Title

Anonymous Author(s)

Affiliation Address email

Abstract

- The abstract paragraph should be indented ½ inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points.

 The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.
- 5 1 Introduction
- Rare diseases have different definitions from countries. In Europe, a disease or disorder has less than 1 affection in 2,000 people called rare. In USA, a disease affects less than 200,000 Americans defined as rare [1]. Unlike influence of a single rare disease, the rare diseases in a whole affect around 8 6%-8% of the general population and the 50% of patients are children [1, 2]. Based on the rarity 9 of these diseases, a general physician may not likely have a single patient in his/her whole career 10 [3]. Furthermore, the diagnosis of these diseases are relatively hard, because the symptoms often 11 vary between individual causes [4]. Therefore, it is important to refer patients to an appropriate 12 healthcare expert who are familiar with the specific disease and symptom set. In order to create a good match between rare disease patients and experts, we created a novel machine learning application by 14 implementing SVM regression on OMIM publication data [5, 6]. By using this machine learning 15 approach, we distinguished 2,1224 experts over 209,110 people with publications in 1,292 diseases. 16

17 2 Data Processing

- In order to evaluate if a person is an expert or not, we have to establish a metrics to measure expertise first. In our project, the metrics are based on people?s publications. On OMIM, which is an online catalog of rare diseases and disorders, each rare disease or symptom has a specific OMIM ID [6]. Under each OMIM ID, related research or clinical publications are included in the Reference Section.
- 22 We scraped the publication data by using OMIM API and stored them into Python dictionaries [6, 7].
- 23 Figure 1: Flow of Experiment Process
- 24 The original OMIM data are not formatted. In order to use OMIM data, we have to reorganized them
- 25 into a data structure which is a dictionary in this case and extract useful information. The useful
- 26 information includes OMIM ID, disease name, publication list. Each publication contains a unique
- 27 PubMed ID and we stored publication information, such as paper title, author names and journal
- 28 name, under each PubMed ID. In the part of author name processing, we found GeneReviews author
- 29 names had a different format from OMIM author names. Therefore, GeneReviews author names had
- 30 been parsed into the same format as OMIM author names which only include full last name, first
- name initial and middle name initial.
- Features are the keys for drawing classification boundaries or fitting regression lines [8, 9]. In our
- 33 algorithm, there were 8 features had been chosen: number of publications on the disease1, normalized
- number of publications on the disease2, number of diseases that the author has publication on3,
- number of publications as first author on the disease4, number of publications as last author on the

- disease 5 and the number of publications on the disease in 3 years 6, 5 years 7 or 10 years 8. Feature 1
- 37 represents the person?s research ability on the disease in general. Feature 2 shows a relative research
- 38 ability compared with other researchers on the same disease. Feature 3 reveals if a person?s research
- 39 interest focus on one specific disease or for many diseases in general. Feature 4 and 5 provide weights
- to the impact from authorship. The last three features represent the effect from time of publications.
- 41 After extracting features from the meta data, mapping from GeneReviews to OMIM ID has been
- 42 made. GeneReviews is an authorized resource which provides relative information for inherited
- 43 conditions. The articles are written by one or more expert on the specific condition or disease [10].
- 44 We used GeneReviews expert data as our positive dataset. We filtered people with no publications,
- because our analysis was based on publications. After mapping, there were 2,160 positive data points
- 46 and 206,950 unknown data points. We randomly marked 2,160 unknown points as negative. Then
- we used both positive and negative data to train different learning algorithms, which include SVM
- classification [11], random forest [12] and naive Bayes [13], and compared the results. We used 10
- times tenfold cross-validation to make the results more accurate [14].

50 3 Logistic Regression

51 4 Experiments and Results

52 5 Discussion

53 6 Conclusion

7 General formatting instructions

- 55 The text must be confined within a rectangle 5.5 inches (33 picas) wide and 9 inches (54 picas) long.
- The left margin is 1.5 inch (9 picas). Use 10 point type with a vertical spacing (leading) of 11 points.
- 57 Times New Roman is the preferred typeface throughout, and will be selected for you by default.
- Paragraphs are separated by ½ line space (5.5 points), with no indentation.
- 59 The paper title should be 17 point, initial caps/lower case, bold, centered between two horizontal
- 60 rules. The top rule should be 4 points thick and the bottom rule should be 1 point thick. Allow 1/4 inch
- space above and below the title to rules. All pages should start at 1 inch (6 picas) from the top of the
- 62 page.
- 63 For the final version, authors' names are set in boldface, and each name is centered above the
- 64 corresponding address. The lead author's name is to be listed first (left-most), and the co-authors'
- 65 names (if different address) are set to follow. If there is only one co-author, list both author and
- 66 co-author side by side.
- 67 Please pay special attention to the instructions in Section 9 regarding figures, tables, acknowledgments,
- 68 and references.

9 8 Headings: first level

- All headings should be lower case (except for first word and proper nouns), flush left, and bold.
- 71 First-level headings should be in 12-point type.

72 8.1 Headings: second level

Second-level headings should be in 10-point type.

74 8.1.1 Headings: third level

Third-level headings should be in 10-point type.

76 **Paragraphs** There is also a \paragraph command available, which sets the heading in bold, flush

177 left, and inline with the text, with the heading followed by 1 em of space.

78 9 Citations, figures, tables, references

79 These instructions apply to everyone.

9.1 Citations within the text

- 81 The natbib package will be loaded for you by default. Citations may be author/year or numeric, as
- long as you maintain internal consistency. As to the format of the references themselves, any style is
- 83 acceptable as long as it is used consistently.
- 84 The documentation for natbib may be found at
- http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf
- 86 Of note is the command \citet, which produces citations appropriate for use in inline text. For example,
- 88 \citet{hasselmo} investigated\dots
- 89 produces
- 90 Hasselmo, et al. (1995) investigated...
- 91 If you wish to load the natbib package with options, you may add the following before loading the 92 nips_2016 package:
- 93 \PassOptionsToPackage{options}{natbib}
- 94 If natbib clashes with another package you load, you can add the optional argument nonatbib95 when loading the style file:
- 96 \usepackage[nonatbib]{nips_2016}
- 97 As submission is double blind, refer to your own published work in the third person. That is, use "In
- 98 the previous work of Jones et al. [4]," not "In our previous work [4]." If you cite your other papers
- 99 that are not widely available (e.g., a journal paper under review), use anonymous author names in the
- citation, e.g., an author of the form "A. Anonymous."

9.2 Footnotes

101

106

- Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number 1
- in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote
- with a horizontal rule of 2 inches (12 picas).
- Note that footnotes are properly typeset *after* punctuation marks.²

9.3 Figures

- All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction.
- The figure number and caption always appear after the figure. Place one line space before the figure
- caption and one line space after the figure. The figure caption should be lower case (except for first
- word and proper nouns); figures are numbered consecutively.
- You may use color figures. However, it is best for the figure captions and the paper body to be legible
- if the paper is printed in either black/white or in color.

¹Sample of the first footnote.

²As in this example.

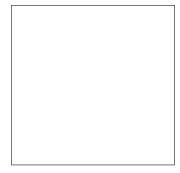


Figure 1: Sample figure caption.

Table 1: Sample table title

	Part	
Name	Description	Size (μm)
Dendrite Axon Soma	Input terminal Output terminal Cell body	~ 100 ~ 10 up to 10^6

113 **9.4 Tables**

121

131

132

133

134

135

136

137

All tables must be centered, neat, clean and legible. The table number and title always appear before the table. See Table 1.

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

Note that publication-quality tables *do not contain vertical rules*. We strongly suggest the use of the booktabs package, which allows for typesetting high-quality, professional tables:

https://www.ctan.org/pkg/booktabs

This package was used to typeset Table 1.

123 10 Final instructions

Do not change any aspects of the formatting parameters in the style files. In particular, do not modify the width or length of the rectangle the text should fit into, and do not change font sizes (except perhaps in the **References** section; see below). Please note that pages should be numbered.

127 11 Preparing PDF files

Please prepare submission files with paper size "US Letter," and not, for example, "A4."

Fonts were the main cause of problems in the past years. Your PDF file must only contain Type 1 or Embedded TrueType fonts. Here are a few instructions to achieve this.

- You should directly generate PDF files using pdflatex.
- You can check which fonts a PDF files uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program pdffonts which comes with xpdf and is available out-of-the-box on most Linux machines.
- The IEEE has recommendations for generating PDF files whose fonts are also acceptable for NIPS. Please see http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf

- xfig "patterned" shapes are implemented with bitmap fonts. Use "solid" shapes instead.
- The \bbold package almost always uses bitmap fonts. You should use the equivalent AMS Fonts:

141 \usepackage{amsfonts}

followed by, e.g., \mathbb{R} , \mathbb{N} , or \mathbb{C} . You can also use the following workaround for reals, natural and complex:

144 \newcommand{\RR}{I\!\!R} %real numbers

145 \newcommand{\Nat}{I\!\!N} %natural numbers

146 \newcommand{\CC}{I\!\!!\!C} %complex numbers

Note that amsforts is automatically loaded by the amssymb package.

148 If your file contains type 3 fonts or non embedded TrueType fonts, we will ask you to fix it.

149 11.1 Margins in L^AT_EX

- 150 Most of the margin problems come from figures positioned by hand using \special or other
- commands. We suggest using the command \includegraphics from the graphicx package.
- 152 Always specify the figure width as a multiple of the line width as in the example below:

```
\usepackage[pdftex]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.pdf}
```

- See Section 4.4 in the graphics bundle documentation (http://mirrors.ctan.org/macros/
- 156 latex/required/graphics/grfguide.pdf)
- A number of width problems arise when LATEX cannot properly hyphenate a line. Please give LaTeX
- 158 hyphenation hints using the \- command when necessary.

159 Acknowledgments

- 160 Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end
- of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper.

162 References

- References follow the acknowledgments. Use unnumbered first-level heading for the references. Any
- choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font
- size to small (9 point) when listing the references. Remember that you can use a ninth page as
- long as it contains only cited references.
- 167 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In
- 168 G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), Advances in Neural Information Processing Systems 7, pp.
- 169 609–616. Cambridge, MA: MIT Press.
- 170 [2] Bower, J.M. & Beeman, D. (1995) The Book of GENESIS: Exploring Realistic Neural Models with the
- 171 GEneral NEural SImulation System. New York: TELOS/Springer-Verlag.
- 172 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent
- 173 synapses and cholinergic modulation in rat hippocampal region CA3. Journal of Neuroscience 15(7):5249-5262.