Interactive Image Search for Clothing Recommendation

Zhengzhong Zhou, Yifei Xu, Jingjin Zhou and Liqing Zhang

MOE-Microsoft Key Laboratory for Intelligent Computing and Intelligent Systems Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shang, China {tczhouzz, fei960922, zhoujingjin}@sjtu.edu.cn, zhang-lq@cs.sjtu.edu.cn

Abstract

This demo delivers a novel retrieval system which meets users' multi-dimensional requirements in clothing image search. In this system, users are able to use both image and keywords as query inputs. We employ the color, texture, shape and attributes as additional descriptors to further refine the requirements. We propose the Hybrid Topic (HT) model, a probabilistic network integrating the multi-channel descriptors into a unified framework, to learn the intricate semantic representation of the descriptors above. The proposed model provides an effective multi-modal representation of clothes. Our experiments show that the HT method significantly outperforms the CNN-based deep search methods.

Introduction

HT Model for the Second-level Abstraction. As noted, we present a novel topic model, *i.e.* Hybrid Topic (HT) to integrate both visual and attribute information of clothing images. Fig. 2(d) shows the graphical model of HT. Images are represented as a list of "phrases". A phrase is a set of codewords w^j in different descriptor channels $(j \in \{t(TEXT), h(HOG), l(LBP), c(COL)\})$. For each image, we draw a topic mixture proportion θ to describe its high-level concept. For each phrase, we assign a topic z to determine its meaning. We assume that w^j s and z follow the Multinomial distribution with prior φ^j s and θ . While φ^j and θ follow the Dirichlet distribution with prior β^j s and α . The above parameters of HT model are estimated by Gibbs sampling. To achieve an effective image representation, we train 17 HT models (one for each local region) on the dataset, and denote the topic proportion



Using image retrieval technology to search for clothing products provides great convenience in apparel e-commerce [7, 5, 6]. When a user gravitates to someone's attire in a photo, he can simply submit the photo to the retrieval system to get a group of garments similar to the photo for selection. However, there still exist two major problems. The first one is "semantic gap". Due to the intraclass variation and inter-class ambiguity in image space, low-level visual features could not describe the semantic concept like category or style explicitly. The second one is "multi-dimensional requirements". Since the appropriate query images are not always available, users tend to further refine their demands by modifying the color, texture, shape and attribute descriptors of the original queries.



Figure 1: Illustration of System Interface.

In this demo, we present a novel interactive retrieval system (see Fig. 1) to tackle both problems utilizing topic model [1, 2, 9]. Our feature extraction procedure is described as a two-level abstraction. For the first level abstraction, we extract multiple clothing descriptors from images and convert them into the bag of words (BOW) formats. The BOW format enables us to use topic model for semantic analysis. Meanwhile, it benefits us to modify the local details of image descriptions. For the second level abstraction, we present a probabilistic network to combine the multi-channel descriptors into a unified framework, *i.e.* Hybrid Topic (HT) model. In this model, each image is depicted as a linear combination of hybrid topics, while each topic is described as distributions over various codewords. The HT model learns the intricate relationship among different descriptors, which can derive an effective representation of clothing images. We develop an interactive retrieval system on the server with 2 Intel Xeon 2.6GHz processors and 128GB memory. It contains about 120,000 upper-body clothing images crawled from online shopping platform JD.com. Experiments show that our system not only achieves better performance than the CNN-based deep search method [6], but also provides a more personalized interaction for online customers. The average retrieval time for a query is within 1 second. θ s as local descriptions of images. Thus, we concatenate the θ s extracted from 17 local regions of a clothing image as its retrieval feature.

Image annotation. To further facilitate the users, we automatically add tags to the query image in our retrieval system. The HT model is able to evaluate the missing values of attribute descriptors by the Bayesian reasoning. Meanwhile, users could modify these recommended tags to confirm their requirements.

Experiments

To evaluate the performance of our approach, we realize a CNN based method called Deep Search [6] as the baseline. We invite 10 subjects to perform 100 retrieval tasks in our retrieval system. For each task, we upload a clothing image and modify its color, shape, texture or attributes to refine the requirement. We measure the performance by Normalized Discounted Cumulative Cain: $NDCG@k = \sum_{j=1}^{k} \frac{2^{r(j)}-1}{\log(j+1)}$, where r(j) = 1 iff the j^{th} returned image satisfies the subject's demand, otherwise 0.

	NDCG@20
Deep Search Hybrid Topic	0.42 0.66
- –	1

 Table 1: The Retrieval Accuracy of HT.

amples of our retrieved results.



Original Query





Our experiment result is shown in Table 1. As we see,

the hybrid topic method improves the precision by 24%

compared with the Deep Search. It integrates the de-

mand of modification into image representation, which

reduces the distances between a query image and user

desired images in feature space. Fig. 3 shows some ex-

FRAMEWORK

Fig. 2 shows the flowchart of our system. Given a query image and corresponding search conditions, the region of clothes is cropped using object detection algorithm [10, 8] as a preprocessing stage. We extract multiple visual descriptors from the detected region and quantize them into BOW representation. Furthermore, we adjust the codeword weights of the modified feature to meet users' requirement. Using the pre-trained hybrid topic model, we integrate the above descriptors into effective retrieval features and infer the attributes of query image to improve the user interaction. Finally, we calculate the similarity between query image and images in database. Those top-ranked images are returned as the retrieval result.

The First-level Abstraction. In the pre-processing stage, we adopt Faster R-CNN [8] to localize the torso and sleeve regions in clothing images. We resize the torso region to 192×128 and use two slide windows (128×64 and 64×128) to sample it (stride 16). In this way, we split the image into 17 local regions (2 sleeve regions, 1 torso region and 14 overlapping regions cropped from the torso region). For each local region, we extract 3 types of visual descriptors: The histogram of oriented gradients (HOG), Local binary patterns (LBP), and color histogram (COL). We also extract the keywords from the text description as descriptor (TEXT). Using sparse coding, we convert all these descriptors into BOW formats, so that we can simply modify the codeword weights of descriptors instead of image contents.



Figure 3: Several Groups of User Requirements and Corresponding Results.

References

- [1] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. In NIPS, pages 601–608, 2001.
- [2] W. Chong, D. Blei, and F.-F. Li. Simultaneous image classification and annotation. In *CVPR*, pages 1903–1910, 2009.
- [3] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, volume 1, pages 886–893, 2005.
- [4] R. Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 2015.
- [5] J. Huang, R. S. Feris, Q. Chen, and S. Yan. Cross-domain image retrieval with a dual attributeaware ranking network. In *ICCV*, pages 1062–1070, 2015.



Figure 2: The Framework of Our System.

- [6] J. Huang, W. Xia, and S. Yan. Deep search with attribute-aware deep network. In ACM MM, pages 731–732, 2014.
- [7] S. Liu, Z. Song, G. Liu, C. Xu, H. Lu, and S. Yan. Street-to-shop: Cross-scenario clothing retrieval via parts alignment and auxiliary set. In *CVPR*, pages 3330–3337, 2012.
- [8] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015.
- [9] Z. Sun, C. Wang, L. Zhang, and L. Zhang. Query-adaptive shape topic mining for hand-drawn sketch recognition. In *ACM MM*, pages 519–528, 2012.
- [10] C. Wang, L. Zhao, S. Liang, L. Zhang, J. Jia, and Y. Wei. Object proposal by multi-branch hierarchical segmentation. In *CVPR*, pages 3873–3881, 2015.

Acknowledgements

The work was supported by the Key Basic Research Program of Shanghai, the National Basic Research Program of China, and the National Natural Science Foundation of China.