

A Novel Framework for Molecular Co-Expression Pattern Analysis in Multi-Channel Toponome Fluorescence Images

Ahmad Humayun¹, Shan-e-Ahmed Raza¹, Christine Waddington², Sylvie Abouna³,
Michael Khan³, Nasir M. Rajpoot^{1†}

¹Department of Computer Science, University of Warwick, UK

²Molecular Organisation & Assembly in Cells (MOAC), University of Warwick, UK

³School of Life Sciences, University of Warwick, UK

Abstract—Bioimage computing is rapidly emerging as an important area in image based systems biology with an emphasis on spatiotemporal localization of subcellular bio-molecules, most importantly proteins. A key problem in this domain is analysis of protein co-localization or co-expression of protein molecules. Imaging techniques, such as the Toponome Imaging System (TIS) [1], with the ability to localize several different proteins in the same tissue specimen are only becoming available recently. Traditional co-localization studies and some of the modern co-expression studies have serious limitations when analyzing this kind of data. Here we present a framework for the analysis of molecular co-expression patterns (MCEPs) in TIS image data.

Index Terms—Protein co-localization, molecular co-expression, multi-fluorescence imaging, bioimage computing, clustering analysis, MELC, TIS

I. INTRODUCTION

In recent years, bioimage computing is emerging as a cornerstone of hypothesis-driven research in life sciences with an emphasis on spatiotemporal localization [2]. A major focus in the post-genomic era is on analyzing subcellular protein patterns using the knowledge of spatiotemporal distribution of key proteins expressed in a given cell type [3]. As proximity of proteins located within similar compartments of a cell provides a powerful surrogate for functional complexes, functional studies involve proteins which appear to be key players in multiple cancer specific complexes. Imaging techniques with the ability to localize several different proteins in the same tissue specimen such as the Toponome Imaging System (TIS) [1], MALDI imaging [4], Raman spectroscopy [5], or multi-spectral imaging methods [6] are only becoming available recently. Of these, TIS is an automated fluorescence technique shown to have the ability to co-map hundreds of different proteins or other TAG-recognizable bio-molecules on a single tissue section. This results in a multi-tag fluorescence image stack with phase contrast images before and after incubation of the corresponding antibodies. For each antibody tag, four images of the tissue section are obtained, resulting in a stack of images for selected visual fields in a tissue specimen. First, an initial phase image is captured followed by a fluorescence image taken under ultraviolet light. After a waiting time for tag

attachment and specimen rinsing, another phase image is taken with its matching fluorescence image. Having taken these 4 images, a bleaching cycle occurs before the process can begin for the next tag. Before any significant conclusions can be drawn about co-localization of proteins, the image stack must be accurately aligned or registered. The phase images can be used for alignment purposes since (a) these are not expected to vary throughout a TIS run, and (b) each phase image is taken a fraction of a second before its matching fluorescence image and so it is assumed that the phase/fluorescence image pair for a particular antibody is perfectly aligned.

Traditional protein co-localization studies involve three key steps: staining the tissue specimen with two or three different dyes (such as DAPI for nuclei and a green or red fluorescent proteins, also known as GFP or RFP, binding to specific antigens), taking fluorescence images of each dye with different wavelengths of the incident laser beam, and treating the individual fluorescence channels as one of the Red, Green, or Blue channels to construct a color image. In this color image, pixel locations with yellow color, for instance, indicate the simultaneous presence of bio-molecules corresponding to red and green channels. A major shortcoming of this method of studying protein co-localization is that it does not take into account the varying levels of expression of different proteins. Another disadvantage is that often this kind of simplistic analysis is limited to the study of co-localization of 2-3 proteins at the same time.

A few such studies discuss measuring the quality of co-localization in an image pair, which is made non-trivial by the presence of structures with multiple morphologies. Bolte and Cordelières [7] and Zinchuk *et al.* [8], [9] survey many such measures including Pearson's and Manders' correlation coefficients. Interestingly, to align the phase images, we also need to employ a similar kind of measure. We resort to mutual information which has the advantage of not being limited to parametric linear relationships between pixels unlike the Pearson's coefficient [10].

Multi-tag bioimaging methods, such as the ones mentioned above, allow us to study combinatorial protein patterns in specific types of tissues and to characterize and differentiate different kinds of cells within the same specimen. However,

† Corresponding author: Nasir M. Rajpoot, email: nasir@des.warwick.ac.uk

recent developments in machine learning and computer vision are yet to impart their influence on the data analysis pipeline of standard software used with most of these imaging systems. For instance, much of the reported analysis of TIS image data such as [13] is based on binarization using manually selected thresholds.

In this paper, we describe a framework to pre-process and analyze multi-channel fluorescence microscopy images obtained using TIS [1]. The framework is composed of three key components: pre-processing for robust alignment of TIS images using a modified version of the recently proposed RAMTaB algorithm [14], segmentation of nuclei, and clustering analysis of protein patterns without binarizing individual fluorescence images. Avoiding binarization in the last step allows analysis and discovery of combinatorial patterns of a whole range of protein expression levels or molecular co-expression patterns (MCEPs). The proposed framework is generic in nature and should be applicable to other multi-channel imaging methods such as MALDI [4], Raman spectroscopy [5], or multi-spectral imaging methods [6].

II. MATERIALS & METHODS

The image data obtained in this study was acquired using a TIS microscope installed at the University of Warwick. The human colon tissues were surgically removed from cancerous patients. For each cancerous tumor, distal normal section was also taken from the same patient's colon. Patient consent and appropriate ethics approval were obtained to remove and handle these tissues for research. A library of 26 antibody tags, some of which are known tumor markers and others cancer stem cell markers, were used based on a previous study [13].

The molecular co-expression analysis framework proposed in this paper is based on three major components: aligning the images corresponding to individual antibodies in a TIS stack with each other, segmentation of nuclei and surrounding pixels to consider only molecular expression in the cellular areas, and analysis & visualization of molecular patterns using a clustering method.

A. Image Registration

We employ the RAMTaB algorithm [14] for aligning TIS images. Apart from performing registration based on mutual information, this method is capable of finding the optimal reference image¹ from a TIS stack while giving a confidence value over the estimated alignment. However, a limitation of this algorithm is that it does not explicitly deal with out-of-focus phase images. It is noteworthy that the TIS microscope autofocuses before capturing the first image and uses this focus over the whole run. Due to this autofocus feature, if there is variation between focusing planes of the visual fields this could cause a drift of the focal plane as the run progresses. This can manifest itself in the phase images either as out-of-focus/blurry images, or as focusing at a different level causing differences in locations of cell walls. Furthermore, where the plane of focus changes, some cells appear smaller and some larger. This sort of focal error is difficult for the RAMTaB registration algorithm to deal with, as it is not a simple shift or shear combination, and so needs to be avoided. Therefore, the sample needs to be very flat on the cover slip and as thin as possible. Choosing visual fields with similar z coordinates for the plane of focus may be effective in reducing focusing errors.

Experiments have shown that such images can be aligned by using recent image deblurring techniques which use a normalized sparsity measure [12]. An example of an out-of-focus phase image is shown in Figure 1b - a corresponding reference phase image with the right focus is shown in 1a. Figure 1c shows deblurred version of 1b using a standard blind deconvolution technique [11], and 1d shows deblurred image with the help of blind deconvolution using normalized sparsity measure. The blind deconvolution algorithm minimizes the scale-invariant cost function, l_1/l_2 norm ratio to estimate the kernel blur. The kernel is estimated in a multi-scale approach

¹An optimal reference image is one which minimizes the total amount of missing pixel intensities due to removal of pixels near image boundaries during the alignment.

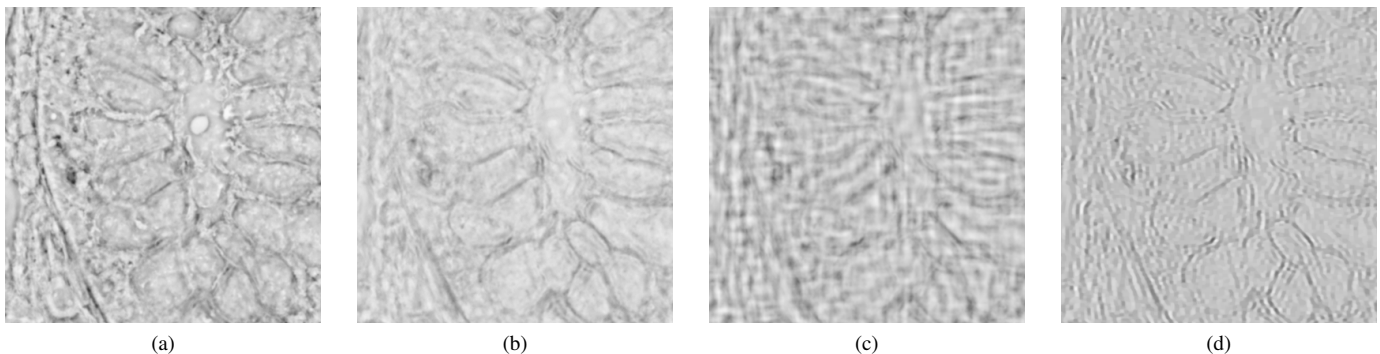


Fig. 1: Phase image 1a shows a reference image with normal focus, phase image 1b shows a floating image with out-of-focus capture during the imaging process. Image in 1c shows a deblurred version of 1b using the standard blind MATLAB® deconvolution of Holmes *et al.* [11], while 1d shows the result of blind deconvolution using normalized sparsity measure [12]. All phase images are shown as their complements here for the clarity of display.

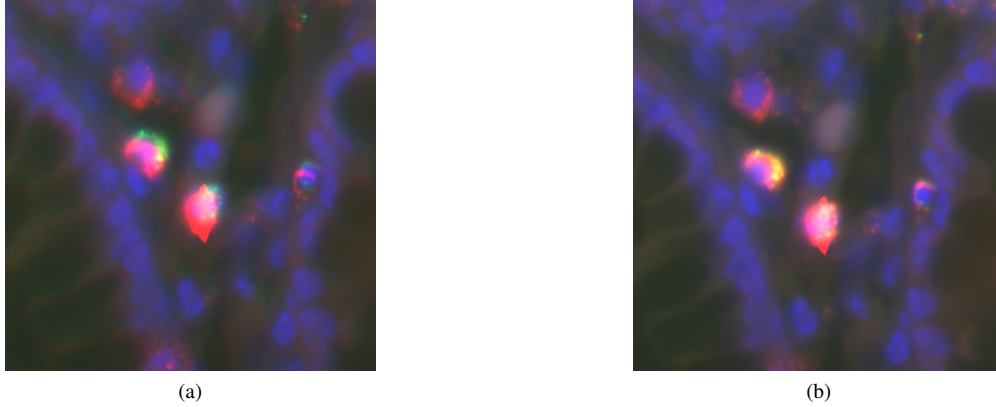


Fig. 2: Pseudo-color images generated by taking fluorescence images for three antibodies Muc1, CD57, and DAPI as red, green, and blue channels: 2a before registration 2b after registration. The co-localization of Muc1 and CD57 in compartments of one of the cells is much more obvious as a yellow in 2b.

from coarse to fine image resolutions. Once the kernel is estimated, the image is deblurred using the method proposed by Krishnan and Fergus [15]. Images in Figure 1a and 1d can be aligned using the RAMTaB algorithm [14]. Figure 2 demonstrates the effectiveness of alignment using this method.

B. Nuclei Segmentation

Once the image data from a TIS multi-fluorescence image stack have been aligned, we normalize the intensity values in each of the aligned TIS images to the range $[0, 1]$. The aligned DAPI channel is then segmented in order to extract pixel locations corresponding to the cell nuclei and their immediate neighborhood only. This step ensures that only molecular patterns localized to cell nuclei and cytoplasm are considered. This removes signal from stroma and lumen in the case of colon, for example, which may add noise to the process of pattern analysis. This segmentation of pixels into nuclei and their immediate neighborhood is achieved using Gaussian mixture modeling (GMM) over the normalized intensity values of the fluorescence channel images. Bayesian information criterion (BIC) [16], [17] is used for selecting the model with the number of clusters K taken from the set $\{2, 4\}$.

C. Clustering Analysis for Extracting Molecular Co-Expression Patterns (MCEPs)

After registration and segmentation of the stack of TIS images, we obtain protein expression vectors \mathbf{T}_i of length L at each segmented pixel location i . Since each dimension in this vector encodes the difference in expression levels after adding a particular antibody, the vector can be used to cluster pixels based on responses to L antibodies. Our goal is to study how well unsupervised clustering can reveal the differences in inter and intra-tissue antibody responses. Our approach is based on a simple hierarchical clustering method, which is a bottom-up clustering method [18]. It starts with each pixel as a cluster and iteratively merges these clusters to form bigger ones. Existing

clusters are merged to create new ones, reducing the number of clusters by 1 at each iteration until there is only one cluster containing all the data points. For instance, if initially there are N pixels (and an equal number of clusters), the first iteration merges two pixels to give $N - 1$ clusters. This process can be better represented as a dendrogram tree structure, where cutting across the tree at level k would give $N - k$ clusters². Like many other unsupervised methods, hierarchical clustering can also be provided with the number of clusters desired. We aim for $C = 20$ clusters, which we will call molecular co-expression patterns (MCEPs). These C clusters are produced by cutting the tree at level $k = N - C$.

As mentioned above, two clusters are merged at each iteration. The criterion we employ to select these two clusters aims to minimize the increase in the variance of clusters [19]. Mathematically, at each tree level k , we have clusters $S_j = \{\mathbf{T}_1, \dots, \mathbf{T}_{n_j}\}$ where $n_j = |S_j|$ and $j \in \{1, \dots, N - k\}$. Here, we can define the within-class variance of cluster S_j as follows:

$$\sigma(S_j) = \sum_{m=1}^{n_j} (\mathbf{T}_m - \bar{S}_j) (\mathbf{T}_m - \bar{S}_j)^T \quad (1)$$

where \bar{S}_j is the centroid vector for cluster S_j . To make clusters for level $k + 1$, we seek to combine vectors in $S_{\hat{u}}$ and $S_{\hat{v}}$ such that:

$$\hat{u}, \hat{v} = \underset{u, v \in \{1, \dots, N - k\}}{\operatorname{argmin}} [\sigma(S_u \cup S_v) - \sigma(S_u) - \sigma(S_v)] \quad (2)$$

or,

$$\hat{u}, \hat{v} = \underset{u, v \in \{1, \dots, N - k\}}{\operatorname{argmin}} \frac{n_u n_v (\|\bar{S}_u - \bar{S}_v\|_2)^2}{n_u + n_v} \quad (3)$$

where $\|\cdot\|_2$ is the Euclidean norm. This step will result in a new cluster $(S_{\hat{u}} \cup S_{\hat{v}})$ formed by merging $S_{\hat{u}}$ and $S_{\hat{v}}$, hence reducing the number of clusters by 1.

²The algorithm starts at level $k = 0$, where there are N clusters. Cutting the tree at level k means truncating the tree after level k .

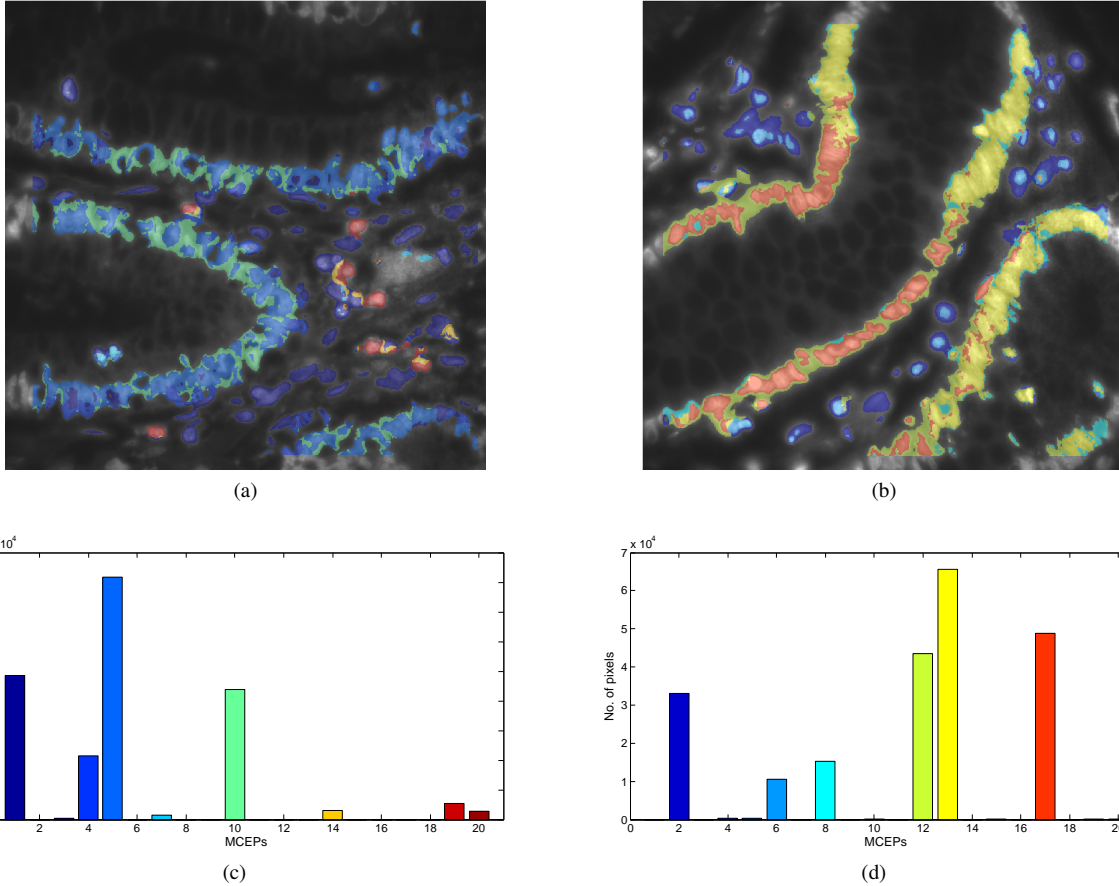


Fig. 3: Pseudo-color overlay of molecular co-expression patterns (MCEPs) on corresponding phase contrast images of two human colon tissue specimens (cancer a and normal b) using the centroids of $C = 20$ clusters. The bar charts in 3c and 3d show a pixel frequency histogram of the 20 MCEPs found in the specimens.

III. EXPERIMENTAL RESULTS

Using the above clustering method, we pick $C = 20$ clusters localized to nuclei and their vicinities. Each of the centroids of these clusters is given a unique color. We use the MATLAB® `jet` colormap, a variation of the `hsv` colormap, which varies from dark blue (for the first MCEP) to dark red (for the last MCEP) passing through the colors cyan, yellow, and orange in between. A pseudo-color overlay of MCEPs on corresponding phase contrast images using the centroids of the 20 clusters for two human colon tissue specimens (cancer on the Left and normal tissue on the Right) is shown in Figure 3. It can be seen from this display of molecular co-expression patterns that there is a clear difference in tissue morphology and molecular expression at sub-cellular level in normal and cancer specimens. This approach is fundamentally different to the standard TIS visualization approaches using thresholds and random colors [1], [13]. Furthermore, by localizing the pattern analysis to DAPI-positive pixels and their surroundings, we are able to filter out any noise due to non-cellular pixel locations such as lumen or stroma.

IV. CONCLUSIONS

In this paper, we have presented a novel framework for the analysis of molecular co-expression patterns in multi-tag fluorescence image stacks generated by the TIS microscope. The framework should be applicable to other multi-tag imaging systems and we hope that it will serve as a critical building block for further analysis of TIS stacks in cancer studies.

V. ACKNOWLEDGEMENTS

This work was partly funded by the Warwick Institute of Advanced Studies (IAS) and the HDC. The authors are grateful to W. Schubert, the inventor of TIS, who helped us establish a TIS machine at University of Warwick, and members of his team at ToposNomos and the University of Magdeburg, especially A. Krusche and R. Hillert. The authors would also like to thank D. B. A. Epstein for many fruitful discussions. Special thanks go to S. Bhattacharya for contributions to design of the antibody library.

REFERENCES

- [1] W. Schubert, B. Bonnekoh, A.J. Pommer, L. Philipsen, R. Bockelmann, Y. Malykh, H. Gollnick, M. Friedenberger, M. Bode, and A.W.M.

- Dress, "Analyzing proteome topology and function by automated multidimensional fluorescence microscopy," *Nature Biotechnology*, vol. 24, no. 10, pp. 1270–1278, 2006.
- [2] S.G. Megason and S.E. Fraser, "Imaging in Systems Biology," *Cell*, vol. 130, no. 5, pp. 784–795, 2007.
- [3] R.F. Murphy, "Putting proteins on the map," *Nature Biotechnology*, vol. 24, no. 10, pp. 1223–1224, 2006.
- [4] D.S. Cornett, M.L. Reyzer, P. Chaurand, and R.M. Caprioli, "MALDI imaging mass spectrometry: molecular snapshots of biochemical systems," *Nature Methods*, vol. 4, no. 10, pp. 828–833, 2007.
- [5] H.J. Van Manen, Y.M. Kraan, D. Roos, and C. Otto, "Single-cell Raman and fluorescence microscopy reveal the association of lipid bodies with phagosomes in leukocytes," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 102, no. 29, pp. 10159–10164, 2005.
- [6] E. Barash, S. Dinn, C. Sevinsky, and F. Ginty, "Multiplexed Analysis of Proteins in Tissue Using Multispectral Fluorescence Imaging," *IEEE Transactions on Medical Imaging*, vol. 29, no. 8, pp. 1457–1462, 2010.
- [7] S. Bolte and F.P. Cordelières, "A guided tour into subcellular colocalization analysis in light microscopy," *Journal of Microscopy*, vol. 224, no. 3, pp. 213–232, 2006.
- [8] V. Zinchuk, O. Zinchuk, and T. Okada, "Quantitative Colocalization Analysis of Multicolor Confocal Immunofluorescence Microscopy Images: Pushing Pixels to Explore Biological Phenomena," *Acta Histochemica et Cytochemica*, vol. 40, no. 4, pp. 101–111, 2007.
- [9] V. Zinchuk and O. Zinchuk, "Recent advances in quantitative colocalization analysis: Focus on neuroscience," *Progress in Histochemistry and Cytochemistry*, vol. 44, no. 3, pp. 125–172, 2009.
- [10] J. Numata, O. Ebenhöf, and E.W. Knapp, "Measuring correlations in metabolomic networks with mutual information," *Genome Informatics*, vol. 20, pp. 112–122, 2008.
- [11] T.J. Holmes, S. Bhattacharyya, J.A. Cooper, D. Hanzel, V. Krishnamurthi, W. Lin, B. Roysam, D.H. Szarowski, and J.N. Turner, "Light Microscopic Images Reconstructed by Maximum Likelihood," in *Handbook of Biological Confocal Microscopy*, pp. 389–402. Plenum Press, New York, 1995.
- [12] D. Krishnan, T. Tay, and R. Fergus, "Blind Deconvolution using a Normalized Sparsity Measure," in *Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition (CVPR '11)*, 2011, pp. 233–240.
- [13] S. Bhattacharya, G. Mathew, E. Ruban, D.A.B. Epstein, A. Krusche, R. Hillert, W. Schubert, and M. Khan, "Toponome Imaging System: In Situ Protein Network Mapping in Normal and Cancerous Colon from the Same Patient Reveals More than Five-Thousand Cancer Specific Protein Clusters and Their Subcellular Annotation by Using a Three Symbol Code," *Journal of Proteome Research*, vol. 9, no. 12, pp. 6112–6125, 2010.
- [14] S. Raza, A. Humayun, S. Abouna, T.W. Nattkemper, D.B.A. Epstein, M. Khan, and N.M. Rajpoot, "RAMTaB: Robust Alignment of Multi-Tag Bioimages," Submitted to BMC Biophysics, 20xx.
- [15] D. Krishnan and R. Fergus, "Fast image deconvolution using hyper-laplacian priors," *Advances in Neural Information Processing Systems*, vol. 22, pp. 1033–1041, 2009.
- [16] C. Fraley and A.E. Raftery, "How Many Clusters? Which Clustering Method? Answers Via Model-Based Cluster Analysis," *The Computer Journal*, vol. 41, no. 8, pp. 578–588, 1998.
- [17] N.M. Rajpoot and M. Arif, "Unsupervised Shape Clustering using Diffusion Maps," *The Annals of the BMVA*, vol. 2008, no. 5, 2008.
- [18] A.K. Jain, M.N. Murty, and P.J. Flynn, "Data Clustering: A Review," *ACM Computing Surveys (CSUR)*, vol. 31, no. 3, pp. 264–323, 1999.
- [19] J.H. Ward, "Hierarchical Grouping to Optimize an Objective Function," *Journal of the American Statistical Association*, vol. 58, no. 301, pp. 236–244, 1963.