rules to be followed while working in a clinical trial in Portugal include:

European Commission Directives:

- 2001/20/EC: implementation of good clinical practice in the conduct of clinical trials on medicinal products for human use (37).
- 2005/28/EC: establishment of principles and detailed guidelines for good clinical practice as regards investigational medicinal products for human use (38).

## Internship Report in Biostatistics

### National Law:

- Law no 46/2004: transposition of Directive 2001/20/CE (39).
- Decree-Law no 102/2007: transposition of Directive 2005/28/CE (40).
- Law no 67/98: protection of personal data (41).

### Other essential Documents:

- Declaration of Helsinki (2008 version): definition of ethical principles for medical research involving human subjects (42).
- ICH guidelines: although all applicable guidelines should be followed, one should be highlighted:
  - ICH-E6 – GCP: international pattern of requirements of ethical and scientific quality that must be respected in the development, conducting, submission and reporting of clinical trials (13).

Clinical trials conducted in Portugal need to be previously approved by the National Authority of Medicines and Health Products (INFARMED) (39), the Commission of Ethics for Clinical Research (CEIC) (39) and the National Committee for Data Protection (CNPD) (43).

In May 2011, the Clinical Trials department was presented to me by a lead clinical research associate (CRA). A CRA is responsible for ensuring that the trial is conducted, recorded and reported in accordance with the study protocol, SOPs and all the applicable ethical and regulatory requirements defined above (13).

In this presentation, the role of Eurotrials in clinical trials was explained. This is the biggest sector of the Company and it may follow the trial from start to finish. The CRA assists the sponsor in selecting appropriate centres for the studies and performs pre-study visits to assess the investigator's interest and if the site is adequate for the trial (1). After the centre is selected, the CRA follows it from start to finish, working in the office and in the centre, by carrying out on-site visits. On-site visits can be divided, according to the time progression of the trial, as:

- Site initiation visit: the purpose of this visit is to orient the study staff to the requirements of the protocol (1). At this time point, the centre must be ready to start the trial (1). Study medication and materials must be available and required documentation complete and available (1). This meeting should be attended by all the study staff (1). All major points of the protocol are reviewed and discussed (1). Finally, the CRA must also train the study staff in the proper way to conduct the study and to handle the study medication and materials (1)

## Internship Report in Biostatistics

- **Monitoring visit:** these visits allow an in-process quality control of the data (1). CRAs must also continually ensure that the study is conducted, recorded and reported in accordance with the protocol (1). This is carried out by several processes (1). One of them is Source Data Verification (SDV) that consists in comparing the filled CRF with the source documents in the study centre (1). The objective of this activity is to confirm the consistency of the data in the CRF, solving of queries and confirming the completeness of the data (1). But the CRA role is much more than this. Monitoring includes confirming that the following activities are performed (1):
  - Proper obtaining of the informed consent.
  - Compliance to the protocol and study entry criteria.
  - Identification of any safety issues (adverse events and serious adverse events).
  - Proper accountability and reconciliation of study medication.
  - Continued adequacy of facilities and study staff.

After this visit, the CRA prepares a monitoring report for the sponsor and a follow-up letter to the centre.

- **Close-up visit:** after the last patient completed the trial, the study may be closed (1). In this visit, the CRA assures that the study medication is reconciled, the integrity of double blind codes is confirmed and that any queries left are solved and documented (1). The CRA must also confirm the best way to keep the source data and the investigator must notify the Ethics Committee of the completion of the study (1).

### *a) Activities developed inside the Clinical Trials department*

In May 2011, I enrolled in a common activity in the Clinical Trials department: sending of clinical study documents to the Central Study File, at the sponsor Headquarters. The study objective is to compare the safety and effectiveness of two different therapies for the treatment of pulmonary arterial hypertension. It is a multicentre, double-blind, randomised, placebo-controlled, parallel, prospective Phase IV study.

Activities performed included:

- Tagging all documents that were assigned to me (in order to identify them).
- Making a photocopy for each one of these documents (photocopies were archived in the Company File for the study).
- Filling a form identifying what documents would be sent to the Central Study File.
- Request for document shipping and ensure that the files are properly sent to the sponsor.

### **3. Monodisciplinary activities within the Biostatistics department**

#### **3.1. Use of statistics in health research**

Statistics is defined as the science of collecting, summarising, presenting and interpreting data, with the purpose of estimating the magnitude of associations and testing hypothesis (44). This science is vital for health research, as it allows information to be organised on a wider and more formal basis, instead of relying in the exchange of anecdotal evidence and personal experience (44). Statistics is also important for accounting biological variation processes and helps the researcher to prioritise the importance of the data obtained (7).

Biological data are highly variable (45). Measurements of human subjects rarely give exactly the same results from one occasion to the next (45). Therefore, it is expected that any comparison made in this context shows differences (45). For example, it cannot be expected that a certain drug induces a similar response in the same patient, under the same condition in two different situations (1).

It is often believed that the job of the statistician is simply data analysis (45). Contrary to this belief, statistical input is highly valuable since the beginning of the study (45). In fact, statistical thinking may be useful in the very beginning, at the definition of the study objective(s) (45). Each analysis for a specific objective is done through statistical analysis on the variables applied to such objective (45). The type of analysis depends on the type of variables to be analysed. In consequence, the definition of the objective and the study variables should consider the input of the statistician, for optimal analysis and result interpretation (7). But this is not the only phase of the project where statistics may be useful.

##### ***a) Study design***

Research can be divided into observational and experimental studies (9).

In observational studies information is collected about one or more groups of subjects, but nothing is done to affect them (9). These studies can be prospective, where subjects are recruited and data are collected about subsequent events, or retrospective, where information is collected about past events (9). They may also be cross-sectional, which are studies that measure the variables of interest in a defined time spot (e.g. prevalence studies) (7).

In experimental studies, the researcher affects what happens to all or some of the individuals (9). This may involve the administration of a therapeutic intervention, whether pharmacological or not, or even other types of intervention, such as a new approach to a medical consult (7). Randomised clinical trials are classical examples of experimental studies (7). In these studies, two randomised groups are compared, in which one group is given a study intervention, while the other is given a placebo or a standard/comparing treatment (7).

The definition of the statistical analysis methodology depends on the definition of the study type and variables of interest to be measured (44).

Another important aspect of the study design is the sample (9). The purpose of most medical research is to generalise the study results from the study sample to the population (9). For this objective, two aspects need to be considered (9):

- Sample should be representative of the population of interest.
- Groups being compared should be as alike as possible, in comparison studies.

The statistician needs to demonstrate that the sample is representative of the population (44). A sample size calculation needs to be performed, justifying the size estimated and demonstrating that this sample is able to answer the study questions (44). For clinical trials, ICH recommends that the basis for sample size calculation, statistical considerations and practical limitations should be provided in the clinical study report, as well as methods used for the calculation, along with any source of reference (14).

The sampling scheme is also relevant to assure the representativeness of the sample. In theory, a truly representative sample can only be obtained by choosing subjects at random (9). In practice, samples are nearly always chosen systematically and the subjects' characteristics are defined in a set of inclusion and exclusion criteria, so that their representativeness can be judged (9).

Homogeneity between groups is vital to assure that no systematic differences arise during allocation between groups (44). Randomisation is a technique that ensures an equal distribution of characteristics between treatment groups (46). It is the best way to assure homogeneity of study groups and reduce bias (9). Blinding is another important tool to reduce bias. The judgment of the investigator may be affected by knowing the treatment that a subject is getting (9). A double blind-design, in which neither participants nor study personnel know what treatment is being administered, should be done whenever possible (44).

### ***b) Statistical Analysis***

The aim of statistical analysis is to use the information gathered from the study sample to make inferences about the relevant population.

Statistics may be classified as descriptive, traditional analytic or Bayesian analytic (46):

- Descriptive statistics: data analysis using non-comparative techniques.
- Analytical statistics: mathematical formulas, models and tests to analyse data (also called inferential statistics).
- Bayesian statistics: considers other events, data and information from the past and incorporates that information when analysing data. This type of statistical analysis is not included in this report, as no activity was performed with it.

In most research studies, data are collected for descriptive purposes, such as demographic and clinical characteristics of the study subjects (9). Some observational studies only use descriptive analysis. Interventional studies always require an inferential component (9).

## Internship Report in Biostatistics

There are two basic approaches to statistical analysis: Estimation and hypothesis testing (9).

Estimation relies on generalisation of data from the study sample to the corresponding population (9). This is done through the determination of a range of values within the Confidence Interval. This range is where the true value is, under a certain level of confidence (9).

Hypothesis testing is used for comparative analysis (9). The majority of statistical analyses involve comparison (9). The numerical value corresponding to the comparison of interest is often called the effect (9). In this approach, two hypotheses can be defined (9):

- Null hypothesis: the effect of interest is zero. This statistical null hypothesis is often the negation of the research hypothesis that generated the data.
- Alternative hypothesis: usually, the effect is not zero.

Having set up the null hypothesis, the probability of the observed data being consistent with the null hypothesis is assessed (9). This probability is called the p-value. This is defined as the probability of having the observed data when the null hypothesis is true (9). When the p-value is below a pre-determined cut-off (e.g. 0,05), the result is called statistically significant. Above it, it is called non-significant (9).

The use of a cut-off p-value leads to treating the analysis as evidence for decision making (46). Two possible errors may be made when using p-value to make a decision (46):

- Type I error ( $\alpha$ ) is the probability of declaring a statistical difference when there is none. It may be viewed as the significance level necessary for the statistical test to detect a difference between treatments that is clinically meaningful (e.g.  $\alpha = 0,05$ ).
- Type II error ( $\beta$ ) is the probability of not detecting the difference that is looked for if it is present. Power ( $1-\beta$ ) is the probability of detecting this difference. It is desirable to have a high power of probability of detecting a difference between the treatments.

The goal of any clinical study is to have small type I and II errors (46). Since there is a trade-off between the two types of errors, it will be necessary to decide which goal is more important (46). During the calculation of the study sample, the statistician needs to consider the acceptable cut-off values of type I and II errors, in order to properly calculate the sample (7).

It can be seen that statistical analysis uses a probabilistic approach (45). This happens due to the high variability of biological data (45). Variations between comparison groups may be due to real effects, random variation, or both (45). The statistician decides how much variation should be ascribed to chance, so that the remaining can be assumed to be due to a real effect (45).

### 3.2. Overview of the statistical methods performed

This section gives an overview of the purpose and rationale of all statistical methods that I performed during my internship. It should be noted that there are many other statistical tests and methods that were not covered in this report. I only present those that I used for the study projects I was assigned to.

#### a) *Descriptive statistics*

In all studies, the first analysis to be performed should be a descriptive analysis of each study variable (univariate descriptive analysis) (7). This analysis will give the statistician an overview of the main characteristics of the study subjects (7).

For comparative studies, a bivariate descriptive analysis may be performed. This is defined as the descriptive statistics within each comparative group (7). This analysis is useful in determining the homogeneity of study groups, as well as giving an insight in differences among variables of interest/effect (7).

Descriptive statistics must account the specific aspect that is being measured (variable). Variables are (7):

- Numerical:
  - Discrete: limited number of discrete values.
  - Continuous: values within a continuous scale.
- Categorical:
  - Nominal: non-ordered or dichotomic categories (e.g. sex).
  - Ordered: values within ordered categories (e.g. age group).

It is also important to distinguish dependent and independent variables (7):

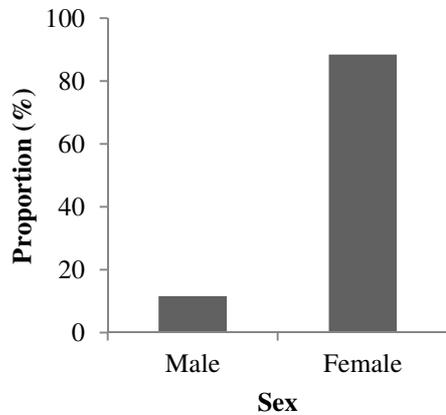
- Independent: variable that may determine the occurrence of an outcome (e.g. study medication).
- Dependent: variable that represents an outcome influenced by the independent variable (e.g. adverse event caused by administration of the study medication).

#### 1. *Categorical variables*

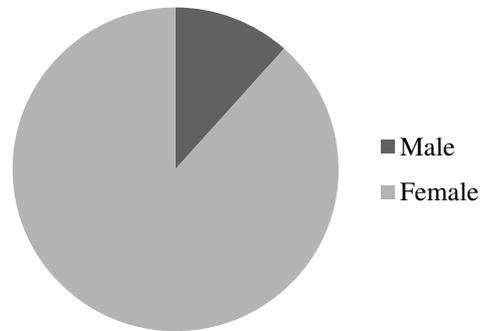
Descriptive statistics for categorical variables require frequency tables (7). For each variable, the corresponding absolute frequency and relative frequency are calculated and shown in these tables (7). Absolute frequency (n) is represented by the count of subjects that possess the category of the analysed variable, while relative frequency (%) is represented by percentage values, giving the proportion of subjects with the category analysed (7).

Absolute frequencies and relative frequencies can be illustrated by bar charts or pie charts (7). In a bar chart, the length of the bars are proportional to the count or percentage of the corresponding category (**Figure 1**) (44). In pie charts, a circle is divided in sectors, each one proportional to the count or percentage of the

corresponding category (**Figure 2**) (44). These illustrations may be used to point out results that are considered more relevant for the specific study (7).



**Figure 1. Bar Chart**



**Figure 2. Pie Chart**

2. *Numerical variables*

For numerical variables, measures of central tendency and dispersion are usually determined.

- Central tendency
  - Mean ( $\mu$ : population;  $\bar{x}$ : sample): the sum of the values divided by the number of values (44). This is the most frequent measure to define the central tendency of the data (7).
  - Median: the midway value that delimitates 50% of the ordered sample (7).
  - Mode: the value which occurs more often (44).
- Dispersion
  - Variance ( $\sigma^2$  population;  $s^2$ : sample): function of the deviations of all observations from the mean ( $x - \bar{x}$ ) (**Figure 3**).

$$s^2 = \frac{\sum(x - \bar{x})^2}{(n - 1)}$$

**Figure 3. Formula for variance**

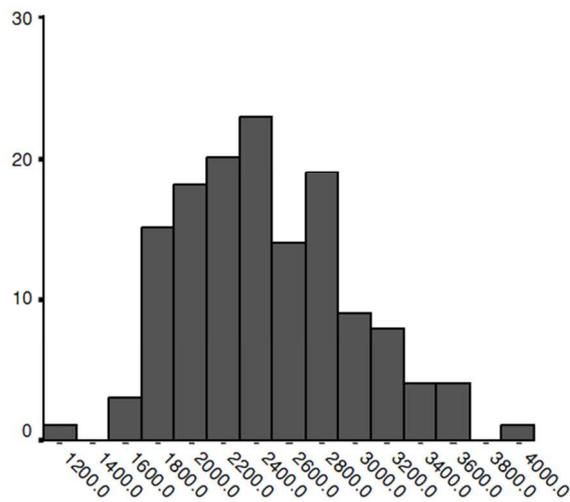
- Standard deviation ( $\sigma$ : population;  $s$ : sample): this is the square root of the variance (44). Variance has the disadvantage of presenting the measures in the square of the units of the observation variable, so this is a way to present dispersion values in the original units (44).
- Amplitude values: maximum and minimum values (7).

For illustration of numerical variables, two charts are most commonly used: the histogram and the boxplot (7).

The histogram is a bar chart for the numerical value, in which each bar represents a class of values (**Figure 4**) (7). Its main objective is to assess the distribution of the variable values (7).

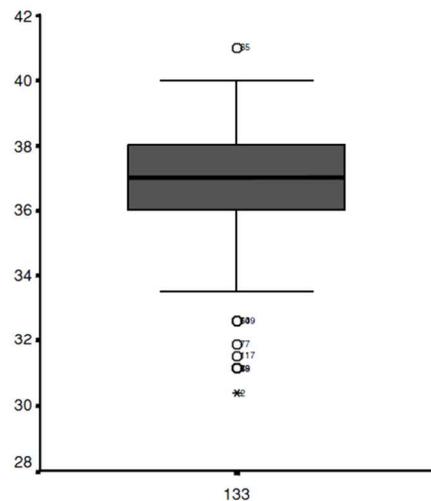
The boxplot (**Figure 5**) allows the observer to see the central tendency of the data, as well as its dispersion (7). Some aspects that need to be described for this illustration:

- The box is delimited by the first quartile (bottom limit), representing 25% of the ordered sample and third quartile (top limit), representing 75% of the ordered sample. In consequence the size of the box characterises the interquartile amplitude (7).
- The central line represents the median (7).
- Top and bottom lines characterise the maximum and minimum values when they are within 1,5 times or below the interquartile range (box length) (47).
- Outlier values between 1,5 and 3 box lengths from the upper or lower edge of the box are shown as open circles and identified with the corresponding database row (47). Extreme values more than 3 box lengths from either edge of the box are shown as asterisks (47).



Source: Peat, 2005 (47)

**Figure 4. Histogram**

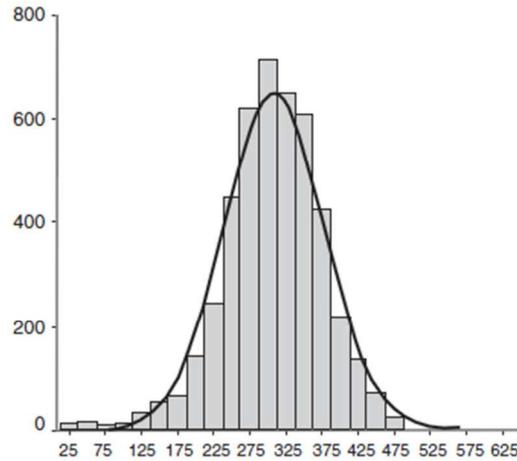


Source: Peat, 2005 (47)

**Figure 5. Box Plot**

**b) The Normal Distribution**

The normal distribution (also called Gaussian distribution) is a type of distribution characterised by a histogram with a normal curve, symmetrical about the mean and bell shaped (**Figure 6**) (44). This bell is tall and narrow for small standard deviations and short and wide for large ones (44).



Source: Bowers, 2008 (48)

**Figure 6. Normal Distribution**

The normal distribution plays a central role in statistical analysis (44). The reason is that sample means, provided that sample size is not too small, will be normally distributed around the true population mean (44). The larger the sample selected, the closer the sample mean will be to being normally distributed (44). This is known as the central limit theorem (44).

It is also known that the standard deviation of the distribution of the sample mean corresponds to the standard error (9). This measure is the quotient of the standard deviation and the square root of the sample size ( $\sigma/\sqrt{n}$ ) (9). Sample standard deviation (s) may be used in this calculation (9).

The application of the central limit theorem allows the statistician to determine the confidence interval for a population mean, from a sample standard normal distribution (44).

The standard normal distribution is a normal curve whose mean is 0 and standard deviation is 1 (44). This curve is created from a normal distribution by subtracting the mean from each observation and dividing by the standard deviation (44). This originates the Standard normal deviate (SND, expressed by z), also known as the z-score (**Figure 7**) (44).

$$z = \frac{(x - \mu)}{\sigma}$$

Source: Kirkwood, 2001 (44)

**Figure 7. Formula for z-score**

The z-score expresses the number of standard deviations an observation is away from the mean (e.g. z=1 corresponds to a value that is one standard deviation ahead of the mean:  $\mu+\sigma$ ) (44).

Exactly 95% of the distribution lies between z-scores -1,96 and 1,96. So, for a standard normal distribution, z=1,96 is the reference value for 95% confidence (7).

With the z-score, standard error and sample mean, the confidence interval (CI) for the population mean can be determined (**Figure 8**) (7). For 95% CI, z-score=1,96 (7). This result is valid for large samples, by replacing  $\sigma$  with s (44).

$$CI = \bar{x} \pm \left( z \times \frac{\sigma}{\sqrt{n}} \right)$$

**Source: Silva, 2009 (49)**

**Figure 8. formula for confidence interval**

By observing this formula, it can be inferred that for the same z, a higher standard error leads to a wider CI, and vice versa, meaning that a more precise sample (with less wide deviations from the mean) leads to a narrower CI (7).

Normal distributions are also important to help choosing the correct hypothesis test for statistical analysis (44). Parametric tests are used when assumptions can be made about the distribution of the sample, so they can be used in normal distribution samples (44). Non-parametric tests are used as an alternative, when normality of the distribution is not assumed (44).

Kolmogorov-Smirnov and Shapiro-Wilk tests may be used to assess normality.

The Kolmogorov-Smirnov test compares the cumulative frequencies of the sample distribution with the cumulative frequencies of what would be expected if the distribution was normal (50). In my internship it was used for larger samples ( $n \geq 50$ ). As for Shapiro-Wilk, it compares the ordered sample values with those that would be expected if the distribution was normal (44). In my case, it was used for samples with sample size below 50.

The null hypothesis in both tests state that there is no difference between the expected distribution and observed distribution, meaning that with a p-value larger than 0,05 (or another defined threshold), the statistician does not reject the hypothesis that this distribution is normal (47).

### **c) Inferential statistics**

#### **1. Estimation**

As stated before, estimation is one method used to generalise observations from a sample of subjects to the respective population (9).

The first measure to calculate, when there is need to estimate a result to the population is the standard error (9). For comparative analysis, the determination of this measure differs for two situations:

i. *Difference between two sample means*

When comparing two independent samples, the variance of the difference between their means is the sum of the separate variances, meaning the standard error (SE) is the square root of this sum (**Figure 9**) (9).

$$\begin{aligned}
 SE(\bar{x}_1 - \bar{x}_2) &= \sqrt{\text{var}(\bar{x}_1) + \text{var}(\bar{x}_2)} \\
 &= \sqrt{\{se(\bar{x}_1)\}^2 + \{se(\bar{x}_2)\}^2} \\
 &= \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}
 \end{aligned}$$

**Source: Altman 1990 (9)**

**Figure 9. Formula for standard error of the difference between two means**

ii. *Difference between two proportions*

The standard error of the proportion  $p$  in a sample of size  $n$  is  $\sqrt{p(1-p)/n}$  (9). The calculation rationale is similar than the difference between two means (variance of the difference of proportions is the sum of the separate variances) (9). So, by having two proportions  $p_1$  and  $p_2$ , the standard error can be determined using a similar formula (**Figure 10**) (9).

$$\begin{aligned}
 SE(p_1 - p_2) &= \sqrt{\text{var}(p_1) + \text{var}(p_2)} \\
 &= \sqrt{\{se(p_1)\}^2 + \{se(p_2)\}^2} \\
 &= \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}
 \end{aligned}$$

**Source: Altman 1990 (9)**

**Figure 10. Formula for standard error of the difference between two proportions**

iii. *Confidence intervals*

A confidence interval is a range of values on which the true value is included, under a certain level of confidence (9). A confidence interval for an estimated mean extends either side of the mean by a multiple of the standard error (9). For example, the 95% confidence interval (the most commonly used) is the range of values from  $\bar{x} - 1,96SE$  to  $\bar{x} + 1,96SE$  (9). For this reason, the determination of the standard error is essential to calculate the confidence interval (9).

For a 95% confidence interval, it is expected that the true population value will not be included 5% of the time (9).

2. Comparing categorical data of two or more independent samples

i. Chi-Square test

The Chi-Square test is used to verify if there is an association between categorical variables (7).

In order to perform the Chi-Square test, data must be arranged in contingency tables (44). A contingency table is a way to examine a relationship between two categorical variables. Usually, the rows of the table correspond to the exposure values (independent variables) and the columns to the outcomes (dependent variables) (44). These tables come with absolute and relative frequencies that allow the statistician to perform the Chi-Square test, while percentage values help orienting the way of these differences (7). An example of a contingency table crossing two binary variables (2x2 table) is presented in **Table 2**.

**Table 2. Association between smoking habits and diagnosis of lung cancer**

| Study variable | Variable categories | Statistical measures | Lung cancer diagnosis positive, n=50 | Lung cancer diagnosis negative, n=50 |
|----------------|---------------------|----------------------|--------------------------------------|--------------------------------------|
| Smoking habits | Non-Smoker          | n (%)                | 14 (28%)                             | 47% (94%)                            |
|                | Smoker              | n (%)                | 36 (72%)                             | 3% (6%)                              |

**Source: Aguiar, 2007 (adapted) (7)**

For this test, the null hypothesis is the non-existence of population differences between the groups (with lung cancer *versus* without lung cancer, according to the example stated in **Table 2**) to the independent variable (7).

This test compares the observed values (O) in each of the categories in the contingency table with the values to be expected if there were no differences between the groups defined by the dependent variable (44). In order to determine the Chi-Square statistic, expected frequencies (EF) must first be calculated (**Figure 11**) (46).

$$EF = \frac{\text{Row Total} \times \text{Column Total}}{\text{Total Sample}}$$

**Source: Kirkwood, 2001 (44)**

**Figure 11. Formula for expected frequencies**

After the EF for each table cell is calculated, Chi-Square can be obtained, (**Figure 12**) (46). Calculation of Chi-Square statistic needs to account the degrees of freedom (d.f.), defined as the number of values in the final calculation of a statistic that are free to vary (7).

$$\chi^2 = \sum \frac{(O - EF)^2}{EF}; d.f. = (Rows - 1) \times (Columns - 1)$$

Source: Kirkwood, 2001 (44)

Figure 12. Formula for Chi-Square statistic

Using the degrees of freedom, it is possible to determine if the Chi-square statistic obtained is statistically significant, using **Table A 1** (44). If the test is significant, the null hypothesis that there are no population differences between groups is rejected (7).

ii. *Fisher's exact test*

If more than 20% of the EF calculated is below 5, Fisher's exact test should be used for 2x2 tables (association analysis between two binary variables) and no test should be performed for other contingency tables (9).

This test assesses the probability associated with all possible 2x2 tables with the same row and column totals as the observed data (9). All possible sets of frequencies are determined (one of which corresponds to the O), as well as the probability of each combination arising if the null hypothesis (no population differences between the groups) is true (9). The p-value of this test corresponds to such probability for O. The lower the p-value, the lower is the probability of the O combination arising when the null hypothesis is true (9).

3. *Comparing numerical data of two independent samples*

i. *Two independent samples t-test*

The two independent samples t-test is a hypothesis test that compares data of the means of two independent normal distributions obtained from a numerical variable (e.g. post-treatment score in two groups) (7). The null hypothesis in this test states that there are no differences between the two groups (7). The formulas needed to determinate the t-test statistic (t) are illustrated in **Figure 13** and explained as follows (46):

- The t statistic is the quotient of the difference between the means of the two groups ( $\bar{x}_1 - \bar{x}_2$ ) and the standard error of this difference (SED).
- The variances  $s_1^2$  and  $s_2^2$  are combined to create  $s^2$ . This is done with the assumption that both variances are homogenous (7). If the variances are heterogeneous, the SED is calculated with  $\sqrt{s_1^2/n_1 + s_2^2/n_2}$  instead of the common variance (7).
- SED is determined by the square root of the sum of the two standard errors of both means.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{SED} \qquad SED = \sqrt{s^2/n_1 + s^2/n_2}$$

$$s^2 = [(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2]/(n_1 + n_2 - 2)$$

$$d.f. = n_1 + n_2 - 2$$

**Source: Kirkwood, 2001 (44)**

**Figure 13. Calculation of two sample t-test statistic**

An approximate p-value corresponding to the t statistic may be obtained from Student’s t distribution (**Table A 2**).

ii. *Mann-Whitney U test*

Mann-Whitney U test is a non-parametric analysis, used alternatively to the t-test to compare independent samples that do not follow a normal distribution (44).

Many non-parametric methods are based on rank attribution (44). Each value of the outcome variable is replaced by a rank, after the variable is sorted into ascending order of magnitude (e.g. values 200, 3 and 1 are given ranks of 3, 2 and 1, respectively) (46). The Mann-Whitney U test relies on the sequential placement of these ranks (44).

The sequential placement index (U), has two sets of values (50):

- Number of times the values of Group B precede those of Group A.
- Number of times the values of Group A precede those of Group B.

The null hypothesis states that both U indexes are not different (50).

For a practical example, **Table 3** should be considered.

**Table 3. Values and respective ranks for two groups A and B**

|                |    |    |    |    |    |    |    |    |    |
|----------------|----|----|----|----|----|----|----|----|----|
| Ordered values | 45 | 51 | 53 | 64 | 70 | 75 | 78 | 82 | 93 |
| Rank           | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
| Group          | A  | B  | B  | A  | B  | A  | A  | A  | B  |

**Source: Feinstein, 2002 (Adapted) (50)**

Being U the number of times the values of Group B precede Group A values and U’ the reverse (**Figure 14**):

$$U = 0 + 2 + 3 + 3 + 3, U = 11$$

$$U' = 1 + 1 + 2 + 5, U' = 9$$

**Figure 14 Example of calculation of U index**

To calculate the p-value of this test, **Table A 3** should be consulted (50). This table states the p-value for the smaller of the U indexes determined (in this example,  $U'=9$ ). Also, and following **Table 3**, considering that, for  $n_1=4$  and  $n_2=5$  the smallest U value for  $P=0,05$  is 1, it can be concluded that, for this example,  $P>0,05$  and the null hypothesis is not rejected at a 5% significance level.

4. *Comparison of numerical data in two paired samples*

i. *Paired t-test*

A paired samples hypothesis test aims to compare two observations made on the same subject in two separate occasions (7). Therefore, this study is applicable to longitudinal studies (7). One example is the assessment of pain in one patient before and after a treatment is administrated.

The null hypothesis states that there are no differences between the two observations (7).

Calculating the t-test statistic (t) for a sample with size n involves calculating the variable of the difference between the two paired observations (7). After this, the mean ( $\bar{x}_{dif}$ ) and variance ( $s_{dif}^2$ ) of this variable are determined (7). Now all measures are available for calculation of the standard error of the difference ( $SD_{dif}$ ) and consequently, the t statistic (**Figure 15**) (7).

$$t = \frac{\bar{x}_{dif}}{SD_{dif}}$$

$$SD_{dif} = \sqrt{\frac{s_{dif}^2}{n}}$$

$$d.f. = n - 1$$

**Source: Silva, 2009 (51)**

**Figure 15. Calculation of Paired t-test**

P-value determination is similar to the two sample t-test. P-value is determined with the consulting of Student's t distribution (**Table A 2**) (7).

It should be noted that the variable of the difference between the two observations should be normally distributed, as this is a parametric test (7). Wilcoxon test should be used as a non-parametric alternative (7).

ii. *Wilcoxon signed rank test*

This test is a non-parametric alternative for analysis of a non-normal, but symmetrical distribution of the variable of the difference between the two paired observations (7). The null hypothesis states that there are no differences between the two paired variables (7). Like Mann-Whitney U, it ranks the values of the difference variable by ascending order of magnitude (52). Then, positive and negative difference ranks are summed separately (52). The test statistic T is the lesser of these two sums. If the null hypothesis is true, it is expected that these two sums are the same (52). After calculating T, **Table A 4** is consulted to obtain the p-value (53). T must be equal or less to the table value for significance at the respective level (53).

iii. *Sign test*

In case of a non-normal and asymmetrical distribution of the variable of the difference between the two paired observations, the sign test is performed (7). The null hypothesis states that there are no differences between the two paired variables (7). This test calculates the differences between the two paired variables for all cases and classifies these differences as positive (+), negative (-) or none (0) (54). A large difference between the number of signs points out a median significantly different from the null hypothesis (54). After this, the test statistic S is the smaller number of – or + signs (54). With this statistic and sample size, **Table A 5** should be used to determine the p-value (53).

5. *Comparison of binary data in two paired samples*

i. *McNemar test*

The McNemar test is used to compare binary data in two separate occasions or two paired samples in the same subject (7). Two typical health related situations where this test can be used include (7):

- Assess statistically significant differences between two diagnostic tests used on the same person.
- Assess statistically significant differences between the presence of a clinical sign/symptom in two different occasions.

The null hypothesis in this test states that there are no differences between the two paired samples (7).

According to its formula (**Figure 16**) and considering **Table 4**, this test is a chi squared statistic with 1 degree of freedom. Only the discordant pairs (Yes-No pairs: B and C) will be used in this calculation (44).

**Table 4. Symptoms before and after treatment**

|                      |     | Symptom<br>(After) |    |
|----------------------|-----|--------------------|----|
|                      |     | Yes                | No |
| Symptom<br>(Before ) | Yes | A                  | B  |
|                      | No  | C                  | D  |

$$X_{paired}^2 = \frac{(B - C)^2}{B + C}, d.f. = 1$$

**Source: Kirkwood, 2001 (44)**

**Figure 16. Calculation of McNemar Statistic**

As this is a Chi-Square statistic, the p-value is obtained the same way as in the Chi-Square test (using **Table A 1** to find the p-value correspondent to the result obtained).

6. *Comparison of numerical data in three or more independent samples*

i. *Kruskal-Wallis test*

This is a non-parametric test used to compare more than two independent samples (7). The parametric alternative for normal distribution samples is the Analysis of Variance (ANOVA) (7) that will not be discussed in this report, as it was not performed during my internship. Like Mann-Whitney U test, it gives a rank to each value, in ascending order of magnitude and proceeds to sum all ranks from each group (55).

The null hypothesis for this test states that all populations have identical distributions (55).

The test statistic H is calculated according to the formula in **Figure 17(55)**.

$$H = \left[ \frac{12}{n_t(n_t + 1)} \sum_{i=1}^t \frac{T_i^2}{n_i} \right] - 3(n_t + 1), d.f. = t - 1$$

**Source: Njuho, 2002 (55)**

**Figure 17. Calculation of Kruskal-Wallis statistic (H)**

In this formula,  $n_t$  corresponds to the total sample size (considering all groups t),  $n_i$  is the total number of items in sample i, the summation operator represents the sum of each group i ( $i=1, \dots, t$ ) and  $T_i^2$  represents the sum of the ranks for group i (55).

For example, considering **Table 5**, where 15 values are ranked (with the lowest value ranked 1 and the largest ranked 15), the calculation in **Figure 18** can be applied.

**Table 5. Distribution of values and respective ranks for groups A, B and C**

| Group A                      | Rank      | Group B                      | Rank      | Group C                      | Rank      |
|------------------------------|-----------|------------------------------|-----------|------------------------------|-----------|
| 50                           | 4         | 80                           | 11        | 60                           | 7         |
| 62                           | 8         | 95                           | 14        | 45                           | 2         |
| 75                           | 10        | 98                           | 15        | 30                           | 1         |
| 48                           | 3         | 87                           | 12        | 58                           | 6         |
| 65                           | 9         | 90                           | 13        | 57                           | 5         |
| <b>Sum (T<sub>1</sub>) =</b> | <b>34</b> | <b>Sum (T<sub>1</sub>) =</b> | <b>65</b> | <b>Sum (T<sub>3</sub>) =</b> | <b>21</b> |

Source: Njuho, 2002(51, 55)

$$H = \frac{12}{15(16)} \left[ \frac{34^2}{5} + \frac{65^2}{5} + \frac{21^2}{5} \right] - 3(16) = 10,22$$

**Figure 18. Example of calculation of Kruskal-Wallis statistic**

For this test, the same **Table A 1** for the Chi-Square analysis is used. Considering that there are 2 degrees of freedom, and that the Chi Square statistic for 2 degrees of freedom and p=0,05 is 5,99, it can be seen that 10,22>5,99, therefore p<0,05. The null hypothesis is rejected at 5% significance level.

7. *Correlation analysis between two numerical variables*

In order to quantify the magnitude the strength of the association between two numerical variables x (independent) and y (dependent), the Pearson coefficient of correlation (r) is determined (**Figure 11**) (7).

$$r = \frac{\sum(x - \bar{x})(y - \bar{y})}{\sqrt{[\sum(x - \bar{x})^2 \sum(y - \bar{y})^2]}}$$

**Source: Kirkwood, 2001 (44)**

**Figure 19. Formula for Pearson’s correlation coefficient**

The correlation coefficient is always a number between -1 and 1. The closer it is to one of these numbers, the stronger the association between the two variables. If the variables are not associated, this coefficient is 0 (7). Values close to 1 show an association where high values of x tend to go with high values of y, while values close to -1 suggest that high x values are associated with low y values and vice-versa (44).

When distribution is not normal, or two categorical ordered variables need to be correlated, Spearman coefficient is calculated, as a non-parametric alternative to Pearson’s (7). Like the Mann-Whitney U test, this method attributes a rank to each value, following an ascending order of magnitude (44). In order to determine Spearman’s rank correlation coefficient, first these ranks are attributed independently to both variables x and y (44). Then, the same formula is applied, only instead of values, the assigned ranks are used. Interpretation is the same as in Pearson’s (44).

8. *Prediction of the effect of exposure variables on a binary outcome - Logistic Regression*

i. *Logistic regression model*

Logistic regression is used to explain or predict the outcome of a binary variable, in function of the effect of various independent variables (co-variables, predictor variables), whether they are numerical, categorical or binary (7). Practical applications in health research include determination of the effect of a certain risk factor or treatment on a binary variable, such as the presence of a certain disease or occurrence of an adverse event, adjusted for the effect of possible confounding variables (7) (variables related both to the outcome and to the exposure groups being compared (44)).

In order to perform this analysis, a regression model must be built and assessed for its quality (7). For this analysis, the estimation of Odds Ratios (OR) needs to be done (7). The Odds in favour of an event represent the possibility of occurrence of a certain characteristic (56). It is the quotient between the probability of the outcome occurring (p) and the inverse of such probability (**Figure 20**) (56).

$$Odds = \frac{P}{(1 - P)}$$

**Source: Silva, 2009 (56)**

**Figure 20. Formula for Odds**

The OR is the quotient between the two Odds (56). For example, the exposure OR for a disease is given by the quotient of the Odds of the affected exposed subjects and the Odds of the affected non-exposed subjects.

Back to the logistic regression analysis, the model for logistic multiple regression is given by the formula in **Figure 21** (7):

$$P = \frac{e^L}{(1 + e^L)}, \text{ where } L = B_0 + B_1 \times X_1 + B_2 \times X_2 + \dots + B_k \times X_k$$

**Source: Aguiar, 2007 (7)**

**Figure 21. Model of logistic multiple regression**

Where (7):

- $P$  is the estimated probability of occurrence of the outcome.
- $e$  is the value usually used in exponential function  $\approx 2,718$ .
- $B_k$  are the k regression coefficients estimated and related to k independent variables.
- $B_0$  estimated constant for the model.
- L function (also called LOGIT function) is the determination of the natural logarithm of the Odds (**Figure 22**).

Estimation of OR in function of regression coefficients B is an important part of this model that allows for the prediction of binary outcomes (7). OR is obtained through  $e^B$  (7).

$$L = LOGIT$$

$$LOGIT = \ln\left(\frac{P}{(1-P)}\right) = B_0 + B_1 \times X_1 + B_2 \times X_2 + \dots + B_k \times X_k$$

**Source: Aguiar, 2007 (7)**

**Figure 22. Calculation of LOGIT Function**

An example of the application of this model can be the prediction of a certain health event, with the effect of the new treatment (nt), adjusted for smoking habits (sh) and age. Considering **Figure 23** (with  $B_0=-13,106$ ) (7):

$$P = \frac{e^L}{(1 + e^L)}, \text{ where } L = -13,106 - 3,208 \times nt + 2,343 \times sh + 0,245 \times age$$

**Source: Aguiar, 2007 (7)**

**Figure 23. Example of application of logistic multiple regression model**

The following odd ratios may be estimated for the three independent variables (**Figure 24**)

$$OR(nt) = e^{-3,208} = 0,04$$

$$OR(sh) = e^{2,343} = 10,399$$

$$OR(age) = e^{0,245} = 1,278$$

**Source: Aguiar, 2007 (7)**

**Figure 24. Example of calculation of Odds Ratios**

The negative value of the correlation coefficient for the new treatment states that the risk of occurrence of the health event is lower for patients receiving this treatment, in comparison with the comparator (7). On the other hand, considering smoking habits, the risk of the occurrence of the health event in smokers is 10,4 times higher, comparing to non-smokers (7). Also, there is an incremental risk of 1,3 by increasing year of age (7).

If  $B=0$ , then  $OR=1$  (7). This result means that there are no differences between exposed and non-exposed (7).

ii. *Confounding variables and modification of effect*

In order to analyse for possible confounding, first a statistical analysis that associates the confounding variable with the exposure variable, as well as other analysis associating the confounding variable with the outcome, should be performed (7). If there was an association statistically significant, then it can be assumed that it is indeed a confounding variable (7).

If there is a need to determine if a categorical variable is confounding, a stratified analysis could also be performed (7). In other words, the OR can be calculated for each stratum of the possible confounding variable. A combined OR can be calculated considering the strata OR (7). This is called the Mantel-Haenszel common Odds Ratio (7).

If the strata OR are similar and all similarly distant to the combined OR, then it can be assumed that this is a confounding variable (7).

Other situation that requires attention is if the strata OR are assumed as different from each other, thus unable to be combined in a single adjusted OR (7). This is a case of modification of effect, and the strata OR should be presented separately (7).

Logistic regression may be used to adjust the effect of confounding variables in the outcome (7).

The null hypothesis of this analysis states that there is no association between the outcome and the independent variable adjusted for the effect of other independent variables within the model (7). P-value of the effect of the independent variable, adjusted for other independent variables in the model is determined with the Wald statistic that is the quotient between the regression coefficient  $B$  and its standard error  $SE$   $(B/SE)^2$  (7). This is a Chi-Square statistic for 1 degree of freedom, so, **Table A 1** is also used for this test to determine the p-value (7).

### 3.3. Work developed at Eurotrials

My internship at Eurotrials was mainly focused in the Biostatistics department. I assisted in various observational research projects, as well as projects involving bibliographic research and sample stratification. Each project where I was involved is discussed below, as well as my participation in each one. Results are not displayed, due to confidentiality considerations.

During my work, I did not need to calculate any statistical test by hand, since I was using statistical software that automatically runs all calculations needed. This software was SPSS® 16.0.

Before enrolling in any project, I underwent a brief training period, in September 2010. This training involved the following activities:

- Reading of the SOPs of the Biostatistics department.
- Reading of the book *Guia Prático Climepsi de Estatística em Investigação Epidemiológica: SPSS* (Climepsi Practical Guide on statistics in Epidemiologic Research).
- Resolution of a set of exercises provided by the Biostatistics department. These exercises included definition of study variables, calculation of confidence intervals and performance of several hypothesis tests.

#### a) *Statistical analysis for a Geriatric observational study*

Period: September 2010.

This is a national observational, cross-sectional study on Geriatrics. Its objective is to assess the functional dependence level on a geriatric population, as well as to characterise a set of functional and laboratory factors associated to human aging.

For this study I was invited to perform additional statistical analysis requested by the sponsor, as the main analysis had already been performed and submitted. This analysis focused entirely on the assessment of cognitive function (binary variable: favourable/non-favourable) of the individuals studied, in association with demographic characteristics, social aspects and functional dependence level.

Considering that the variables susceptible for additional analysis were all categorical, descriptive analysis consisted only on frequency tables, with absolute and relative frequencies. As for inferential statistics, the following tests were performed:

- Chi-Square test/Fisher exact test: for comparison of categorical data in two independent variables.
- Odds Ratio: for significantly associated variables, in order to measure the magnitude of association.
- Multiple logistic regression: analysis of all statistically significant variables, in order to assess the effect of each exposure variable adjusted to the effect of other associated variables.

## Internship Report in Biostatistics

I also wrote the statistical report for this additional analysis. This report included the study objectives, statistical methodology and rationale, a section for results, where significant results were listed and a tabulation for all results of the analysis performed.

This analysis was performed and the report was submitted to the sponsor.

### ***b) Statistical analysis for an observational study on Public Health***

Period: September 2010 – April 2011.

This is a national, observational, prospective study on Public Health. The objective is to assess the effect of intense heat exposure on several clinical outcomes (need for healthcare, occurrence of symptoms or occurrence of previously defined medical conditions). Demographic characteristics, presence of certain diseases and residence status also were assessed for association with these clinical outcomes.

My role in this study was to perform the statistical analysis and develop the statistical report.

Descriptive statistics were carried out for categorical (frequency tables for absolute and relative frequencies) and numerical variables (mean, median, standard deviation, minimum and maximum).

As for inferential statistics, both comparisons were made for independent study groups and paired samples of before-after exposure. All groups were exposed to intense heat, diverging only on their geographic location. Before comparing numerical data, normality tests (Kolmogorov-Smirnov and Shapiro-Wilk) were performed. All distributions to be analysed were non-normal, so only non-parametric tests were performed.

- Chi-Square/Fisher exact test: for comparison of categorical data in two independent samples.
- Mann-Whitney U: for comparison of numerical data in two independent samples.
- Wilcoxon signed-rank test: for comparison of numerical data before and after exposure.
- McNemar: for comparison of a binary outcome before and after exposure.
- Multiple logistic regression: assessment the effect of each independent variable that showed statistically significant association with health outcomes, adjusted to the effect of other associated independent variables.

The report I developed contains a rationale for the study, objectives, primary and secondary endpoints, statistical methodology and rationale, a result section, describing descriptive results with the aid of charts, statistically significant results, statistical conclusions where statistical results are interpreted and tabulation values for all analyses performed.

The first draft of this study was sent in January 2011. After feedback from the sponsor a second draft was submitted in April 2011. At the end of this internship this draft was still on review.

*c) Statistical analysis for a study on clinical practice regarding interventions in Angiology*

Period: November 2010.

This is an observational, cross-sectional study in the area of Angiology. The objective is to assess the effect of angioplasty effectiveness and safety in current practice. The type of procedure applied, risk factors and severity of illness are associated with a pre-defined set of health outcomes.

For this study, I assisted another statistician. My role was to review the statistical results and copy them to the tables on the statistical report. I also had to adapt the statistical report to the results found.

Descriptive statistics were carried out for categorical (frequency tables for absolute and relative frequencies) and numerical variables (mean, median, standard deviation, minimum and maximum).

The only hypothesis test performed was the Chi-Square/Fisher exact test, to compare categorical data in 2 independent samples. Only categorical data were requested to be compared by the sponsor.

The sponsor requested a very simple report. Only result tables and a brief description of each table, stating what data are being analysed and highlighting statistically significant associations were carried out.

This analysis was performed and the report was submitted to the sponsor.

*d) Application of the Multi Attribution Decision Model (MADM) in two studies for Rheumatology and Dermatology*

Period: January 2011 – April 2011.

Two expert panels were organised to discuss the therapeutic options for two different diseases, in Rheumatology (rheumatoid arthritis) and Dermatology (plaque psoriasis). These studies are described together in this report, since the methodology used was the same for both.

The objective of the panel is to determine the most appropriate therapeutic option for each of a set of pre-defined clinical cases, considering that each clinical case is different (e.g. an acute case is different than a remission case). A questionnaire is provided to each physician invited and sets a number of clinical attributes to be considered (e.g. short term efficacy, long term safety, convenience).

The physician gives a score for each clinical attribute applied to each clinical case and therapeutic alternative, in ascending order of importance. The scores assigned to the therapeutic alternatives are generic (do not consider any of the clinical cases in study).

For example, a physician interprets that short term efficacy is more important for an acute clinical case and gives it a higher score, in comparison to a remission case. This physician also assigns a higher score for short term efficacy to the therapeutic alternative that acts sooner.

Scores attributed allow to (57):

- Assess the relative performance of each alternative with respect to each attribute.
- Assess the relative importance of each attribute for each clinical case.

The MADM is a decision-support method widely used to select between different solutions to a particular problem on the basis of multiple attributes (58). In this study, the mathematical model used to conduct this method is called Technique for Ordered Preference by Similarity to the Ideal Solution (TOPSIS).

This method integrates data from the performance of each treatment alternative with the attribute importance weights for each clinical case (57). After this integration, the model identifies a hypothetical ideal therapeutic alternative for each case, as well as an anti-ideal, which is the worst possible alternative (57). Finally, it measures the distance of each real alternative from the hypothetical ideal and anti-ideal (57).

A TOPSIS score is calculated scaled from 0% (identical to the anti-ideal) and 100% (identical to the ideal). This way, the alternatives may be ranked for each clinical case (57).

The mathematical model is not included in this report.

My work in these two studies was to build this model (according to the indications of the statistics Consultant), run it and provide the results to the medical writer, who wrote the study reports.

Both studies required submission of multiple reports, the last ones being sent to the sponsor in April 2011.

### *e) Statistical analysis of an observational study in Pulmonology*

Period: February 2011.

This is an observational, retrospective, case-control study in Paediatric Pulmonology. The objective is to identify individual risk factors for the development of bronchiolitis by studying the association of the presence of various possible risk factors with typical disease clinical and laboratorial outcomes.

For this study I performed statistical analysis and assisted in writing the statistical report.

Descriptive statistics were carried out for categorical (frequency tables for absolute and relative frequencies) and numerical variables (mean, median, standard deviation, 25<sup>th</sup> and 75<sup>th</sup> percentiles, minimum and maximum).

For inferential analysis, comparative analysis between groups exposed and non-exposed to the risk factors were carried out. Before comparing numerical data, normality tests (Kolmogorov-Smirnov and Shapiro-Wilk) were performed. Tests performed include:

- Chi-Square/Fisher exact test: for comparison of categorical data in two independent samples.

## Internship Report in Biostatistics

- Two independent samples t-test/Mann-Whitney U: for comparison of numerical data in 2 independent samples.
- Odds Ratio: for all statistically significant results, Odds Ratios were determined in order to assess the magnitude of association.
- Multiple logistic regression: assessment of the effect of each risk factor statistically significant, adjusted to the effect of other associated risk factors, as well as demographic factors.

Some of the statistical analysis was already performed by the sponsor. Therefore, a part of the report was already written. I only wrote the results on the already created tables and gave a small interpretation of statistically significant data, by request of the sponsor.

This analysis was performed and the report was submitted to the sponsor.

### *f) Statistical analysis of an observational, cross-sectional study in Oncology*

Period: March 2011

This is an observational, international, cross-sectional study in the area of Oncology. The objective of this study is to assess the prevalence of a specific genetic mutation in Lung Cancer, by studying a sample of patients with this illness.

My role in this study was to perform the statistical analysis and write the statistical report.

Descriptive statistics were carried out for categorical (frequency tables for absolute and relative frequencies) and numerical variables (mean, median, standard deviation, minimum and maximum).

Descriptive statistics were enough for the study objective, as this is a prevalence study. However, some hypothesis tests were also carried out, in order to associate the effect of the presence of the genetic mutation and clinical outcomes. Before comparing numerical data, normality tests (Kolmogorov-Smirnov and Shapiro-Wilk) were performed. Tests performed include.

- Chi-Square/Fisher exact test: for comparison of categorical data in two independent samples.
- Two independent samples t-test/Mann-Whitney U: for comparison of numerical data in two independent samples.

A template was available for the report. I wrote the results, statistical conclusions and tabulations sections.

This analysis was performed and the report was submitted to the sponsor.

### *g) Statistical analysis of an observational, prospective study, in Oncology*

This is an observational, national, prospective clinical cohort study in the area of Oncology. Its objective is to evaluate the effectiveness of a study medication in patients with Prostate Cancer, as well as to compare the effectiveness of two doses of this medication, evaluating clinical and laboratorial outcomes.

## Internship Report in Biostatistics

My role was to perform the statistical analysis and write the statistical report.

Descriptive statistics were carried out for categorical (frequency tables for absolute and relative frequencies) and numerical variables (mean, median, standard deviation, minimum and maximum).

For inferential statistics, comparisons were made for independent study groups (corresponding to the two doses studied) and for paired samples of before-after treatment. Before comparing numerical data, normality tests (Kolmogorov-Smirnov and Shapiro-Wilk) were performed. Tests carried out included:

- Chi-Square/Fisher exact test: for comparison of categorical data in two independent samples.
- Two independent samples t-test/Mann-Whitney U: for comparison of numerical data in two independent samples.
- Paired samples t-test/Wilcoxon signed-rank test/Sign test: for comparison of numerical data before and after treatment.
- McNemar: for comparison of binary outcomes before and after treatment.

As for the clinical report, a template for was already available. I wrote the results, statistical conclusions and tabulations sections.

This analysis was performed and the report was submitted to the sponsor.

### *h) Sample stratification in a study on Pulmonology*

Period: May 2011 - June 2011

This is a national, observational, cross-sectional study in Paediatric Pulmonology. The objective of this study is to obtain prevalence data on RSV infection, in hospital setting. It also has a prospective component for clinical assessment of infected individuals.

My role in this study was to distribute the previously estimated study sample for each hospital involved, considering the population covered by the specific centre. For this, the following steps were taken:

- Using population data provided by the National Statistics Institute, stratified by region and age, as well as considering the population covered by each hospital, estimates were made about the paediatric population covered.
- After this, the sample was distributed, using cross multiplication. The method is described in **Figure 25**, where total population (TP), total sample (TS), estimated population for centre x (PCx) and estimated sample for centre x (SSx) are shown:

|   |                                  |
|---|----------------------------------|
| $\begin{array}{l} TP \leftrightarrow TS \\ PCx \leftrightarrow SSx \end{array}$ | $SSx = \frac{PCx \times TS}{TP}$ |
|---|----------------------------------|

**Figure 25. Cross multiplication method**

The sample was stratified and the stratification proposal was sent to the sponsor.

***i) Statistical analysis of an observational, in Endocrinology***

Period: May 2011 – June 2011

This is an observational, national, cross-sectional study in the area of Paediatric Endocrinology. Its objective is to analyse the health status, demographic, socio-economic and lifestyle factors in children with obesity, as well as their parents, in order to assess the correlation of health status and lifestyle factors between parents and their children. My role in this study was to perform the statistical analysis.

Descriptive statistics were carried out for categorical (frequency tables for absolute and relative frequencies) and numerical variables (mean, median, standard deviation, minimum and maximum).

For inferential statistics, comparisons were made between parents and children, as well as correlation tests were performed for health and lifestyle factors between these two groups. Before comparing numerical data, normality tests (Kolmogorov-Smirnov and Shapiro-Wilk) were performed and all distributions were non-normal. Tests carried out included:

- Chi-Square/Fisher exact test: for comparison of categorical data in 2 independent samples.
- Mann-Whitney U: for comparison of numerical data in 2 independent samples.
- Kruskal-Wallis: for comparison of numerical data in more than 2 independent samples.
- Spearman Correlation: for assessment of correlation between the health status and lifestyle factors of parents and children.

This study started in the end of May 2011 and was still ongoing when my internship ended.

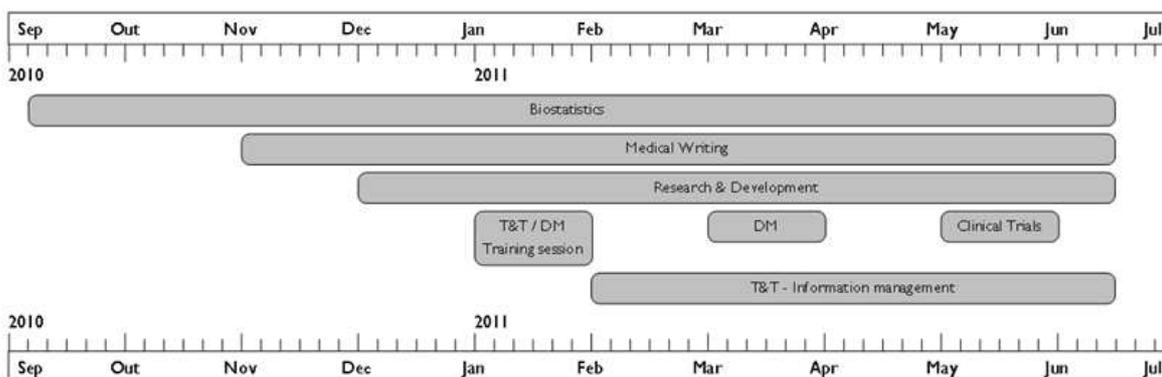
***j) Bulletin “Saúde em Mapas e Números” (Health in Maps and Numbers)***

As stated in Chapter 1, this document focuses on one health subject. I collaborated in the development of 3 bulletins during all my internship period: no 32 (Sleep Disturbances: December 2010) (11), no 33 (Occupational Health: March 2011) (59) and no 34. This bulletin is still to be released (due to June 2011).

The bulletin is divided in 3 sections: Portugal, Europe and World. The objective is to search and compile prevalence data (and other data deemed relevant, such as quality of life related information) for each of these sections. Information should be retrieved preferably from reliable and updated sources (such as Governmental portals, World Health Organization reports, or scientific articles with less than 10 years). Use of illustrations such as maps, bar charts or pie charts is also common to point out relevant information. Comparison of prevalence data between regions/countries, demographics, socio-economic status and other relevant factors are shown. However, only descriptive statistics are carried out. There is also a section named *Sabia que...* (Did you know...) where epidemiological trivia and other relevant information are highlighted. This bulletin serves as a marketing tool to promote the statistical services of Eurotrials.

#### 4. Discussion and conclusions

A timeline illustrating the duration of my internship activities is displayed in **Figure 26**. The segments of this chart represent each department, being divided in the T&T activities. This is due to the internal training session being a shared activity with the Data Management department (DM).



**Figure 26. Timeline of internship activities**

##### 4.1. Multidisciplinary Activity

By observing this timeline, it can be perceived that, although the Biostatistics department was my main point of activity, as planned, there were also other departments/activities where I dedicated a substantial amount of time. Medical writing was indeed a significant activity performed in my internship. This may partially be explained by the link between this activity and Biostatistics. The development of a study synopsis and protocol need the input of the statistician that plans or knows what statistical methodologies are planned for the specific study. The report also requires the attention of the statistician, which will help developing and/or review the statistical methodologies, results and discussions sections of the document, in order to ensure that the study data are properly reported and interpreted. Regardless, I also performed many activities out of the statistician scope, such as preparation and submission of manuscripts. The experience gathered by performing these activities gave me some sensibility on how to properly prepare several documents, such as study synopses, protocols and reports, as well as how to properly prepare scientific manuscripts for submission to a scientific journal. The latter was very important, as I had never carried out such activity before and did not fully realise the steps that needed to be taken to perform a submission. I had already learned that each journal sets specific rules and provides a specific pathway for submission (one journal may require tables and figures to be sent separately from the article, while other may state that the article should be sent in one piece), but the performance of these activities gave me a much more solid idea of how this process works.

The R&D department also deserved a significant amount of time, even though I did not participate in a large quantity of projects. This may be explained due to the fact that the projects involved in this department usually take a long time to be completed and require a very thoughtful and careful planning and development. The activities on this department required a large amount of time dedicated to bibliography research in order

## Internship Report in Biostatistics

to properly develop a study or a product development plan. This was a very challenging department as translational medicine requires the professional to understand both the basic scientific concept and the adequate clinical research plan that needs to be applied, if a transfer of knowledge from basic science to clinical application is intended. Indeed, these activities made me seek and interiorise large amounts of regulatory and scientific information that I did not have previous contact with. I also have now a better perspective on what can be done to transfer knowledge from pure to applied science.

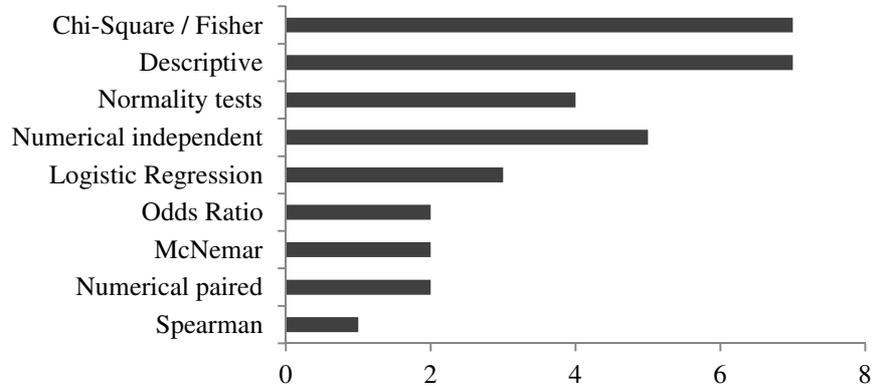
The timeline may lead the wrong impression by showing a relatively small amount of time dedicated to data management. What is displayed shows how much time I solely dedicated to this department, but as a statistician intern I was in constant contact with the data manager of any study I was working on. I always had to perform a quality control of the database that was sent and notify the data manager of any inaccuracies. Besides, other activities such as discussion of the arrangement of the database and formatting of specific variables, or discussion of the annotated CRF to help the analysis process were carried out frequently, showing that data management and statistics interact in a constant basis. My experience with data management allowed me to understand the process of database creation, data retrieval, entry and quality control. With this department I also learned important requirements for data organisation, namely the CDISC guidelines. These guidelines are not only very useful for the data manager, but also for the statistician, that must perform statistical analysis in accordance with the ADaM standard, if the study follows the CDISC recommendations.

Development and lecturing of an internal training session, in the scope of the T&T department, allowed me to understand what should be done to prepare these sessions, as well as what should be done after the lecture. It also gave me an opportunity to interiorise new information (notably CDISC and FDA information for clinical data management) and to perform a lecture. As for the information management task, bibliographic search was the most relevant activity. Here I learned new important sources for statistical information applied to clinical research.

The activities carried out in the scope of the Clinical Trials department showed me how complex can a clinical research project be. The activity I performed showed me how important it is for the sponsor to keep all relevant study documents and how everything sent to the sponsor must be identified in a standardised fashion and properly archived, thus showing the importance of proper handling of study documentation and knowledge of study procedures and good practices. It gave me a practical overview of the CRA working environment and how this professional is vital not only for the proper handling of study documentation, but also for the adequate performance of the trial, ensuring that the safety and rights of the patients are protected and that the study is carried out in accordance with the protocol and with all applicable regulatory and/or ethical requirements.

#### 4.2. Monodisciplinary activity: Biostatistics department

The statistical methods and analysis performed during this internship are illustrated in **Figure 24**.



**Figure 27. Statistical tests performed**

The high frequency of descriptive methods reveals the importance of data characterisation. Regardless of the objective of the study, these statistics are generally used to give the statistician an overview of the data, so that he/she can assess the data for inaccuracies, such as ages out of the defined limit (7). The high frequency of Chi-Square/Fisher analyses reveals that most of my studies involved comparison of categorical data. It can also be inferred that correlation and longitudinal studies were not so frequent.

Before I started to perform a test for the first time, I needed to recall it. A senior statistician would assess if I was able to work with the test, before I could start. With time, I gathered enough experience to run all these tests and correctly interpret the results almost independently.

The development of the statistical reports for my studies was also very important for my training. With this activity I learned how to present statistical data (using tables and illustrations, where relevant) and how to properly interpret this data, so that I could transmit it correctly to the investigator and/or sponsor.

The experience gathered by working with SPSS® 16.0 was also helpful. I started working with the syntax of the software, something that I was unaware of before entering this internship. This allowed me to work much faster with a high level of accuracy.

As stated before, interaction with other departments/activities, mainly the Data Management department and medical writing helped me understand the multidisciplinary network that needs to be arranged in order to successfully conduct a proper data analysis and develop a study report.

Overall, this experience in the Biostatistics department allowed me to better understand the essential role of statistics in health research and its importance not only for data analysis, but also for proper study planning and design. It also gave me practical knowledge on how to properly plan, perform and report the statistical analysis of health research data, in order to assure maximum accuracy and representativeness of the study population, so that it can be useful in improving healthcare and Public Health.



## References

1. Edwards LD, Fletcher AJ, Fox AW, Stoiner PD, editors. Principles and Practice of Pharmaceutical Medicine. 2 ed: JohnWiley & Sons Ltd; 2007.
2. Carlson PE. Clinical Research Industry Trends. National Center on Education and the Economy; 2007.
3. Nahler G. Dictionary of Pharmaceutical Medicine. 2 ed: Springer; 2009.
4. Getz K. CRO contribution to drug development is substantial and growing globally. Tufts Center for the Study of Drug Development Impact Report. 2006;8(1):4.
5. Almeida L. Contract Research Organizations (CRO). Aveiro: Training Programme in Pharmaceutical Medicine; 2010.
6. Eurotrials | Scientific Consultants. Lisbon2011 [updated 2011; cited 2011 7/May/2011]; Available from: <http://www.eurotrials.com/>.
7. Aguiar P. Guia Prático Climepsi de Estatística em Investigação Epidemiológica: SPSS. 1 ed: Climepsi; 2007.
8. Arnold R, J. G., editor. Pharmacoeconomics: From Theory to Practice. 1 ed: CRC Press; 2010.
9. Altman DG. Practical Statistics for Medical Research. 1 ed: Chapman and Hall/CRC; 1990.
10. Plano de Formação 2010/2011. Eurotrials, Scientific Consultants; 2010.
11. Silva T. Perturbações do Sono. In: Alves C, Silva C, Negreiro F, Vicente V, editores. Saúde em Mapas e Números: Eurotrials, Scientific Consultants; 2010.
12. Foote M. Medical Writing as a Career Choice. Drug Information Association. 2003.
13. ICH. Guideline for Good Clinical Practice. Step 5. E6 (R1): ICH; 1996.
14. ICH. Structure and Content of Clinical Study Reports. E3: ICH; 1995.
15. About ClinicalTrials.gov. US National Institutes of Health; 2008; Available from: <http://clinicaltrials.gov/ct2/info/about>.
16. ClinicalTrials.gov "Basic Results" Data Element Definitions. US National Institutes of Health; 2011; Available from: [http://prsinfo.clinicaltrials.gov/results\\_definitions.html](http://prsinfo.clinicaltrials.gov/results_definitions.html).
17. Littman BH, Krishna R, editors. Translational Medicine and Drug Discovery: Cambridge University Press; 2011.
18. Littman B, Di Mario L, Plebani M. What's next in translational medicine? Clinical Science. 2007;112.
19. Regulation (EC) No 1394/2007 of 13 November 2007, (2007).
20. Regulation (EC) No 668 / 2009 of July 2009, (2009).
21. EMA. Guideline on safety and efficacy follow-up – Risk management of advanced therapy medicinal products (Draft). EMA; 2008.
22. EMA. Guideline on the minimum quality and non-clinical data for certification of advanced therapy medicinal products. EMA; 2010.
23. EMA. Procedural advice on the certification of quality and nonclinical data for small and medium sized enterprises developing advanced therapy medicinal products. EMA; 2010.
24. UCLH. Guide on Good Practice for Data Management for Chief investigators of research sponsored by University College London. NHS; 2010.
25. Jolles E. Electronic CRF Design. International Clinical Sciences Support Center.
26. Silva T. Bases de dados e análise estatística: Recomendações para ensaios clínicos: Eurotrials, Scientific Consultants; 2011.
27. CDISC - Mission & Principles. CDISC; 2011; Available from: <http://www.cdisc.org/mission-and-principles>.
28. Gibson B. Using the CDISC Standards End-to-End in Clinical Trials: SAS; 2011.
29. CDISC. Clinical Data Acquisition Standards Harmonization (CDASH). CDISC CDASH Core and Domain Teams; 2008.
30. CDISC. Study Data Tabulation Model (SDTM) Final Version 3.1.1. CDISC Submission Data Standards (SDS) Team; 2008.
31. CDISC. Analysis Data Model (ADaM) Final Version 2.1. CDISC Analysis Data Model Team; 2009.
32. ICH. Statistical Principles for Clinical Trials. ICH-E9: ICH; 1998.
33. FDA. Study Data Specifications. FDA; 2009.
34. ICH. Guideline for Good Clinical Practice. Step 5. E6 (R1): ICH; 2002.

## Internship Report in Biostatistics

35. Griffin JP, O'Grady J, editors. The Textbook of Pharmaceutical Medicine. 5 ed: Blackwell Publishing; 2006.
36. ICH. General Considerations for Clinical Trials. ICH; 1997.
37. Directive 2001/20/EC of the European Parliament and of the Council of 4 April 2001, (2001).
38. Commission Directive 2005/28/EC of 8 April 2005, (2005).
39. Lei n.º 46/2004, de 19 de Agosto, (2004).
40. Decreto-Lei n.º 102/2007, de 2 de Abril, (2007).
41. Lei n.º 67/98 de 26 de Outubro, (1998).
42. Association WM. World Medical Association Declaration of Helsinki. 2008.
43. Deliberação N.º 333 / 2007, (2007).
44. Kirkwood BR, Sterne JAC. Essential Medical Statistics. 2 ed: Blackwell Publishing; 2001.
45. Campbell M, J, Machin D. Medical Statistics: A Commonsense Approach. 3 ed: Wiley; 1999.
46. Spilker B. Guide to Clinical Trials. 1 ed: Lippincott Williams & Wilkins; 1991.
47. Peat J, Barton B. Medical Statistics: A Guide to Data Analysis and Critical Appraisal. 1 ed: BMJ Books; 2005.
48. Bowers D. Medical Statistics from Scratch: An Introduction for Health Professionals Wiley-Interscience; 2008.
49. Silva C, Aguiar P, Vicente V. Complementos de Bioestatística e Epidemiologia - Módulo 2: Eurotrials, Scientific Consultants; 2009.
50. Feinstein A. Principles of Medical Statistics. 1 ed: Chapman and Hall / CRC; 2001.
51. Silva C, Aguiar P, Vicente V. Complementos de Bioestatística e Epidemiologia - Módulo 4: Eurotrials, Scientific Consultants; 2009.
52. Crichton N. Information Point: Wilcoxon Signed Rank Test. Journal of Clinical Nursing. 2000;9(4):1.
53. Eysenck MW. Fundamentals of Psychology. 3 ed: Psychology Press; 2009.
54. Larson R, Farber B. Elementary Statistics: Picturing the World 4ed: Prentice Hall; 2008.
55. Njuho PM. Statistical Analysis Research: University of Kwazulu-Natal Pietermaritzburg Campus 2002.
56. Silva C, Vicente V, Aguiar P. Complementos de Bioestatística e Epidemiologia - Módulo 3: Eurotrials, Scientific Consultants; 2009.
57. Ferrari M, Goadsby P, Lipton R, Dodick D, Cutrer F, McCrory D, et al. The use of multiattribute decision models in evaluating triptan treatment options in migraine. Journal of Neurology. 2005;252(9):1026-32.
58. Guibal F, Iversen L, Puig L, Strohal R, Williams P. Identifying the biologic closest to the ideal to treat chronic plaque psoriasis in different clinical scenarios: using a pilot multi-attribute decision model as a decision-support aid. Current Medical Research and Opinion. 2009;25(12):2835-43.
59. Silva T. Saúde Ocupacional. In: Alves C, Silva C, Negreiro F, Vicente V, editors. Saúde em Mapas e Números: Eurotrials, Scientific Consultants; 2010.

**Appendix 1**

**Table A 1. P-values for the Chi-Square distribution**

| d.f. | Level of significance for two-tailed test |       |       |       |       |
|------|---|-------|-------|-------|-------|
|      | 0,20                                      | 0,05  | 0,02  | 0,01  | 0,001 |
| 1    | 1,64                                      | 3,84  | 5,41  | 6,64  | 10,83 |
| 2    | 3,22                                      | 5,99  | 7,82  | 9,21  | 13,82 |
| 3    | 4,64                                      | 7,82  | 9,84  | 11,34 | 16,27 |
| 4    | 5,99                                      | 9,49  | 11,67 | 13,28 | 18,46 |
| 5    | 7,29                                      | 11,07 | 13,39 | 15,09 | 20,52 |
| 6    | 8,56                                      | 12,59 | 15,03 | 16,81 | 22,46 |
| 7    | 9,80                                      | 14,07 | 16,62 | 18,48 | 24,32 |
| 8    | 11,03                                     | 15,51 | 18,17 | 20,09 | 26,12 |
| 9    | 12,24                                     | 16,92 | 19,68 | 21,67 | 27,88 |
| 10   | 13,44                                     | 18,31 | 21,16 | 23,21 | 29,59 |
| 15   | 19,31                                     | 25,00 | 28,26 | 30,58 | 37,70 |
| 20   | 25,04                                     | 31,41 | 35,02 | 37,57 | 45,32 |
| 25   | 30,68                                     | 37,65 | 41,57 | 44,31 | 52,62 |
| 30   | 36,25                                     | 43,77 | 43,49 | 50,89 | 59,70 |
| 40   | 47,27                                     | 55,76 | 60,44 | 63,69 | 73,40 |

Source: Eysenck, 2009 (53)

**Table A 2. Student's t distribution**

| d.f. | Two-sided p-value |       |       |       |        |        |        |        |        |
|------|-------------------|-------|-------|-------|--------|--------|--------|--------|--------|
|      | 0,500             | 0,400 | 0,200 | 0,100 | 0,050  | 0,025  | 0,010  | 0,005  | 0,001  |
| 1    | 1,000             | 1,376 | 3,078 | 6,314 | 12,706 | 25,452 | 63,657 | —      | —      |
| 2    | 0,816             | 1,061 | 1,886 | 2,920 | 4,303  | 6,205  | 9,925  | 14,089 | 31,598 |
| 3    | 0,765             | 0,978 | 1,638 | 2,353 | 3,182  | 4,176  | 5,841  | 7,453  | 12,941 |
| 4    | 0,741             | 0,941 | 1,533 | 2,132 | 2,776  | 3,495  | 4,604  | 5,598  | 8,610  |
| 5    | 0,727             | 0,920 | 1,476 | 2,015 | 2,571  | 3,163  | 4,032  | 4,773  | 6,859  |
| 6    | 0,718             | 0,906 | 1,440 | 1,943 | 2,447  | 2,969  | 3,707  | 4,317  | 5,959  |
| 7    | 0,711             | 0,896 | 1,415 | 1,895 | 2,365  | 2,841  | 3,499  | 4,029  | 5,405  |
| 8    | 0,706             | 0,889 | 1,397 | 1,860 | 2,306  | 2,752  | 3,355  | 3,832  | 5,041  |
| 9    | 0,703             | 0,883 | 1,383 | 1,833 | 2,262  | 2,685  | 3,250  | 3,690  | 4,781  |
| 10   | 0,700             | 0,879 | 1,372 | 1,812 | 2,228  | 2,634  | 3,169  | 3,581  | 4,587  |
| 15   | 0,691             | 0,866 | 1,341 | 1,753 | 2,131  | 2,490  | 2,947  | 3,286  | 4,073  |
| 20   | 0,687             | 0,860 | 1,325 | 1,725 | 2,086  | 2,423  | 2,845  | 3,153  | 3,850  |
| 25   | 0,684             | 0,856 | 1,316 | 1,708 | 2,060  | 2,385  | 2,787  | 3,078  | 3,725  |
| 30   | 0,683             | 0,854 | 1,310 | 1,697 | 2,042  | 2,360  | 2,750  | 3,030  | 3,646  |
| 40   | 0,681             | 0,851 | 1,303 | 1,684 | 2,021  | 2,329  | 2,704  | 2,971  | 3,551  |
| 50   | 0,680             | 0,849 | 1,299 | 1,676 | 2,008  | 2,310  | 2,678  | 2,937  | 3,496  |

Source: Feinstein, 2001 (50)

**Table A 3. Critical Values for U**

| Two-tailed test significance level: 0,05 |                |   |   |   |    |    |    |    |    |    |    |    |    |
|--|----------------|---|---|---|----|----|----|----|----|----|----|----|----|
| n <sub>2</sub>                           | n <sub>1</sub> |   |   |   |    |    |    |    |    |    |    |    |    |
|  | 1              | 2 | 3 | 4 | 5  | 6  | 7  | 8  | 9  | 10 | 11 | 12 | 13 |
| 1  | —              | — | — | — | —  | —  | —  | —  | —  | —  | —  | —  | —  |
| 2  | —              | — | — | — | —  | —  | —  | 0  | 0  | 0  | 0  | 1  | 1  |
| 3  | —              | — | — | — | 0  | 1  | 1  | 2  | 2  | 3  | 3  | 4  | 4  |
| 4  | —              | — | — | 0 | 1  | 2  | 3  | 4  | 4  | 5  | 6  | 7  | 8  |
| 5  | —              | — | 0 | 1 | 2  | 3  | 5  | 6  | 7  | 8  | 9  | 11 | 12 |
| 6  | —              | — | 1 | 2 | 3  | 5  | 6  | 8  | 10 | 11 | 13 | 14 | 16 |
| 7  | —              | — | 1 | 3 | 5  | 6  | 8  | 10 | 12 | 14 | 16 | 18 | 20 |
| 8  | —              | 0 | 2 | 4 | 6  | 8  | 10 | 13 | 15 | 17 | 19 | 22 | 24 |
| 9  | —              | 0 | 2 | 4 | 7  | 10 | 12 | 15 | 17 | 20 | 23 | 26 | 28 |
| 10                                       | —              | 0 | 3 | 5 | 8  | 11 | 14 | 17 | 20 | 23 | 26 | 29 | 33 |
| 11                                       | —              | 0 | 3 | 6 | 9  | 13 | 16 | 19 | 23 | 26 | 30 | 33 | 37 |
| 12                                       | —              | 1 | 4 | 7 | 11 | 14 | 18 | 22 | 26 | 29 | 33 | 37 | 41 |

Source: Eysenck, 2009 (adapted) (53)

**Table A 4. Critical values for the Wilcoxon signed rank test**

| n  | Two-sided p-value |      |      |       |
|----|-------------------|------|------|-------|
|    | 0.1               | 0.05 | 0.02 | 0.002 |
| 5  | T≤0               | —    | —    | —     |
| 6  | 2                 | 0    | —    | —     |
| 7  | 3                 | 2    | 0    | —     |
| 8  | 5                 | 3    | 1    | —     |
| 9  | 8                 | 5    | 3    | —     |
| 10 | 11                | 8    | 5    | 0     |
| 11 | 13                | 10   | 7    | 1     |
| 12 | 17                | 13   | 9    | 2     |
| 13 | 21                | 17   | 12   | 4     |
| 14 | 25                | 21   | 15   | 6     |
| 15 | 30                | 25   | 19   | 8     |
| 16 | 35                | 29   | 23   | 11    |
| 17 | 41                | 34   | 27   | 14    |
| 18 | 47                | 40   | 32   | 18    |
| 19 | 53                | 46   | 37   | 21    |
| 20 | 60                | 52   | 43   | 26    |
| 25 | 100               | 89   | 76   | 51    |
| 30 | 151               | 137  | 120  | 86    |

Source: Eysenck, 2009 (adapted) (53)

**Table A 5. Critical values for the sign test**

| <b>n</b> | <b>Two-sided p-value</b> |             |             |             |              |
|----------|--------------------------|-------------|-------------|-------------|--------------|
|          | <b>0.10</b>              | <b>0.05</b> | <b>0.02</b> | <b>0.01</b> | <b>0.001</b> |
| 5        | 0                        | —           | —           | —           | —            |
| 6        | 0                        | 0           | —           | —           | —            |
| 7        | 0                        | 0           | 0           | —           | —            |
| 8        | 1                        | 0           | 0           | 0           | —            |
| 9        | 1                        | 1           | 0           | 0           | —            |
| 10       | 1                        | 1           | 0           | 0           | —            |
| 11       | 2                        | 1           | 1           | 0           | 0            |
| 12       | 2                        | 2           | 1           | 1           | 0            |
| 13       | 3                        | 2           | 1           | 1           | 0            |
| 14       | 3                        | 2           | 2           | 1           | 0            |
| 15       | 3                        | 3           | 2           | 2           | 1            |
| 16       | 4                        | 3           | 2           | 2           | 1            |
| 17       | 4                        | 4           | 3           | 2           | 1            |
| 18       | 5                        | 4           | 3           | 3           | 1            |
| 19       | 5                        | 4           | 4           | 3           | 2            |
| 20       | 5                        | 5           | 4           | 3           | 2            |
| 25       | 7                        | 7           | 6           | 5           | 4            |
| 30       | 10                       | 9           | 8           | 7           | 5            |
| 35       | 12                       | 11          | 10          | 9           | 7            |

**Source: Eysenck, 2009 (adapted) (53)**