

Automatically Generating Color Palettes for Categorical Values

Chinmay Kulkarni

Stanford University HCI Group
Computer Science Department
Stanford, CA 94305
chinmay@cs.stanford.edu

Julie Fortuna

Stanford University HCI Group
Computer Science Department
Stanford, CA 94305
jfortuna@stanford.edu

ABSTRACT

Is it possible to automatically create a color palette relevant to a topic? Could such a palette be used to guide color choices while visualizing data? We envision a tool that automatically creates aesthetically pleasing and topic-relevant palettes for a large class of topics. In order to do this, we must first extract palettes from color pixel values of images from Google Images via clustering and topic models.

ACM Classification: H5.2 [Information interfaces and presentation]: User Interfaces. - Graphical user interfaces.

General terms: Design, Human Factors, Experimentation

Keywords: Information visualization, colors, crowdsourcing, user study

INTRODUCTION

blah blah blah

RELATED WORK

Prior work exists on automatic creation of color palettes. This work falls broadly in two categories. The first focuses on finding representative colors from images, that can be used as color palettes. The most recent of these is [8]. This line of research has so far focused only on extracting colors from a single image. This project extends this work by extracting colors from multiple, related images. I believe that some of the techniques used by [8], such as a weighted histogram that uses color saturation and neighborhood color coherence, can be adapted for multiple images too. Depending on constraints of time, I plan to explore some of these techniques.

The second category of research on palette generation focuses on optimizing visual properties, such as color saliency and perceptual color distance, both manual or rule-based, as pioneered by Brewer [3]; and with varying degrees of automation [6, 9]. I believe most such optimization research

is complementary to this project, and can be used as a post-extraction step to optimize the colors chosen. Statistical work on color saliency is valuable, even if it hasn't been directly applied as an optimization objective; color saliency in the context outside data-visualization in [4, 1].

Topic models have been shown to be effective in information retrieval. Latent semantic analysis (and later, LDA), for instance, has been used to find "latent" similarities between concepts [5, 2]. Similar similarity-measures have been computed for nodes in a graph [7]. While these similarity measures may help to better cluster color-values, they don't target the domain of color recommendations directly.

SYSTEM DESIGN

Query System

Google Images is queried for images from the category.

Statistical summarization

We assume that the images from the category are a random sampling from the concept-space of the category. Taking this assumption further, we look at the *average* frequencies of the different colors as a metric of how concepts are shared across the values in a category.

$$Old = \alpha * average + (1 - \alpha) * new \quad (1)$$

$$new = \frac{(old - \alpha * average)}{(1 - \alpha)} \quad (2)$$

Since we are interested in the colors specific to a category value, we subtract a fraction of the average color frequency.

Clustering

We cluster the result to get relevant colors in LAB space. We found that low saturation colors are less likely to be relevant, so we reweight more saturated colors to be more relevant.

SYSTEM EVALUATION

We evaluated the system on three related metrics: the likability of the generated color palettes, how topic-relevant the palettes were perceived to be, and how the colors in the palette affect understanding of the data they represent. For all three metrics, the algorithmically generated palettes were compared against a randomly generated palette, and

one generated by experts. For the likability and understanding metrics, the random palette was chosen from the set of palettes generated for other topics by our system. This was to ensure that only the relevance, not the base quality of the colors was considered. For all topics tested, we limited the number of specific items represented in the palette to four. This also allowed us to compare the algorithmically and randomly generated palettes to the randomly generated palettes. We ran a small laboratory study of X participants recruited through school mailing lists, in addition to a large-scale crowdsourced study on Amazon's Mechanical Turk.

Likability

To measure likability, the automatically, expert, and randomly generated color palettes for a given topic are presented in a random order. Participants rate each palette on a seven-point Likert scale based on how much they like each palette for a given topic.

Relevance

For relevance, an association task is used: given a topic (e.g. "US Politics") and one of the topic terms (e.g. "Democrat"), the participant chooses which color, among a set of displayed swatches, is relevant to it.

Understanding

For understanding, users will be shown differently-colored infographics, and participants will be timed while they answer conceptual questions related to the infographic. Since the three metrics may interact strongly, they will be studied in a within-subjects design.

RESULTS

blah blah

DISCUSSION

blah

REFERENCES

1. R. Benavente, F. Tous, R. Baldrich, and M. Vanrell. Statistical modelling of a colour naming space. In *Proceedings of the 1st European Conference on Colour in Graphics, Imaging, and Vision (CGIV2002)*, pages 406–411.
2. D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. *The Journal of Machine Learning Research*, 3:993–1022, 2003.
3. C. Brewer. Color use guidelines for data representation. In *Proceedings of the Section on Statistical Graphics, American Statistical Association*, pages 55–60, 1999.
4. J. Chuang, M. Stone, P. Hanrahan, and S. Consulting. A Probabilistic Model of the Categorical Association Between Colors. 2008.
5. S. Dumais, G. Furnas, T. Landauer, S. Deerwester, and R. Harshman. Using latent semantic analysis to improve access to textual information. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 281–285. ACM, 1988.
6. C. Healey. Choosing effective colours for data visualization. In *Visualization'96. Proceedings.*, pages 263–270. IEEE, 1996.
7. G. Jeh and J. Widom. SimRank: a measure of structural-context similarity. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 538–543. ACM, 2002.
8. B. Morse, D. Thornton, Q. Xia, and J. Uibel. Image-Based Color Schemes. In *Image Processing, 2007. ICIP 2007. IEEE International Conference on*, volume 3, pages III–497. IEEE, 2007.
9. A. Zeileis, K. Hornik, and P. Murrell. Escaping RGB-land: Selecting colors for statistical graphics. *Computational Statistics & Data Analysis*, 53:3259–3270, 2009.