



LISBON
SCHOOL OF
ECONOMICS &
MANAGEMENT
UNIVERSIDADE DE LISBOA

Master in Actuarial Science

Master's Final Work Internship Report

Health Insurance Pricing with Generalised Linear Models

Ana Beatriz Marques Cabral Valente

October - 2020

Master in
Actuarial Science

Master's Final Work
Internship Report

Health Insurance Pricing with Generalised Linear Models

Ana Beatriz Marques Cabral Valente

Supervision:

João Manuel de Sousa Andrade e Silva

Maria Isabel Teixeira da Silva Pimenta Ribeiro

October - 2020

PREVIEW

Acknowledgments

The realization of this internship report was a pleasant challenge and this accomplishment was only possible with all the support I received.

First of all, I would like to express my deepest appreciation to my thesis advisor, professor João Andrade e Silva, who provided valuable and constructive suggestions. His extensive knowledge and commitment were essential for the materialization of this report.

I would also like to acknowledge my colleagues at Allianz for their patient guidance and assistance. Despite all the barriers created by the pandemic situation and the home confinement, your support was fundamental and is greatly cherished. Special gratitude should be given to my company's supervisor, Isabel Ribeiro, for embracing this project by my side with an enthusiastic encouragement present at all times.

Finally, I must express my gratitude for the unconditional love of my family, friends and dogs. Thank you for always being there, your moral support and comfort throughout this confinement phase were necessary to maintain the focus and finish this report.

Abstract

Generalized Linear Models (GLMs) are being broadly used in the Non-Life Insurance Pricing. The premium charged by the insurance company is calculated based on a tariff. The most standard procedure to estimate the pure premium is by assuming that the claim counts and claim amounts are independent. From this independence, the claim frequency and severity can be forecasted by distinct GLMs and the Tariff is obtained by combining both models.

The present report gives a brief introduction on the methodology and describes how we prepared the data prior to the GLM application. The models obtained for the Stomatology Treatments and Appointments, one of the many coverages that can be included in a Health Insurance policy, are analyzed in this report. The SAS software was used to construct the datasets and to properly organize the data and R was the software used for the modelling process. Once the models were estimated, the pure premium was calculated and a tariff for the mentioned coverage was constructed.

Finally, we compared the results obtained by modelling the coverage in R with the output obtained by my colleagues, using the software implemented by the company. We conclude that both models are not significantly different, despite having some structural distinctions.

Keywords: Health Insurance, Insurance Pricing, Tariff, Generalized Linear Models, Claim Frequency, Claim Severity

Resumo

Os Modelos Lineares Generalizados (GLMs) são amplamente utilizados na precificação de seguros do ramo Não Vida. O prêmio cobrado pela seguradora é calculado com base em uma tarifa. A abordagem clássica para estimar o prêmio é feita assumindo a independência entre o número de sinistros e o seu custo. A partir desta independência, a frequência e a severidade dos sinistros são estimados através de GLMs separados e a tarifa é obtida combinando os dois modelos.

O presente relatório fornece uma breve introdução sobre a metodologia e descreve como preparámos os dados antes da aplicação do GLM. Os modelos obtidos para os Tratamentos e Consultas de Estomatologia, uma das muitas coberturas que podem ser incluídas numa apólice de Seguro Saúde, são analisados neste relatório. O software SAS foi utilizado para construir as bases de dados e para organizar adequadamente a informação e o software R foi utilizado para o processo de modelagem. Uma vez estimados os modelos, o prêmio puro foi calculado e a tarifa, para a cobertura mencionada, foi construída.

Por fim, comparámos os resultados obtidos em R com as conclusões obtidas pelos meus colegas, utilizando o software implementado pela empresa. Concluímos que ambos os modelos não são significativamente diferentes, apesar de apresentarem algumas distinções estruturais.

Palavras-Chave: Seguro de Saúde, Precificação de Seguros, Tarifa, Modelo Linear Generalizado, Frequência de Sinistros, Severidade de Sinistros

List of Tables

2.1	Equivalence between implementing the exposure in the offset or adding it as a weight variable, when modelling the number of claims and the claim frequency, respectively.	10
3.1	Summary of the number of observations per calendar year.	16
3.2	Absolute and relative frequency of the number of people, total cost and average claim cost, per claim number.	17
3.3	Description of the variables used in our GLMs.	20
4.1	Regression estimates of the Poisson GLM for the Claim Frequency in the training set.	24
4.2	Regression estimates of the Gamma GLM for the Claim Severity in the training set.	26
4.3	Tariff for the Stomatology Appointments and Treatments model.	28
5.1	Comparison between the coefficients obtained from the two different models.	30
5.2	The Risk premium obtained by using the models from R and Emblem. . .	31

List of Figures

4.1	Relative frequency of the total years of exposure (dark blue) and the total number of claims (lighter blue) distributed through the different levels of each relevant variable.	23
4.2	Claim cost histogram for the data without the zero claims.	25
4.3	The average claim cost (blue bars) and the total number of claims (stars) corresponding to each level of the relevant variables.	25
4.4	Graphical representation of the tariff coefficients from table 4.3.	28
A.1	Empirical claim frequencies of each original age level used for grouping the age levels in the Frequency Model.	34
A.2	Relative frequency of the total years of exposure (dark blue) and the total number of claims (lighter blue) distributed through the new levels obtained for the Frequency Model.	34
A.3	Empirical claim frequencies of each original age level used for grouping the age levels in the Severity Model.	34
A.4	The average claim cost (blue bars) and the total number of claims (stars) corresponding to the new levels obtained for the Severity Model.	35
A.5	The first two digits of the postal code associated with each district.	35

Contents

Acknowledgements	i
Abstract	ii
Resumo	iii
1 Introduction	1
1.1 Motivation and Goals	1
1.2 Context	2
1.3 Report Organization	3
2 Basic Concepts of Non-Life Insurance Pricing	4
2.1 Premium Concepts	4
2.2 Tariff	5
2.3 GLM Overview	6
2.4 Exposure, Offset and Weight	8
2.5 Goodness of fit	10
2.5.1 Log-Likelihood and Deviance	10
2.5.2 AIC and BIC	11
2.5.3 Test MSE	12
3 Data Processing	13
3.1 Main datasets	13
3.2 Partitioning the data	14
3.2.1 Train and Test	14
3.2.2 Cross Validation	15
3.3 Data Overview	16
3.4 Variables	17

3.4.1 Variables Selection 17

3.4.2 Grouping Categorical Variables 18

3.4.3 Variables Introduction 19

4 Models 21

4.1 Claim Frequency 21

4.2 Claim Severity 24

4.3 Pure Premium 27

5 Analysis of the Results and Further Developments 29

5.1 Analysis of the Results 29

5.2 Further Developments 32

Appendix A 34

PREVIEW