

A Graphical Model Approach for Matching Partial Signatures

Xianzhi Du¹, David Doermann¹ Wael AbdAlmageed²

¹UMIACS, University of Maryland, College Park. ²Information Sciences Institute, University of Southern California.

For signature matching we distinguish between two problems—*full signature matching*, which assumes an accurate segmentation and an author who has produced a complete and consistent signature, and *partial signature matching*, where we enroll a full reference signature, but may only have a partial signature or even initials to match against the reference signature. A lot of work has shown great promise for addressing the full signature matching problem using feature descriptors such as shape context [5] [2]. However, partial signature matching remains an open problem.

To address the partial signature matching problem, we developed a method based on the combination of supervised latent Dirichlet allocation (sLDA) [1] and hierarchical Dirichlet processes (HDP) [3]. In our approach, sLDA is first used to discover the salient regions in all training signatures. A salient region is a distribution over the features in the visual vocabulary, which groups similar co-occurring observations. Each author is modeled as a combination of all salient regions with different proportions. For a query signature, classification is performed by computing the salient region proportions for the signature based on observations. Further, instead of guessing the number of salient regions empirically, HDP is used to estimate the number needed for the given dataset.

First, we model signatures as a group of observations. Shape context features, which describes the relations between nearby points while tolerating slight shape distortion, are used to describe the 2-D binary shapes. To build observations, we first extract contour points from one signature proportional to the total length of the contour. For each contour point, a log-polar space is formed around it with uniform bins. A histogram is built by calculating the number of nearby points that fall in each bin in a certain order, based on the relative distance and angle of the two points. K-means clustering is used to build the visual vocabulary for all signatures. We cluster all contour points into different clusters and treat each cluster label as one observation. The vocabulary consists of all the cluster labels. For each signature, the number of contour points being classified into one cluster is regarded as the appearance frequency of this observation.

Then, we build a supervised topic model for the partial signature matching problem. The generative process for the n_{th} observation in t_{th} signature is given as follows:

1. For the t_{th} signature, draw salient region proportions θ_t from $Dir(\alpha)$
2. For each observation:
 - (a) Draw a salient region assignment $S_{t,n}$ from $Mult(\theta_t)$
 - (b) Draw an observation $O_{t,n}$ from $Mult(\beta_{S_{t,n}})$
3. Draw authorship variable A_t from $N(\eta^T \bar{S}_t, \sigma^2)$

Where the $Dir(\cdot)$, $Mult(\cdot)$, $N(\cdot)$ represent Dirichlet distribution, Multinomial distribution, and Normal distribution respectively. α is an R -dimensional hyperparameter for Dirichlet distribution with R being the number of salient regions. $\beta = [\beta_1, \beta_2, \dots, \beta_R]$, where each β_r is the distribution of salient region r over the vocabulary, and \bar{S}_t is the mean of the salient regions of the t_{th} signature. With only observations and authorship given, we want to estimate $\alpha, \beta, \eta, \sigma^2$. Variational EM algorithm [4] is used to estimate the parameters of the sLDA model.

To efficiently solve the problem of estimating the number of salient regions, we further use a new topic model structure called hierarchical Dirichlet processes provided by Teh [3], which lets the data estimate the number of salient regions needed.

Our method is tested on two partial signature datasets and one full signature dataset. The top-N rank accuracy is used as the evaluation protocol. To compare with previous methods, our goal is to show that our method

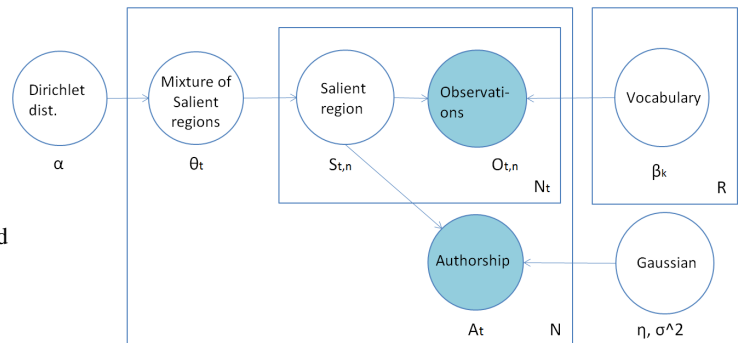


Figure 1: SLDA model for our problem.

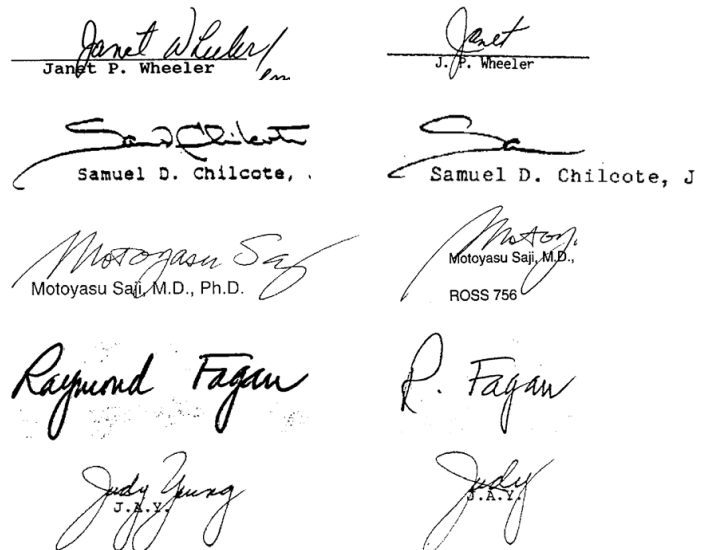


Figure 2: Sample signatures in DS-II partial dataset. Left column shows five full signatures. Right column shows their partial signatures with different kinds of degradations.

works well on both full and partial signatures. Some sample signatures are shown in Figure 2.

The first set of experiments was designed to test only partial signatures. It obtained a top-1 accuracy of 87.8%, a great improvement of 70.3% over the multi-stage LSH method. The second set of experiments was designed to test only full signatures. It obtained a top-1 accuracy of 92.1%, a slight improvement of 4.7% over the previous method. The third set of experiments tested on a more challenging dataset with both full and partial signatures. It obtained a top-1 accuracy of 37.8%, a moderate improvement of 15.6% over the previous method.

- [1] D. Blei and J. McAuliffe. Supervised topic models. *In NIPS*, 2008.
- [2] X. Du, W. AbdAlmageed, and D. Doermann. Large-scale signature matching using multi-stage hashing. *In ICDAR*, pages 976–980, 2013.
- [3] Y. W. Teh, M. Jordan, M. Beal, and D. Blei. Hierarchical dirichlet processes. *Journal of the American Statistical Association*, pages 1566–1581, 2006.
- [4] C. Wang, D. Blei, and F. Li. Simultaneous image classification and annotation. *In CVPR*, 2009.
- [5] G. Zhu, Y. Zheng, and D. Doermann. Signature detection and matching for document image retrieval. *IEEE Trans. PAMI*, 31(11):2015–2031, 2009.