

Improving Image Understanding with Concept Graph

Libo Yin

The Australian National University

October 9, 2015

1 The Original HEX Model

1.1 Structure

The original HEX model [1] is an extension of the baseline flat multiclass classification model. There are two types of multiclass classification: exclusive, where the classifier predicts exactly one out of all possible states to be true; and independent, where each state is assigned true or false independently. HEX finds a balance between these two ends of the spectrum. It models the hierarchical and exclusive relationship between concepts, along with their hypernyms, with a semantic graphical model. Each node in this graphical model corresponds to a concept in the extended concept space being true or false. States of neighbouring nodes are constrained in that if a is a hypernym of b , then it is not possible that $a = 0$, $b = 1$; if a and b are exclusive, then they cannot both be true. HEX model classifies an image into a hierarchy that satisfy the above semantic consistency. A simple HEX graph is shown in [figure 1](#).

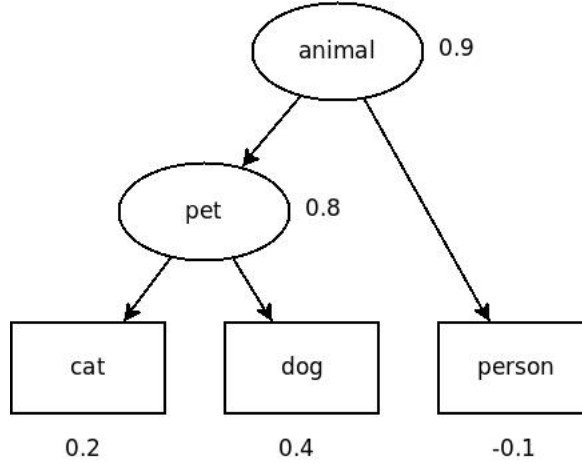


Figure 1: A simple HEX graph with three nodes in the original concept space (denoted by rectangles) and their hypernyms (denoted by ellipsoids). Directed edges denote semantical hierarchy: $(a \rightarrow b)$ if a is a hypernym of b according to WordNet; and undirected edges denote exclusion: $(a - b)$ if a and b cannot be true at the same time. Since exclusive relationship is not covered by WordNet, the exclusive subgraph is initialized greedily: two concepts are exclusive unless they share a common descendant in the hierarchical subgraph. Note that the hierarchical subgraph is in general a DAG rather than a tree. Value beside each node is the confidence on that particular node.

Denoting the set of vertices by V , the set of hierarchical edges by E_h , and the set of exclusive edges by E_e , the joint assignment $y \in \{0, 1\}^V$ can be defined as a CRF:

$$\tilde{p}(y|x) = \prod_{i \in V} \exp\{f_i(x; w)I[y_i = 1]\}I[y \text{ legal}]$$

where unary input $f_i(x; w)$ is the prediction of an arbitrary classifier¹ on concept i , and

$$I[y \text{ legal}] = \prod_{(v_i, v_j) \in E_h} I[(y_i, y_j) \neq (0, 1)] \prod_{(v_i, v_j) \in E_e} I[(y_i, y_j) \neq (1, 1)]$$

formalises local semantic constraints on hierarchical and exclusive edges. Thanks to these constraints, the state space of a HEX graph is much smaller than independent multiclass classification on the same concept space. For example, the valid states of the HEX graph in [figure 1](#) and their respective potentials are listed in [table 1](#):

state	potential
\emptyset	$\exp(0)$
{animal}	$\exp(0.9)$
{animal, pet}	$\exp(0.9 + 0.8)$
{animal, person}	$\exp(0.9 - 0.1)$
{animal, pet, cat}	$\exp(0.9 + 0.8 + 0.2)$
{animal, pet, dog}	$\exp(0.9 + 0.8 + 0.4)$

Table 1: Valid states of the HEX graph in [figure 1](#) and their respective potentials.

1.2 Observations

Note that a valid state does not have to have an active node in the original concept space. This allows an image to be classified to abstract concepts when the classifier is not confident to classify to anything more concrete. This is by no doubt a desirable feature for deployment. However, since all testing images are labelled in the original concept space, classifying to the extended concept space makes performance evaluation troublesome. In this work, accuracy is tested both in the extended concept space and the original one, by limiting legitimate states to those with an active concept in the original concept space.

It is also noticed that in [table 1](#), the competition between state {animal, pet, cat} and {animal, pet, dog} depends entirely on the confidence of node “dog” and “cat”. However, to discriminate between {animal, pet, dog} and {animal, person} requires examining the confidence along the path. From this process it is clear that, in the testing stage, confidence is passed down the hierarchy from abstract concepts to concrete ones. The other side of the same coin is that, during the training stage of the underlying classifier, branching nodes receive more training data, and therefore gain more confidence. This can be seen as confidence being passed up in the hierarchical subgraph.

Listed below are two less important observations:

1. There are no learnable variables in this CRF. In other words, all learning is performed in the underlying classifier.

¹According to [1]. It shall be explained later that the choice of underlying classifier is actually not arbitrary.

2. Mathematically, CRF requires $\forall y : \tilde{p}(y|x) > 0$. However, computationally, assigning zero to $\tilde{p}(y|x)$ can be interpreted as assigning an infinitesimal value. Therefore, the above definition is computationally equivalent to a legitimate CRF.

1.3 Problems

Consider the situation in [table 1](#) again: the underlying classifier is rather confident that this image is an animal and a pet, slightly sure that it's a cat or a dog, and highly sure that it's not a person. In such case, the intuitive decision is to classify it to $\{\text{animal, pet}\}$, but not $\{\text{animal, pet, dog}\}$. The classifier chooses $\{\text{animal, pet, dog}\}$ not because the confidence is higher, but because there are more nodes along the path. This reveals the first problem of the original HEX model: depth is not normalised, and the greedy potential function only considers active nodes, not inactive ones.

The second problem is a combined consequence of the greedy exclusion setup and the first one. It is stated in [\[1\]](#) that the underlying classifier is arbitrary. This requires the range of the output of the underlying classifier to be both positive and negative. Otherwise, for example, a probabilistic classifier is used, then it is guaranteed that a bottom-level node is activated. This means the effective state space is reduced to the size of the original concept space. In such case, inference can be performed by brute force even for ImageNet, which has 1000 labels. Also, the classifier cannot predict to an abstract concept any more. This defeats the purpose of HEX.

2 Improvements

To solve the hierarchy depth problem, we define:

$$\tilde{p}(y|x) = [y \text{ legal}] \prod_i \exp\{f_i(x; w)[y_i = 1] + (1 - f_i(x; w))[y_i = 0]\}$$

There are two issues in Deng's paper. First, the potential function does not handle different depth problem, and is greedy in labelling to bottom-layer nodes. As shall be discussed later, this has similar effects to reducing the state space, and is not a desired thing. Second, the inference system has no learnable part.

To fix the first problem, we redefine the potential function to consider not only the active nodes, but also the inactive ones. However, this fix contradicts with the realistic labelling assumption. Deng's paper assumes that a high proportion of images are actually labelled to their immediate parents, e.g. Husky are labelled as dog. As a result, the bottom layer classifiers (I did not use word "leaf layer" because the hierarchy graph is a DAG in general, and HEX is a loopy CRF. In the case of PASCAL, it happens to be a forest.) have very low confidence, due to the lack of training data. As a result, stopping at the next-to-bottom layer is almost always more preferable.

Of course, this problem can be fixed by limiting the state space to those with one active bottom-layer node. With this fix, the advantage of revised potential function is clear. However, this fix removed one of the core advantages of the HEX model: the possibility to label to an intermediate node, when the classifier is not confident enough to classify to a bottom-layer node.

On this issue, there is one more point to make. All of the ImageNet or PASCAL labels are on the bottom-layer. If the classifier is allowed to label an image to an intermediate layer, then the label space is enlarged. While this can be problematic during the validation and testing stage, it is without doubt a desirable feature during the deploy stage.

To sum up here, with little confidence on the bottom layer classifiers, the problem is to attempt to classify to the bottom layers. However, in case the classifier really cannot make a

decision, it should be allowed to stop at an intermediate layer. In addition, the classifier should consider both active and inactive nodes.

Use all images to train the CNN, and images labelled to leaf nodes to train the CRF. Computation of partition function is by brute force, thanks to the tiny state space of PASCAL. Binary weights in the learned model should not suffer from the depth problem, as weights are able to adjust themselves.

3 Dataset

Not augmented by rotation or flip since the original paper didn't

4 HEX with learning

$$\theta = \arg \min_{\theta} \left\{ -\frac{C}{N} \log \prod_{(x,y) \in D} p_{\theta}(y|x) + \frac{1}{2} \|\theta\|^2 \right\}$$

$$p_{\theta}(y|x) = \frac{1}{Z(x)} \exp \left\{ \frac{1}{|V|} \sum_{i \in V} w_i (x_i \cdot I[y_i = 1] + (1 - x_i) \cdot I[y_i = 0]) \right\} \\ \cdot \exp \left\{ \frac{1}{|E|} \sum_{(i,j) \in E} t_{ij} \cdot x_i x_j \cdot I[y_i = y_j = 1] \right\} \cdot I[y \text{ legal}]$$

$$\log p_{\theta}(y|x) = \frac{1}{|V|} \sum_{i \in V} w_i (x_i \cdot I[y_i = 1] + (1 - x_i) \cdot I[y_i = 0]) \\ + \frac{1}{|E|} \sum_{(i,j) \in E} t_{ij} \cdot x_i x_j \cdot I[y_i = y_j = 1] - \log \sum_{\hat{y}} \tilde{p}_{\theta}(\hat{y}|x)$$

$$\nabla_{\theta} \log \prod_{(x,y) \in D} p_{\theta}(y|x) = \begin{bmatrix} \nabla_w \log \prod_{(x,y) \in D} p_{\theta}(y|x) \\ \nabla_t \log \prod_{(x,y) \in D} p_{\theta}(y|x) \end{bmatrix} \\ = \begin{bmatrix} \sum_{(x,y) \in D} \nabla_w \log p_{\theta}(y|x) \\ \sum_{(x,y) \in D} \nabla_t \log p_{\theta}(y|x) \end{bmatrix}$$

$$\nabla_w \log p_{\theta}(y|x) = \underbrace{\frac{1}{|V|} \left[x_i \cdot I[y_i = 1] + (1 - x_i) \cdot I[y_i = 0] \right]_{i \in V}}_{\phi_u(x,y)} - \nabla_w \log \sum_{\hat{y}} \tilde{p}_{\theta}(\hat{y}|x) \\ = \phi_u(x, y) - \sum_{\hat{y}} p_{\theta}(\hat{y}|x) \cdot \phi_u(x, \hat{y})$$

$$\begin{aligned}
\nabla_w \log \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) &= \frac{1}{\sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x)} \nabla_w \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) \\
&= \frac{1}{Z(x)} \sum_{\hat{y}} \nabla_w \exp \left\{ \frac{1}{|V|} \sum_{i \in V} w_i (x_i \cdot I[y_i = 1] + (1 - x_i) \cdot I[y_i = 0]) \right\} \\
&\quad \cdot \exp \left\{ \frac{1}{|E|} \sum_{(i,j) \in E} t_{ij} \cdot x_i x_j \cdot I[y_i = y_j = 1] \right\} \\
&= \frac{1}{Z(x)} \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) \cdot \frac{1}{|V|} \left[x_i \cdot I[\hat{y}_i = 1] + (1 - x_i) \cdot I[\hat{y}_i = 0] \right]_{i \in V} \\
&= \sum_{\hat{y}} p_\theta(\hat{y}|x) \cdot \phi_u(x, \hat{y})
\end{aligned}$$

$$\begin{aligned}
\nabla_t \log p_\theta(y|x) &= \underbrace{\frac{1}{|E|} \left[t_{ij} \cdot x_i x_j \cdot I[y_i = y_j = 1] \right]_{(i,j) \in E}}_{\phi_t(x,y)} - \nabla_t \log \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) \\
&= \phi_t(x, y) - \sum_{\hat{y}} p_\theta(\hat{y}|x) \cdot \phi_t(x, \hat{y})
\end{aligned}$$

$$\begin{aligned}
\nabla_t \log \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) &= \frac{1}{\sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x)} \nabla_t \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) \\
&= \frac{1}{Z(x)} \sum_{\hat{y}} \nabla_t \exp \left\{ \frac{1}{|V|} \sum_{i \in V} w_i (x_i \cdot I[y_i = 1] + (1 - x_i) \cdot I[y_i = 0]) \right\} \\
&\quad \cdot \exp \left\{ \frac{1}{|E|} \sum_{(i,j) \in E} t_{ij} \cdot x_i x_j \cdot I[y_i = y_j = 1] \right\} \\
&= \frac{1}{Z(x)} \sum_{\hat{y}} \tilde{p}_\theta(\hat{y}|x) \cdot \frac{1}{|E|} \left[t_{ij} \cdot x_i x_j \cdot I[\hat{y}_i = \hat{y}_j = 1] \right]_{(i,j) \in E} \\
&= \sum_{\hat{y}} p_\theta(\hat{y}|x) \cdot \phi_t(x, \hat{y})
\end{aligned}$$

A Attributed HEX

In earlier stages of this project, one of the proposed directions is the joint analysis of concept and attributes. That is, for example, to classify an image of yellow Labrador into “animal, pet, dog, yellow, hasFur”, etc. An example HEX graph is shown as follows:

While the training data is provided by UIUC’s aPascal & aYahoo dataset. [TODO: cite dataset, discuss dataset properties] The possibility of using CNN as feature extractor has been confirmed in [TODO: cite].

We only connected attributes to bottom-level concepts. In other words, the aHEX graph can be seen as a semantic part and an attribute part. Such design is for the sake of reducing loops in the aHEX graph. While junction-tree algorithm could potentially handle such a loopy

graph during inference stage, the loops makes the graph very tricky to learn as a CRF. Similar design has been applied to pixel labelling problem such as [TODO: find references].

However, it was not the loops that caused such idea to be dropped. Under the same non-learning model, potential function is submodular, therefore such extension is trivial. Under the learnable CRF model, the aHEX graph can be learned part-wise: the semantic subgraph can be learned in the same way as discussed in [TODO: cross-ref learnable CRF], whereas for the attributes part, all bottom-level concepts can be grouped into a super-concept with 20 states (corresponding to 20 PASCAL labels), then the attribute graph becomes a tree, which is learnable[TODO: find references]. With such model, the inference is the same as an non-learning system, as the potential function for the attributed part still has a submodular structure.

B Probabilistic HEX

After the ECCV 14 paper, the original authors extended HEX to pHEX by relaxing 0/1 hard constraints to $u \in (0, 1)$, and turning the HEX graph into an Ising model. pHEX managed to beat HEX by [TODO]. However, the crucial u is chosen by cross-validation, and it still applies to all constraints in the graph. In addition, pHEX requires a long time to train, although still much shorter than training the underlying CNN. Also, it uses the same potential function as original HEX. Compared to pHEX, this work focuses on using a more flexible potential function. In terms of running time, these two works are not directly comparable since this work is not scalable.

References

- [1] Jia Deng, Nan Ding, Yangqing Jia, Andrea Frome, Kevin Murphy, Samy Bengio, Yuan Li, Hartmut Neven, and Hartwig Adam. Large-scale object classification using label relation graphs. In *Computer Vision–ECCV 2014*, pages 48–64. Springer, 2014.