Hierarchical Annotation of Images with Two-Alternative-Forced-Choice Metric Learning

Hellinga Niels^{*1} Menkovski Vlado^{*1}

Abstract

Many tasks such as retrieval and recommendations can significantly benefit from structuring the data, commonly in a hierarchical way. To achieve this through annotations of high dimensional data such as images or natural text can be significantly labor intensive. We propose an approach for uncovering the hierarchical structure of data based on efficient discriminative testing rather than annotations of individual datapoints. Using two-alternative-forced-choice (2AFC) testing and deep metric learning we achieve embedding of the data in semantic space where we are able to successfully hierarchically cluster. We actively select triplets for the 2AFC test such that the modeling process is highly efficient with respect to the number of tests presented to the annotator. We empirically demonstrate the feasibility of the method by confirming the shape bias on synthetic data and extract hierarchical structure on the Fashion-MNIST dataset to a finer granularity than the original labels.

1. Introduction

High-dimensional datapoints such as natural images commonly carry complex semantic information. For example to characterize an image of a clothing item it is not enough to simply label it by its type, but we also need to know its color, gender type and size. Fine-grain annotations enable many downstream task on such data. Furthermore, it allows for efficiently organizing it in a hierarchical structure (Nguyen & Rieu, 1989). Therefore, a clothing e-commerce retailer may benefit from a certain hierarchical structure (i.e. gender > type > color > size) such that its customers (or a recommendation algorithm) can find what they are looking for quicker. This is beneficial since humans naturally group and cluster similar objects together in order to form a class or super class (Brown, 2007). As the low-level pixel information in such data is far removed from the semantic meaning and annotations we typically need complex non-linear maps build usually with deep neural networks to map to these annotations. However, train such models we also need a significant amount of annotations. To address this challenge we propose a method that leverages the efficiency of discrimination testing to capture the latent perception of difference between the data points by the annotators. Work in psychometics on measurement of subjective perception of objective stimuli provides strong insights in how such data collection can be effectively developed (Fechner, 1889). Specifically the two-alternative-forced choice (2AFC) method (Ehrenstein & Ehrenstein, 1999), which has been adapted for measurement of complex high-dimensional stimuli such as images and video (Maloney & Yang, 2003; Menkovski & Liotta, 2012). In this paper we present a method that combines 2AFC tests, with active learning methods, deep metric learning and agglomerative clustering to develop a rich embedding of the data that captures semantic relationship between the data points and uncover this semantic structure.

In order to demonstrate the feasibility of the method we empirically confirm the shape vs color bias (Ritter et al., 2017) by using our own created synthetic dataset and extract hierarchical structure on the Fashion-MNIST dataset to a finer granularity than the original labels.

2. Related Work

Extracting hierarchical structure from data is a lively field of study. In (Li et al., 2010), Li et al. present two types of hierarchies studied, namely language based (i.e. WordNet (Miller, 1995; Snow et al., 2006)) and the low level visual

^{*}Equal contribution ¹Eindhoven University of Technology Eindhoven, The Netherlands. Correspondence to: Menkovski Vlado <v.menkovski@tue.nl>, Hellinga Niels <n.hellinga@student.tue.nl>.

This work has received funding from the Electronic Component Systems for European Leadership Joint Undertaking under grant agreement No 737459 (project Productive4.0). This Joint Undertaking receives support from the European Union Horizon 2020 research and innovation program and Germany, Austria, France, Czech Republic, Netherlands, Belgium, Spain, Greece, Sweden, Italy, Ireland, Poland, Hungary, Portugal, Denmark, Finland, Luxembourg, Norway, Turkey.

²⁰¹⁹ ICML Workshop on Human in the Loop Learning (HILL 2019), Long Beach, USA. Copyright by the author(s).

feature based. Even though these approaches work fine and help in tasks like image organization, they lack the visual information that connects images together. Concepts like snowy mountains and skiing are far apart from each other on the WordNet hierarchy, which is a language based hierarchical approach but visually these concepts should be closer. There have been some purely visual feature based hierarchies (Ahuja & Todorovic, 2007; Bart et al., 2008) but they are difficult to interpret. There motivation comes from the fact that the authors belief that an image hierarchy is not following a language hierarchical structure. For example, sharks and whales should be close neighbours on in image hierarchy which is a useful property of tasks such as image classification. One problem of such visual hierarchies is that none of the work was able to evaluate the effectiveness directly. This is why Li et al. (Li et al., 2010) created a meaningful hierarchy for end-tasks such as image annotation and classification. Given the images and their tags (labels) their approach is able to automatically create a hierarchy, which is organizes images from very general to specific attributes. Ge et al. (Ge et al., 2018) propose a hierarchical triplet loss (HTL) which is able to automatically collect insightful training samples by using a predefined hierarchical structure that encodes global context information. They have two main components in their method, the constructions of the hierarchical class tree and a dynamic margin.

Fine-grained image recognition (FGIR) tasks are also closely related to extracting hierarchical structure of image data. In (Lin et al., 2015) the authors introduce a bi-linear model in order to create high-order image representations which are able to compute local pairwise interactions between features of two independent sub-networks. Such approaches have been enabled by the hierarchical representation learning present in modern convolutional neural network models (Chen et al., 2016; Kaiming et al., 2016). However, due to the high dimensionality of the features it becomes impractical for subsequent analysis. In order to reduce the high dimensionality of bilinear model features, Gao et al. (Gao et al., 2016) introduced a model that approximates such bilinear feature by using polynomial kernels. Kong et al. (Kong & Fowlkes, 2016) went a step further and introduced a classifier co-decomposition to further restrict a bilinear model.

There has also been work that is able to capture the slight visual differences between categories (Huang et al., 2016; Zhang et al., 2014) which uses bounding boxes to locate discriminative regions. The big drawback of this approach is that annotating these bounding boxes is a labour intensive process and these methods have therefore not been applicable to large-scale real world problems. In order to overcome this issue, visual attention models (Chen et al., 2018; Liu et al., 2018) where applied to FGIR tasks (Fu et al., 2017; Zheng et al., 2017) in order to automatically search the

regions of interest. It works well since it can behave as a bounding box which where labour intensive to annotate. There have also been works that use extra guidance in order to learn a semantic-related regions which, in return creates a more meaningful region for FGIR tasks. Lui et al. (Chen et al., 2016; Liu et al., 2017) introduced such work which makes use to part-based attribute in order to learn more discriminative features for fine-grained bird recognition. Also He et al. (He & Peng, 2017) used detailed text descriptions in order to mine discriminative parts or characteristics.

The most recent work is that of Chen et al., (Chen et al., 2018) who proposed a Hierarchical Semantic Embedding (HSE) framework which is able to predict categories of different levels in such a hierarchy and simultaneously integrate this structured correlation information which most of the other works, introduced above, overlook. Their HSE framework sequentially predicts category score vectors for each level and at each level of the hierarchy use the highest score vector as prior knowledge to learn a finer grained feature representation.

However, there are two main gaps in the above works which motivates our approach. One is that due to the labels of each data point there is a limitation to the depth of the hierarchy, meaning that non of the work shows finer granularity beyond the labels. The second is the resource intensive collection of labels in order to get a deeper hierarchy. Using our method, the embeddings allow us to extract a hierarchical structure which enables us to effectively circumvent the labour intensive process of labelling individual data points.

3. Method

We approach the hierarchical annotation of images by embedding the data in an embedded space that captures the semantic information that we are interested in and applying agglomerative clustering of the data in that space. To achieve such embedding, we use the 2AFC technique to measure the latent perception of differences by the annotators and use deep metric learning techniques to train an embedding model on these measurements. As 2AFC test can be inefficient in the number of queries to the annotator we optimize the test process by incorporating active learning techniques.

3.1. Two-alternative-forced-choice

Organizing information in a hierarchical structure is a natural and efficient way for the multitude of downstream tasks that we want to enable on this data (Soergel, 1985). We expect that when presented with such data, experts or annotators use a latent structure to produce the annotations. Capturing this latent structure directly is difficult because it requires a significant effort to capture and communicate it.



Figure 1. Triplet Selection Layout using two-alternative-forced choice method

On the other hand, a discirminative comparisons come with much lower cost. This characteristics has been known and utilized in psychometrics specifically in measurement of subjective perception of objective stimuli (Fechner, 1889). More recently such methods have also been developed for measurements of subjective perception of complex stimuli such as images and video (Maloney & Yang, 2003; Menkovski & Liotta, 2012; Menkovski et al., 2011). In this work 2AFC methods have been used to efficiently capture the perception of difference in between pairs of stimuli. This allowed for modeling where an individual data point reside on a relative scale of a particular quantity. In a similar manner we use the 2AFC procedure to capture the relative difference between the pair of images for a specific question.

As given in Figure 1, we select an anchor image and two query images. We ask the annotator to discriminate between the distance given by the anchor and the first query image (option A) and the anchor and the second query image (option B). The distance is with respect to a particular quantity in the image such as: the size of the object, the category of the objects, value of the object. We then store the answers by marking the image which was chosen as closer to the anchor (positive) and the other as further than from the anchor (negative).

3.2. Deep metric learning

Our aim is to embed the high-dimensional input data in to a space that captures the semantic structure that we want to uncover. As the input is high dimensional, we aim to rely on deep neural network models to capture the feature present in the image more effectively as demonstrated in by the advances of these methods in the image analysis domain (Chen et al., 2016; Kaiming et al., 2016). We also recognize that the input produced by the 2AFC test and our goals are perfectly aligned with the advances in deep metric learning and particularly with the triplet training procedure (Schroff et al., 2015).

Triplet training procedure consists of three instances of the same feed forward neural network M_e that share the same parameters. For this we used the highly successful ResNet

model(Kaiming et al., 2016). Depending on the dataset we used a different depths of the ResNets. For images of size 128x128x3 we used a ResNet-110 and for images with size 28x28x1 we had the ResNet-20. For both experiments, the models output an embedding with a dimensionality of 8.

In order to train the model we used the loss function as given in (Schroff et al., 2015). If we define the distances with respect to the anchor (x) as,

$$d(x, x^{+}) = ||M_e(x) - M_e(x^{+})||_2^2$$

$$d(x, x^{-}) = ||M_e(x) - M_e(x^{-})||_2^2$$

Then, the learning objective here is that,

$$d(x, x^+) \le d(x, x^-) - \alpha$$
$$d(x, x^+) - d(x, x^-) + \alpha \le 0$$

where α represents the margin which enforces a distance between $d(x, x^+)$ and $d(x, x^-)$. Note that alpha is also needed such that M_e cannot satisfy this equation with zero vectors for the embeddings (M_e (any image)). We used an alpha of 0.2. During training the loss function will be the following:

$$TripletLoss = Max(d(x, x^{+}) - d(x, x^{-}) + \alpha, 0) \quad (1)$$

3.3. Triplet Selection

Even though answering one of the questions is fairly quick for the annotator, the total number of available questions given a number of images is very large. Furthermore, not all questions are equally valuable for training and improving our embedding model. Such questions have been the focus of the active learning field (Settles, 2010). We used an active learning approach using the pool-based uncertainty sampling approach. Algorithm 1 shows the overall method for the active learning approach. In order to determine Q, we create pools of images where each pool contains close neighbours from a random selected image. From this pool of images we generate new potential questions. We can use the Bayes Factor as an uncertainty sampling method to determine if, for a given question q_i , whether we have a 50-50 change for choosing an answer $(a_0 \text{ vs } a_1)$ or that we have any another ratio/change such that we can be sure either a_0 or a_1 is more likely to be clicked by the annotator. Hence, we would like to compare two similar models for $a_0 \sim Bin(n, \Theta)$ given that model M_1 has a $\Theta = 0.5$ and model M_2 has an unknown Θ . For M_2 we will take the prior distribution for Θ to be uniform on [0, 1].

Using Bayes Factor we can construct the following likelihood ratio $BF = \frac{P(N'_i|M_1)}{P(N'_i|M_2)}$ where N'_i is the set of all neighbouring questions to q_i . If BF > 1 then we can strongly assume, given the data N'_i , that M_1 is supported over M_2 . Any value of BF < 1 we can assume that M_2 is supported by the data. In our case, if BF < 1 then we can assume that we know either a_0 or a_1 will be clicked by the annotator and that we do not need to ask this question again.

In order to calculate BF we need to know $P(N'_i|M_1)$ and $P(N'_i|M_2)$.

$$P(N'_i|M_1) = \binom{n}{k} \Theta^k (1-\Theta)^{n-k}$$

$$= \binom{n}{k} 0.5^k (1-0.5)^{n-k}$$

$$= \binom{n}{k} 0.5^n$$

$$P(N'_i|M_2) = \int_0^1 \binom{n}{k} \Theta^k (1-\Theta)^{n-k} d\Theta$$

$$= \binom{n}{k} \int_0^1 \Theta^k (1-\Theta)^{n-k} d\Theta$$

$$= \binom{n}{k} B(k+1, n-k+1)$$

$$= \binom{n}{k} \frac{\Gamma(k+1)\Gamma(n-k+1)}{\Gamma(k+n-k+2)}$$

$$= \frac{n!}{k!(n-k)!} \frac{k!(n-k)!}{(k+n-k+1)!}$$

$$= \frac{n!}{(n+1)!}$$

$$= \frac{1}{n+1}$$

where n is the total amount of clicks and k is equal the amount of a_0 clicks. Note that it does not matter if we count a_0 or a_1 since the test here is whether the model 'guesses' or not. If either of the two answers is favoured then M_2 will be supported by N'_i . Knowing $P(N'_i|M_1)$ and $P(N'_i|M_2)$ we can calculate BF as,

$$BF = \frac{P(N'_i|M_1)}{P(N'_i|M_2)} \\ = \frac{\binom{n}{k}0.5^n}{1/(n+1)} \\ = \binom{n}{k}0.5^n(n+1)$$

The main idea is that we want to know if, given a current question and all the previous answers, whether we probaAlgorithm 1 Triplet selection Initialize $T = \{q_1, q_2, ..., q_n\}$ - set of n random unanswered triplets. $D = \{\}$ - set of answered triplets $\tau = 0.75$ while not converged do $D \leftarrow$ Have annotators answer TUpdate M_e with D $Q \leftarrow$ select new potential questions $T = \{\}$ for q in Q do if $BF(q) > \tau$ then $T \leftarrow \text{add } q$ end if end for Sort T by highest BF $T \leftarrow \text{top } 0.8 \text{ triplets of } T + 0.2 \text{ random triplets for}$ generality end while

bility of clicking an answer will be a 50% change or not. If there is a high probability of choosing any of the two answers we do not need to ask the question. Whereas, if the probability of choosing an answer is 50% then we need to ask the question to the annotator since we cannot be sure yet.

3.4. Agglomerative clustering

After the utility of asking further questions to the annotators has diminished we conclude that we can now successfully embed the data such that its semantic information is captured by the distance metric of the space. To extract this information we run a complete-linkage agglomerative clustering algorithm (Rokach & Maimon, 2005) and produce a dendrogram that represents the captured structure.

4. Experiments and results

To evaluate the proposed method we develop two empirical studies. In the first one we test whether the method can uncover the well studied shape bias in humans (Landau et al., 1988) on a synthetic dataset. In the second we extract hierarchical structure on the FashionMNIST dataset (Xiao et al., 2017) containing images of clothing items.

4.1. Shape bias on simple shapes

We have created a synthetic simple-shapes dataset which contains 9 different shapes where each shape has 3 different thicknesses and each shape and thickness has 5 different colors. Hence, 135 unique objects that we split into a train and test set (Figure 2). The dimensionality of the images is



Figure 3. Simple Shape dendrogram splits

Figure 5. Fashion-MNIST granularity - blue

128x128x3. In this experiment the annotators give answer to the question "Which object is more similar to the anchor object?".

After collecting 840 triplets, we trained the ResNet-110 model with the specified triplet loss and extracted the data structure using the complete-linkage clustering (Defays, 1977) algorithm. Figure 3 shows the resulting splits. We can clearly see that the resulting clusters are based on the shape and not color or thickness of the objects in the images. The initial three spits are separating the different shapes: circles, triangles and rectangles. The next level the shape is again is the discriminator for the case of the circles and the rectangles, while in the case of the triangles the results are not as clear. This is somewhat expected as the case of the triangle height of the triangle is not connected to a different concept as in the case of the circle vs. oval. In the case of the rectangles are clustered against the horizontal rectangles.

4.2. Fashion-MNIST

Using the 2AFC metric learning method, we are also able to extract a hierarchical structure based on the perception of difference of the annotator. We will be using the Fashion-MNIST dataset (Xiao et al., 2017) with the question "Which object looks more similar to the anchor object?".

Results of the initial splits can be seen in Figure 4. Note that we can clearly see that the first split is based on cloth (left), bags (middle) and shoes (right) which continues further down in more fine-grained detail. Further splits of shoes can be seen in Figure 5. Here we can clearly see that we end up with clusters that present us with a finer granularity than the original Fashion-MNIST labels. We can observe for example that sandals have been split into high-heal sandals and flat sandals. In order to construct this hierarchical structure we used 1700 triplets.

We further contrast these results with clustering on the raw pixel values to form a baseline and demonstrate the value of developing the embedding space using the 2AFC tests. To compare the two sets of clusters we compute the normalized mutual information. Both results are then compared to the true labels of the Fashion-MNIST dataset. Results can be found in Table 1. Note that 'Level' is based on a binary tree level and therefore the nodes are the amount of clusters created at each level.

The results demonstrate empirically that the 2AFC method produces an embedding in which clustering captures the semantic structure in the data. We also show that our proposed method allows us create clusters with finer granularity than the dataset labels.

5. Conclusion

In this work we present a method that leverages the efficiency of discrimination 2AFC testing using to capture the latent perception of difference between data points. We have shown that we are able to capture the shape bias with synthetic data and have shown that it is possible to extract a meaningful hierarchical structure on the Fashion-MNIST

Level	Baseline	2AFC
0	0.000	0.000
1	0.192	0.392
2	0.345	0.491
3	0.426	0.583
4	0.487	0.562
5	0.499	0.520

Table 1. Normalized Mutual Information compared to true labels given

dataset, resulting in a finer granularity than the original labels. We have also achieved this efficiently by incorporating an active learning triplet selection based on Bayesian Factor estimation.

There are wide variety of applications that can benefit from extraction of hierarchical structure of data both in imaging domains such as medical imaging, but also broader in other domains that rely on high dimensional datasets.

References

- Ahuja, N. and Todorovic, S. Learning the taxonomy and models of categories present in arbitrary images. *Proceedings of the IEEE International Conference on Computer Vision*, (October), 2007.
- Bart, E., Porteous, I., Perona, P., and Welling, M. Unsupervised Learning of Visual Taxonomies. *CVPR*, 2008.
- Brown, C. *Cognitive psychology*. Wadsworth Cencage Learning, Belmont, Calif., 5th edition, 2007.
- Chen, T., Lin, L., Liu, L., Luo, X., and Li, X. DISC : Deep Image Saliency Computing via Progressive Representation Learning. *IEEE Trans. Neural Netw. Learning Syst*, 27, 6:1135–1149, 2016.
- Chen, T., Wang, Z., Li, G., and Lin, L. Recurrent Attentional Reinforcement Learning for Multi-label Image Recognition. *Proc. of AAAI Conference on Artificial Intelligence*, pp. 6730–6737, 2018.
- Defays, D. An efficient algorithm for a complete link method. *The Computer Journal*, 20(4):364–366, 1977.
- Ehrenstein, W. H. and Ehrenstein, A. Psychophysical methods. In *Modern techniques in neuroscience research*, pp. 1211–1241. Springer, 1999.
- Fechner, G. T. Elemente der Psychophysik. *Leipzig : Bre-itkopf*, 1889.
- Fu, J., Zheng, H., and Mei, T. Look Closer to See Better : Recurrent Attention Convolutional Neural Network for Fine-grained Image Recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*, 2017.

- Gao, Y., Beijbom, O., Zhang, N., Darrell, T., and Berkeley, U. C. Compact Bilinear Pooling. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 317–326, 2016.
- Ge, W., Huang, W., Dong, D., and Scott, M. R. Deep metric learning with hierarchical triplet loss. 11210 LNCS:272– 288, 2018. doi: 10.1007/978-3-030-01231-1_17.
- He, X. and Peng, Y. Fine-grained Image Classification via Combining Vision and Language. *Proceedings of* the IEEE Conference on Computer Vision and Pattern Recognitions, 2017.
- Huang, S., Xu, Z., Tao, D., and Zhang, Y. Part-Stacked CNN for Fine-Grained Visual Categorization. *Proceedings of* the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1173–1182, 2016.
- Kaiming, H., Ziangyu, Z., Shaoqing, R., and Sun, J. Deep Residual Learning for Image Recognition. *In Proceedings* of the IEEE conference on computer vision and pattern recognition, pp. 770–778, 2016.
- Kong, S. and Fowlkes, C. Low-rank Bilinear Pooling for Fine-Grained Classification. *arXiv preprint*, 2016.
- Landau, B., Smith, L. B., and Jones, S. S. The importance of shape in early lexical learning. *Cognitive development*, 3(3):299–321, 1988.
- Li, L.-j., Wang, C., Lim, Y., Blei, D. M., and Fei-fei, L. Building and Using a semantivisual Image Hierarchy. pp. 3336–3343, 2010.
- Lin, T.-y., RoyChowdhury, A., and Maji, S. Bilinear CNN Models for Fine-grained Visual Recognition. *Proceed*ings of the IEEE International Conference on Computer Vision, pp. 1449–1457, 2015.
- Liu, L., Wang, H., Li, G., Ouyang, W., and Lin, L. Crowd Counting using Deep Recurrent Spatial-Aware Network. *Proc. of International Joint Conference on Artificial Intelligence.*, 2018.
- Liu, X., Wang, J., Wen, S., Ding, E., and Lin, Y. Localizing by Describing : Attribute-Guided Attention Localization for Fine-Grained Recognition. AAAI, pp. 4190–4196, 2017.
- Maloney, L. T. and Yang, J. N. Maximum likelihood difference scaling. *Journal of Vision*, 3(8):5–5, 2003.
- Menkovski, V. and Liotta, A. Adaptive psychometric scaling for video quality assessment. *Signal Processing: Image Communication*, 27(8):788–799, 2012.

- Menkovski, V., Exarchakos, G., and Liotta, A. The value of relative quality in video delivery. *J. Mobile Multimedia*, 7(3):151–162, 2011.
- Miller, G. A. WordNet : A Lexical Database for English. COMMUNICATIONS OF THE ACM, 38(11):39– 41, 1995.
- Nguyen, G. T. and Rieu, D. Schema evolution in objectoriented database systems. *Data and Knowledge Engineering*, 4(1):43–67, 1989.
- Ritter, S., Barrett, D. G. T., Santoro, A., and Botvinick, M. M. Cognitive Psychology for Deep Neural Networks: A Shape Bias Case Study. 2017.
- Rokach, L. and Maimon, O. Clustering methods. In *Data* mining and knowledge discovery handbook, pp. 321–352. Springer, 2005.
- Schroff, F., Kalenichenko, D., and Philbin, J. FaceNet: A unified embedding for face recognition and clustering. *Proceedings of the IEEE Computer Society Conference* on Computer Vision and Pattern Recognition, pp. 815– 823, 2015.
- Settles, B. Active Learning Literature Survey. Technical report, University of WisconsinMadison, WisconsinMadison, 2010.
- Snow, R., Jurafsky, D., and Ng, A. Y. Semantic Taxonomy Induction from Heterogenous Evidence. ACL, (July): 801–808, 2006.
- Soergel, D. Organizing information: Principles of data base and retrieval systems. Elsevier, 1985.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. pp. 1–6, 2017.
- Zhang, N., Donahue, J., Girshick, R., and Darrell, T. Partbased R-CNNs for Fine-grained Category Detection. *European conference on computer vision*, pp. Springer, 834– 849, 2014.
- Zheng, H., Fu, J., Mei, T., and Luo, J. Learning Multi-Attention Convolutional Neural Network for Fine-Grained Image Recognition. *roceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5209–5217, 2017.