

RESEARCH

Association of variants in gene coding regions with clinical data in colombian patients using data mining

Jennifer Velez^{1*†} and Elizabeth Leon²

Abstract

Background:

The need to understand the biological processes that are involved in different diseases, from the available biological data such as genomic sequences, microarrays, protein interactions, biomedical images among others and the rapid adoption of electronic medical records that provides an opportunity for large-scale research. Therefore, data mining techniques for the discovery of knowledge from obtaining information from different sources are increasingly important in biological and medical research.

Results: A group model was implemented for 228 patients, and they were associated with the variants for 4813 genes, obtaining 5 groups with their options available in the rules. As an analysis, an analysis of the CFTR the gene was also carried out by means of association rules and previously obtained groups. This is the search tag the measurements of the frequencies of the population were made in terms of the number of variants present the age, sex, type of variant and the allelic state of the variants. It was found that for the CFTR hay gene without pathogenic variants in the sampled population. A board also created to visualize the groups and the necessary rules for group and a database for variants in exons for Colombian patients.

Conclusions: Data mining techniques applied to disk support allow an inference of genetics structure the Colombian population and the epidemiological follow-up of the variants and their possible effects in patient's phenotypes.

Keywords: Variant; data mining; clustering; association rules

Content

This paper is organized in section Background, Results ,Discussion, Conclusions and Methods.

Background

Biological data mining (seen from bioinformatics) is the process of extracting new knowledge (previously unknown) from biological data, this also allows the use of concepts of data mining and automatic learning in theories and applications in research Biological, by deeding the data that are used to be applied, are genomes that come from DNA sequencing, transcriptome sequences that are RNA or proteins that come from inferences and experimental data from chemistry [1].

Inferences regarding large amounts of genomic data require analysis of computational tools to interpret data, being one of the most active areas where data mining is used (I understand data mining as the method of extracting information through learning automatic, statistics, artificial intelligence, recognition and visualization patterns) to solve biological problems, some examples where mining techniques have been applied is the classification of genes, the analysis of mutations in cancer and gene expression [2].

Clustering techniques of differentially expressed genes have also been applied, vector support machines have been used to associate the interactions between genes and generate biological networks, as well as traditional methodologies for data mining are not precise or efficient and they require new algorithms to be developed and methodologies that respond in a more precise way to a biological question [3]. Without forgetting that it is necessary to evaluate the available platforms,

*Correspondence: jevelezse@unal.edu.co

¹School of Engineering, Universidad Nacional de Colombia, Bogotá D.C., Colombia

Full list of author information is available at the end of the article

[†]Equal contributor

the technological tools that allow the implementation of processes that associate data with research and obtain more generalized results. This should apply to the research requirements to ensure successful implementation [3].

Some of the data mining tasks are: 1. Classification: where the data is classified to a predefined class, 2. Association: see elements that are associated by rules, 3. Grouping or grouping: as the definition of a population of data within a subgroup or group [2].

The use of high performance sequencing techniques together with the application of data mining can contribute to the diagnosis of complex diseases such as cancer [4, 5].

Results

Exploratory analysis

The exploratory analysis of the information contained in the database was carried out. A sample of 250 patients donated by the Genetix SA laboratory was taken, of which only 228 had the informed consent to use the information for research purposes.

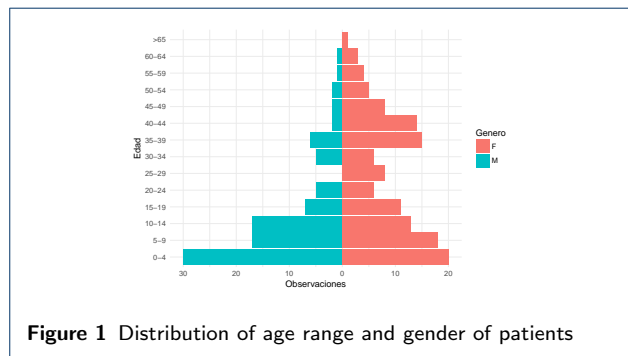


Figure 1 Distribution of age range and gender of patients

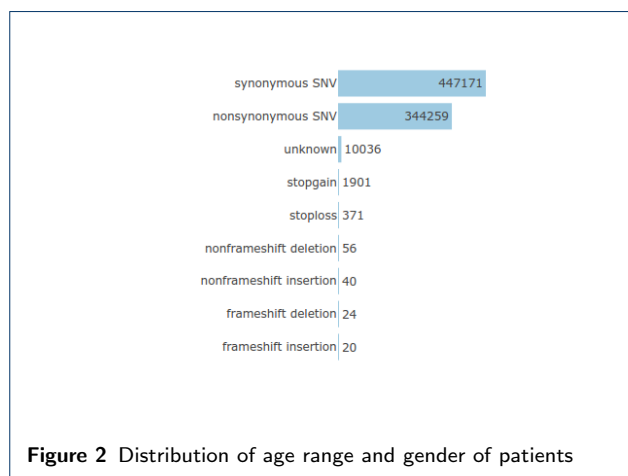


Figure 2 Distribution of age range and gender of patients

The database contains 228 patients of which 133 are female and have a total of 468,485 variants and 95

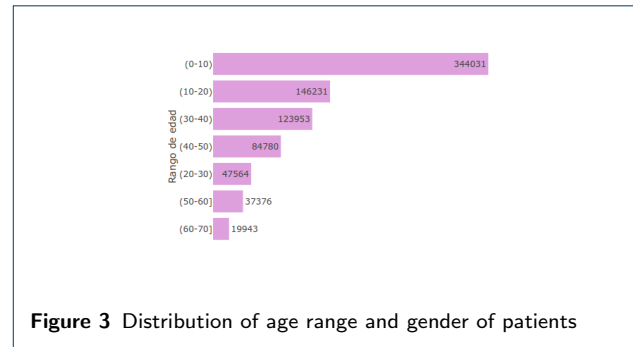


Figure 3 Distribution of age range and gender of patients

male patients with 345,239 variants obtaining a total of 803,878 variants. Figure 1 represents the distribution of patients by age range and Figure 2 represents the distribution of variants according to their type. In Figure 2 shows the number of variants that are synonymous and not synonymous, being the most frequent in the population, at a global level it is known that these are the most frequent types of variants [6].

The unknown variants are the third type of variant more frequent given that there is still the problem of selecting the transcript to perform the appropriate nomenclature of the variants, so the annotator reports that they are unknown [7]. Figure 3 shows the distribution of the variants identified according to the age range, the range with the highest number of variants being patients between the ages of 0 to 10 years, given that it is the most represented population within the database.

The allelic state of the variants (zygosity) found within the database are divided into 458639 heterozygotes corresponding to 57.05% of the total of the variants and homozygous 345239 corresponding to 42.95%. The distribution of the zygosity of the variants can be explained from the error that can be generated in the identification of the variants given that during the call of variants it is possible that a homozygous variant is classified as heterozygous, if during the sequencing process it is identified wrongly nucleotides [8, 9].

Textual analysis of clinical information.

The results obtained for the frequency of words were breast, cancer, syndrome, suspicion and years. In figure 4 shows the first 30 most frequent words and the word cloud of all the documents.

The frequency of words shows us the main characteristics of the clinical information, the words cancer and breast are the main phenotypes, the word syndrome is also found that can be associated with different diseases and the word suspicion refers to ambiguous diagnoses that patients may have, One of the contributions of sequencing is that based on the phenotype can help a diagnosis, between different symptoms and syndromes that can be applied to rare and complex diseases [10].

4. Hannah-Shmouni, F., Seidelmann, S.B., Sirrs, S., Mani, A., Jacoby, D.: The Genetic Challenges and Opportunities in Advanced Heart Failure. *Canadian Journal of Cardiology* **31**(11), 1338–1350 (2015). doi:10.1016/j.cjca.2015.07.735

5. Kawashima, K.: Text Mining and Pattern Clustering for Relation Extraction of Breast Cancer and Related Genes, 1–5 (2017)

6. Fu, W., O'Connor, T.D., Jun, G., Kang, H.M., Abecasis, G., Leal, S.M., Gabriel, S., Altshuler, D., Shendure, J., Nickerson, D.A., Bamshad, M.J., Akey, J.M.: Analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants. *Nature* **493**(7431), 216–220 (2013). doi:10.1038/nature11690

7. McCarthy, D.J., Humburg, P., Kanapin, A., Rivas, M.A., Gaulton, K., Cazier, J.B., Donnelly, P.: Choice of transcripts and software has a large effect on variant annotation. *Genome Medicine* **6**(3) (2014). doi:10.1186/gm543

8. Babraham Bioinformatics: FASTQC manual (2016). [http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/Help/3 Analysis Modules/](http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/) Accessed 2016-06-25

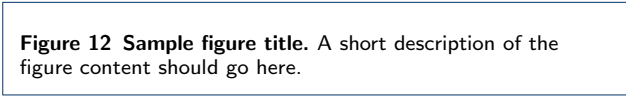
9. Pirooznia, M., Kramer, M., Parla, J., Goes, F.S., Potash, J.B., McCombie, W.R., Zandi, P.P.: Validation and assessment of variant calling pipelines for next-generation sequencing. *Human genomics* **8**(1), 14 (2014). doi:10.1186/1479-7364-8-14

10. Tetreault, M., Bareke, E., Nadaf, J., Alirezaie, N., Majewski, J.: Whole-exome sequencing as a diagnostic tool: Current challenges and future opportunities. *Informa Healthcare* (2015). doi:10.1586/14737159.2015.1039516

11. Salton, G., Buckley, C.: Term-Weighting Approaches in Automatic Text Retrieval. *Information Processing & Management* **24**(5), 513–523 (1988)

12. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, É.: Scikit-learn: Machine Learning in Python. *Journal of Machine Learning* **12**, 2825–2830 (2011). doi:10.1007/s13398-014-0173-7.2. arXiv:1201.0490v2

Figures



Tables

Table 1 Sample table title. This is where the description of the table should go.

	B1	B2	B3
A1	0.1	0.2	0.3
A2
A3

Additional Files

Additional file 1 — Sample additional file title
Additional file descriptions text (including details of how to view the file, if it is in a non-standard format or the file extension). This might refer to a multi-page table or a figure.

Additional file 2 — Sample additional file title
Additional file descriptions text.