
Interactive R App for Carrying Out an Elicitation Process Applied to Prostate Cancer in Colombia

M.Sc. Andrés Felipe Flórez RIVERA (afflorezr@unal.edu.co)

Ph.D Juan Carlos Correa MORALES (jccorrea@unal.edu.co)

Escuela de Estadística, Universidad Nacional de Colombia, Sede Medellín.

Ph. D Manuel García FLÓREZ (garcia@usco.edu.co)

Facultad de Salud, Universidad Sur Colombiana, Neiva (Huila)

ABSTRACT

In this paper is presented an interactive R-application for carried out an elicitation process to the vector of parameters of the Multinomial distribution. The application is developed mainly under two R libraries, Shiny and RSQLite. Shiny is a package that allows developing a web application framework for R and RSQLite package embeds the SQLite database engine in R and provides an interface compliant with the DBI package. The application has two main components, the first stores the records of an elicitation, an expert and a variable. The second components allow that the analyst to ask the experts their judgments about previous variables stored, it return a plot allowing that the expert have feedback about elicited values. The application is tested in the estimation of the prevalence of each level Gleason score in patients who have been diagnosed with prostate cancer in Colombia.

Keywords: *Multinomial Distribution, Dirichlet Distribution, Shiny, RSQLite, Gleason.*

INTRODUCTION

The use of available information is one of the main strengths of the Bayesian statistics. Once this information is recorded in an appropriate form, it could be combined with information obtained in a formal sample process easily. In many occasions it is too costly to obtain information via a sample process and the only solution it is to elicit the information out of experts. The estimation of prior informative probability distributions is an important task, because when these are combined with the sampling distributions provide a better fit to the case study (O'Hagan 1998; Correa 2014).

Elicited probability distribution is commonly used as prior distribution in the Bayesian analysis where it represents the initial beliefs about the parameters of a model. Is common in the elicitation process involving a facilitator who helps the expert in the development of knowledge in a probabilistic way. The "expert" is the person who is being analyzed, this may be the scientist who collected the data and want to do

conclusions from them, or may be a decision maker that wants to make inferences on the basis of the available data and prior knowledge (Garthwaite et al., 2005).

The expert elicitation can also be seen as knowledge engineering, widely used in contexts where elicited distribution is not combined with the evidence of the data, since the expert opinion is essentially all the knowledge available. Unfortunately, the literature on expert elicitation is surprisingly small compared to the vast and extensive literature on Bayesian statistics in general (O'Hagan, 1998 and 2005). However at present the statisticians, especially Bayesian statistician, are building more and more complex model for real applications (O'Hagan, 2005).

With the increased use of computers in different areas of knowledge in recent years, seems natural that elicitation methods should also use it, due to that interactive computing is almost essential to carry out a process of elicitation, because usually, they require a sequence of questions that are determined by the above answer. Within of the benefits of interactive computing are included the ability to provide feedback and identify apparent inconsistencies in the expert judgments, the experts may be left alone to answer questions that the computer asks, saving safe and carefully all your answers. Computers are also useful when is used hypothetical future sampling technique, where the hypothetical values given are usually the result of a function of the values previously elicited, which the hypothetical sample can to be quickly calculate and display to the expert. However, the main advantage of the use of interactive computer is making use of graphical methods, which allows the analyst to provide most reliable expert's feedback. Although interactive computing and software packages can provide significant support, even many of the researchers are not aware of it (Jenkinson, 2005; Devilee and Knol, 2011).

A review of the software to elicit expert judgment is made by Devilee and Knol (2011), which, from experience with projects of the National Institute for Public Health and the Environment (RIVM) which concerned elicitation of expert was learned that the use of software was very useful, it allowed experts see the results and compare their own opinions with others. In its report 4 main phases are distinguished. Characterization of uncertainty, expert selection, design and implementation of elicitation session and packages for elicitation of expert groups:

- **Uncertainty characterization:** At this stage, seems that there is no software available, the only report was made by agency of the Netherlands for environmental assessment, which has developed a line of guidance for estimating and reporting uncertainty and has tools for characterizing uncertainty.
- **Selection of experts:** Devilee and Knol (2011) point out that the analyst can make use of bibliographic databases online (Pubmed, Scopus, Web of Science, Google scholar) as a first step, then conduct online surveys since they offer the possibility to include experts from different geographical locations, some of the packages that allow online surveys are: survey Monkey, Google Docs, Opinions Online and Lime survey.
- **Design and implementation of elicitation session:** Regarding to elicitation process design are highlighted the software to develop conceptual models like MS Visio, Mindjet MindManager, MatchWare MindView, iMindMap, Diagram designer, VUE, Freemind and Xmind. Regarding

to implementation of the elicitation process, the software that can benefit substantially, either graphical support or performing of the required calculations are ELI, Elicitor (v. 2010), Probes, ArcGIS customized, SL Gallery, Probes and Hypo.

- **Elicitation of expert groups:** In order that multiple experts can express their opinions in the form of brainstorming, voting or action plans have been created software that can help these activities: Delphi Blue, System Vanguardia, Facilitate Pro coffee, Thinktank 3.2 and Ynsyte WebIQ.

R IMPLEMENTATION

R is a language and environment for data manipulation, calculation and graphical display. It is a GNU project which is similar to the S language and environment which was developed at Bell Laboratories (formerly AT&T, now Lucent Technologies) by John Chambers and colleagues. R can be considered as a different implementation of S. There are some important differences, but much code written for S runs unaltered under R (R Core Team, 2015). Among other things it has:

- An effective data handling and storage facility.
- A suite of operators for calculations on arrays, in particular matrices,
- A large, coherent, integrated collection of intermediate tools for data analysis,
- Graphical facilities for data analysis and display either directly at the computer or on hardcopy, and
- A well developed, simple and effective programming language (called ‘S’) which includes conditionals, loops, user defined recursive functions and input and output facilities. (Indeed most of the system supplied functions are themselves written in the S language.)

The term “environment” is intended to characterize it as a fully planned and coherent system, rather than an incremental accretion of very specific and inflexible tools, as is frequently the case with other data analysis software. R is very much a vehicle for newly developing methods of interactive data analysis. It has developed rapidly, and has been extended by a large collection of packages. However, most programs written in R are essentially ephemeral, written for a single piece of data analysis (Venables et al, 2014; R Core Team, 2015). Making use of the advantages said above about using R in data analysis, an application for elicitation method proposed by Flórez (2015) is built, the application is developed mainly under two libraries R, Shiny and RSQLite.

Shiny

Shiny is an open source R package developed by RStudio e Inc that provides an elegant and powerful web framework for building web applications using R. Shiny helps you turn your analyses into interactive web applications without requiring HTML, CSS, or JavaScript knowledge (Although these languages can be used to expand, improve and visualize better the results of the app). Shiny makes this task a very easy

thing because it uses a reactive programming model to simplify the development of web applications. The models “reactive” link inputs, outputs and extensive pre-configured players that make it possible to build presentable, useful and powerful applications with minimal effort; additionally it is designed for fully interactive visualization, using JavaScript libraries like d3, Leaflet, and Google Charts (Chang W et al., 2015).

RSQLite

RSQLite embeds the SQLite relational database engine in R, providing a DBI-compliant interface. SQLite is a public-domain, single-user, very light-weight database engine that implements a decent subset of the SQL 92 standard, including the core table creation, updating, insertion, and selection operations, plus transaction management. Relational databases are ubiquitous. “Relational” means that instead of storing data in a “flat” table such as an Excel spreadsheet, data are stored in separate tables and then the separate tables are related to one another. In turn, this means that if the data are normalized, redundancy is reduced. “Embedded” means that the database engine has been designed to coexist inside other applications. This means that there is no negotiating access across networks, nor any administrative setting up (usernames, passwords, dns, etc). It means that the database engine is installed along with the application. This is the case with R. When the RSQLite package is installed, SQLite comes with it. There is no need for separate installation of an SQLite server. It also means that the data are stored in regular files and the files can be stored almost anywhere; they do not have to reside inside a server. (Wickham et al., 2015; Muspratt, 2012).

The APP

The app has 2 main modules registration and elicitation. In the **registration** module the analysts can to store the information about:

An elicitation:

View of how a new elicitation project is stored

Graphic 1

The screenshot shows a web application interface. On the left is a navigation sidebar with a menu structure: 'Introduction', 'Register' (expanded to show 'Elicitation', 'Expert', and 'Variable'), and 'Elicitation'. The main content area displays a form titled 'Register a new Elicitation'. The form contains two input fields: 'Elicitation ID:' and 'Elicitation Name:'. Below the input fields is a 'Send' button.

An Expert:

View of how a new expert is stored

Graphic 2

The screenshot shows a web application interface. On the left is a navigation sidebar with a tree view containing 'Introduction', 'Register' (expanded to show 'Elicitation', 'Expert', and 'Variable'), and 'Elicitation'. The main content area displays a form titled 'Register a new Expert'. The form has three input fields: 'Expert ID:', 'Expert Name:', and 'n Equivalent:'. A 'Send' button is at the bottom of the form.

Note that a new input is needed here; the analyst should store the expert's n-Equivalent. That is, for what sample size expert's knowledge is equivalent.

A Variable:

View of how a new variable is stored

Graphic 3

The screenshot shows a web application interface. On the left is a navigation sidebar with a tree view containing 'Introduction', 'Register' (expanded to show 'Elicitation', 'Expert', and 'Variable'), and 'Elicitation'. The main content area displays a form titled 'Register a new Variable'. The form has four input fields: 'Variable ID:', 'Variable Name:', 'Number of Categories:' (with the value '4' entered), and 'Name of each Category:' (with two sub-inputs). A 'Send' button is at the bottom of the form.

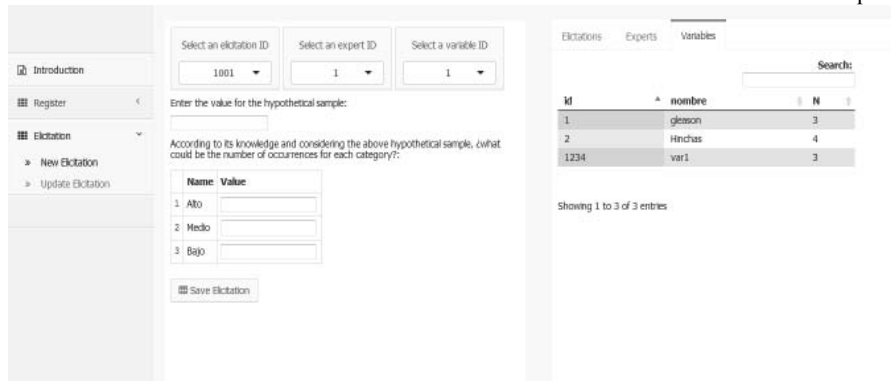
Here the analyst indicates how many levels have the variable and stored too the name of each of these.

In the **elicitation** module the analysts can to store and update the values elicited from the expert:

Stored elicited values:

View of how start with the elicitation process

Graphic 4



On the left side the analyst can stored the values elicited from the expert, he just needs to select at the top the id's values of elicitation, expert and variable then enter N hypothetical and elicited values. Note that on the right side the analyst can to see all elicitation, experts and variables stored.

Update elicited values:

View of how update the elicited values

Graphic 5



On the left side the analyst just needs to select at the top the id's values of elicitation, expert and variable to see the estimated values. On the right side the analyst can to see a plot of the distribution for each level of the elicited variable, the plot is used to give feedback to the expert so that if the expert does not feel comfortable with the results the analyst could to update the values until that the expert is agree with them. Finally the analyst can download the results.

THE APP USED IN THE CANCER PROBLEM

In Colombia cancer represents a growing public health problem: according to mortality ciphers of Globocan (2012), the country 104 people die every day from this disease. This result is a combination of factors including: bad habits of the population and consumption of snuff, overweight, low fruit and vegetable intake and excess fast foods, sugary drinks and alcohol. Another determining factor is related to the problems of health services in the areas of prevention, early detection, treatment and palliative care. The number of sick and dead by this cause has increased in recent years; In 2012 nearly 71,000 new cases and 38,000 deaths by this cause were reported, meaning that 195 people are diagnosed and 104 more die daily by this disease (Heidenreich et al., 2010).

Prostate cancer is the most common no cutaneous tumor in the man, it represents for the urologist a very common condition and its single mention will cause a number of questions both the patient and the doctor himself in relation to prognosis, stage, free survival of disease, free survival of biochemical relapse, which certainly could answer easily depending on the number of parameters to be considered in the structure of the responses. That is why for many years, urologists have made use of all those methods that could scientifically endorse the answers that support our therapeutic decisions where, within the main diagnostic tools used to obtain evidence of prostate cancer are included rectal examination, serum concentration of PSA and trans rectal ultrasound (ETR), although one of the most consistent prognostic factors or prognostic data has been the Gleason system (Heidenreich et al, 2010; Potenziani, 2014).

The Gleason score is a system used to measure the degree of aggressiveness of cancer based on microscopic observation of the characteristics that present the cells from a sample obtained in an organs biopsy. The procedure is to select two areas of the sample and assign each of them a number from 1 to 5. 1 corresponds to a well differentiated tumor and therefore unaggressive, 5 to a poorly differentiated and very aggressive tumor. The values between 2 and 4 are assigned to medium degrees of differentiation. Then are added the values obtained in both areas and is obtained a number comprised between 2 and 10. This value is the Gleason Score (Zollo, 2006; Gleason et al., 2002). The possible outcomes are:

- **Gleason score between 2 and 4:** Cancer with low invasiveness, slow growth and therefore better prognosis.
- **Gleason score from 5 to 7:** Cancer with an intermediate aggressiveness.
- **Gleason score from 8 to 10:** Cancer of high aggressiveness and poor prognosis.

In this application we want to estimate the prevalence of each level of the Gleason score in patients who have been diagnosed with prostate cancer. Knowing the distribution of the different levels of the Gleason score in patients already diagnosed with prostate cancer in Colombia is helpful in public health issues since the importance and significance of Gleason score in this pathology is invaluable, because it not only allows best diagnostic strategy, clinical and surgical, but also allow undoubtedly predict the evolution of the cases studied in time and the possibilities of biochemical relapse which result in strategies for associated treatments and higher alert for doctors responsible for these cases (Potenziani, 2012).

Methodology

As mentioned in Flórez (2015), to carry out elicitation process that satisfies a professional scrutiny addition to contextualize the expert on the subject of interest, is necessary to include too the following steps in the process:

- **Design and Validation of questions:** Since that the proposed method of elicitation uses a technique of indirect elicitation, the expert is asked directly by the number of patients (already diagnosed with prostate cancer) that he considered to be in each of the levels (2-4, 5-7, 8-10) of Gleason score given a hypothetical sample of patients. Before asking for the number of patients in each level Gleason score, the expert was given an introduction of the problem, where is validated that the question was understood correctly and that the expert in effect had an answer for each question.
- **Structuring, decomposition and training in probability:** the expert is given an introduction, initially on multinomial probability distribution and its relationship with Gleason score, then on elicitation of expert judgment and finally is specifies to the expert that accurately the uncertain quantity that will be estimated is the prevalence of each levels of the Gleason score in the population of patients diagnosed with prostate cancer and is determined that the measurement scale is the number of individuals in each level Score.
- **Application of the method**
 1. **n-equivalente estimation:** The analyst based on their knowledge about the expert selected estimates an N-Equivalent.
 2. **Estimation of relative frequency:** Through of questions like “Dr García, if are selected randomly 100 patients diagnosed with prostate cancer in Colombia, how many of them according a its knowledge, are at level 1 Gleason score (Gleason 2 to 4) how many level 2 (Gleason between 5 and 7) and finally how many level 3 (Gleason 8 to 10) This process is repeated 6 times, where the expert is asked to distribute 6 different hypothetical samples in the three Gleason Score levels. Once the expert carried out the distribution of the 6 samples, the relative frequency for each level is calculated. If the frequencies are closer proceeds with the simulation stage taking as relative frequency the average of the six estimated frequencies. If the estimated frequencies are divergent we proceed to validate the consistency of these estimates with the expert (imitating the Delphi method), the 6 hypothetical samples estimated are shown to expert and again the expert is asked to redistribute the hypothetical samples until that the frequencies estimated are as closely as possible.
 3. **Simulation and Parameter Estimation of Vector:** Once the consistency of the estimates made by the expert is validated, was conducted the simulation process and the estimation of the vector parameters. For this process is only necessary to enter the estimated values within the application in R, which is responsible for making the simulation and deliver the outputs of alphas estimated.

Results

Hypothetical samples elicited were:

Hypothetical samples and elicited values

Table 1

Sample	Estimated Values		
	Level 1	Level 2	Level 3
50	10	25	15
150	35	80	35
250	80	150	20
340	80	180	80
510	90	300	120
713	150	400	163

The estimates obtained for the Dirichlet distribution are:

Alpha estimated values

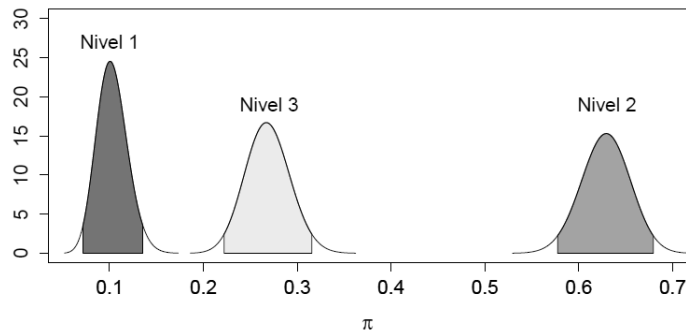
Table 2

	Alpha	Mean	Variance
α_1	141.64	0.221	0.0002
α_2	417.74	0.652	0.0003
α_3	81.186	0.126	0.0001

The marginal probability distribution for each level of Gleason Score:

Score of Gleason distribution

Graphic 6



Level 1= low, Level 2 = medium and Level 3 = high.

Where probability intervals with a region 100 (1- α)% of credibility and a confidence level of $\alpha = 0.05$ for each level of the Score are given by:

IP for each Score Gleason level

Table 3

Level	LI	LS
1	0.073	0.137
2	0.577	0.678
3	0.223	0.316

CONCLUSIONS

In the graph 6 we can see that for level 1 Gleason score greater density is concentrated around 0.1, for level 2 is concentrated around 0.62 and finally for level 3 the highest density is found around 0.28. The fact that patients with minor proportion are the Level 1 gives an indication that Colombia early detection of cancer is a latent problem, on the other hand, the result which about 62% of patients diagnosed are in level 2 is an input to validate whether public health policies are oriented to treating patients with this level of aggression and whether clinics and hospitals are ready to serve them, it is worth noting that the treatment of prostate cancer, even in illness clinically localized, is increasingly complex due to the different treatment options available that present an equivalent oncologic efficacy, but secondary effects associated with significantly different treatment (Heidenreich et al., 2010), by this, to know if clinics and hospitals are prepared to provide adequate treatment for patients with Gleason score at Level 2 can generate a significant impact in the population diagnosed with the disease. Finally we found that a substantial proportion is at level 3, this ratio requires more aggressive treatments so as well as to patients who are at level 2 worth validate whether the orientation of public health policies, the clinics and hospitals are prepared to provide timely care and treatments for these patients, since the treatment decisions at every level of Gleason score should be based on clinical guidelines, clearly indicating the one used in the making process decisions, which is further recommended a multidisciplinary approach (Heidenreich et al., 2010).

Is worth noting that if we compare the elicitation procedure applied here with select a representative random sample of the population of patients diagnosed with prostate cancer and doing a biopsy on them to determine the level of Gleason score, we can see that the process elicitation is relatively cheaper and faster. Now if there are sample data of the distribution of Gleason score, the distribution estimated by the method of elicitation can also be updated by this data, so that the result is a distribution closer to the real.

Undoubtedly this implementation is only the beginning of the way that seek to integrate different methods of elicitation with the set integrated of statistical tools that brings the software R, although previously Goulet et al. (2009) and O'Hagan and Oakley (2010) have used R in elicitation processes, they have used it only to make some calculations and not to capture information, store information and feedback to the experts in real time. At present, the application allows to carry out multiple processes elicitation for different experts and different variables of research; although only are implemented the elicitation method proposed by Flórez (2015) the final goal is to implement the elicitation methods of Judgments of experts most used. The application is now in its beta

version and currently functions of other packages like shinydashboard are being adding to support the visualization and implementation of other elicitation methods within of it.

REFERENCES

1. Chang W., Cheng J., Allaire JJ., Xie Y and McPherson J., 2015, shiny: Web Application Framework for R: R package version 0.11.1.9002, <http://shiny.rstudio.com>.
2. Chang W., 2015, shinydashboard: Shiny Dashboard, R package version 0.2.2, <https://github.com/rstudio/shinydashboard>.
3. Correa J. C., 2014, Elementos de Estadística Bayesiana.
4. Devilee J., Knol A., 2011, Software to support expert elicitation: An exploratory study of existing software packages, National Institute for Public Health and the Environment, Ministry of Health, Welfare and Sport, RIVM Letter Report 630003001.
5. Flórez A., 2015, Elicitación de una distribución subjetiva del vector de parámetros π de la distribución Multinomial, M.sc. tesis, Universidad Nacional de Colombia.
6. Garthwaite P., Kadane J. y O'Hagan A., 2005, Statistical Methods for Eliciting Probability Distributions, Journal of the American Statistical Association, Vol. 100, 680-700, No.470.
7. Gleason D., Mellinger G. and The Veretans Administration Cooperative Urological Research Group, 2002, Prediction of prognosis for prostatic adenocarcinoma by combined histological grading and clinical staging., http://globocan.iarc.fr/Pages/fact_sheets_cancer.aspx, Vol. 167, 953-958.
8. Globocan, 2012, Estimated Cancer Incidence, Mortality and Prevalence Worldwide., The Journal of Urology, update at: 18/08/2014.
9. Goulet V., Jacques M. and Pigeon M., 2009, Expert: Modeling Without Data Using Expert Opinion, The R Journal, Vol. 1/1.
10. Heidenreich A., Bolla M., Joniau S., Mason M., Matveev V., Mottet N., Schmid H., vander T., Wiegel T. and Zattoni F., 2010, Guía clínica sobre el cáncer de próstata, European Association of Urology, <http://www.uroweb.org/gls/pdf/spanish/01-%20GUIA%20CLINICA%20SOBRE%20EL%20CANCER%20DE%20PROSTATA.pdf>, update at: 18/08/2014.
11. Jenkinson D., 2005, The Elicitation of Probabilities - A Review of the Statistical Literature.
12. Muspratt S., 2012, R and SQLite: Part 1, <http://sandymuspratt.blogspot.com/2012/11/r-and-sqlite-part-1.html>
13. O'Hagan A., 1998, Eliciting Expert Beliefs in Substantial Practical Applications, Journal of the Royal Statistical Society, Series D, Vol. 47, 21-35, No. 1.
14. O'Hagan A., 2005, Research in Elicitation, University of Sheffield, UK.
15. O'Hagan A. and Oakley J., 2010, SHELF: the Sheffield Elicitation Framework (version 2.0), School of Mathematics and Statistics, University of Sheffield, <http://tonyohagan.co.uk/shelf/> update at: April 18 of 2014)
16. Potenziari J., 2014, Significado del Grado de Gleason y del Score de Gleason en pacientes con Cáncer Prostático., National Academy of Medicine, http://www.researchgate.net/publication/237050093_Significado_del_Grado_de_Gleason_y_del_Score_de_Gleason_en_pacientes_con_Cncer_Prosttico/, update at: 20/08/2014.
18. R Core Team, 2015, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, <http://www.R-project.org>, update at: 05/03/2015.
19. Venables W., Smith D. and R Core Team., 2014, Notes on R: A Programming Environment for Data Analysis and Graphics Version 3.1.1., <http://cran.r-project.org/doc/manuals/r-release/R-intro.pdf>, update: August 17 de 2014 at 13:45.
20. Wickham H, James D, Falcon S, Healy L, 2015, RSQLite: SQLite Interface for R: R package version 1.0.0, <http://CRAN.R-project.org/package=RSQLite>.
21. Zollo A., 2006, Medicina interna. Secretos, Elsevier España S.A., cuarta edición.