# Image Segmentation Using Bayesian Network and Superpixel Analysis

**Mohammad Akbari**
Graduate School for Integrative Sciences and Engineering

**Shahab Ensafi**
School of Computing

**Jie Fu**
Graduate School for Integrative Sciences and Engineering

National University of Singapore

## Abstract

Unsupervised image segmentation algorithms rely on a probabilistic smoothing prior to enforce local homogeneity in the segmentation results. Tree-structure prior is one such prior which allows important multi-scale spatial correlations that exist in natural images to be captured. This research presents a novel probabilistic unsupervised image segmentation framework called Irregular Tree-Structured Bayesian Networks (ITSBN) which introduces the notion of irregular tree structure. Our method, however, does not update the adaptive structure at every iteration, which drastically reduces the computation required. We derive a time-efficient exact inference algorithm based on a sum-product framework using factor graph. Furthermore, a novel methodology for the evaluation of unsupervised image segmentation is proposed. We also used deep learning approach to extract a set of features for each superpixels. By integrating non-parametric density estimation technique with the traditional precision-recall framework, the proposed method is more robust to boundary inconsistency due to human subjects.

## 1 Background and related works

Unsupervised image segmentation has long been an important subject of research in computer vision and image understanding. Probabilistic modeling of images provides a useful framework to conduct inference from images. There has been much interest in this area over recent years, and this has given rise to the development of a rich and varied suite of models and techniques. Numbered amongst these are wavelet models, elastic template matching, and one which has been particularly prominent is that of Markov Random Field (MRF) approaches. We will prepare an overview of some important points of these works in this section.

### 1.1 Fixed Architecture image models

Markov Random Field (MRF) is one of the most popular methods in image modeling. This class of model has seen widespread use in many areas, and in image segmentation these models have been used by authors such as Chellappa and Chatterie [18], Chellappa and Jain [19]. In MRF models a neighbourhood relation is firstly defined between pixels or groups of pixels in the image. A random variable designating the state is associated with each element in the image, and at a particular site the state value of this variable is not only given by the class of the pixel(s) to which it relates, but is also conditional on the states of its neighbours.

43     Thus the model considers the local correlations within an image, and produces a Gibbs
44     distribution which can be explored by sampling methods [20].

45     However, the nature of MRF algorithms is such that they do not examine global effects
46     directly, but only as a conspiracy of local effects, by virtue of the fact that the state of a node
47     in the model is only influenced directly by those of its neighbors. This enables efficient
48     parallel implementation of the algorithm but leads to a failure to capture explicitly
49     less-localized effects such as a region of sky or a table in an image. MRFs are undirected
50     graphs and non-hierarchical in structure which means that unlocalized or even global
51     information cannot directly be considered. This is clearly disadvantageous in image
52     segmentation.

53     Recently tree-structured belief networks (TSBNs) have been applied to image segmentation.
54     Bouman and Shapiro [11] and Luettgen and Willsky [21] provide two such examples. These
55     TSBNs use fixed balanced tree-structures which have the image (or an encoded
56     representation of the image) instantiated at the leaf level of the network and an algorithm
57     which propagates messages of belief about the status of the network to all the other nodes in
58     the tree up to the root. The resulting equilibrium state, which occurs when all nodes have had
59     their beliefs updated, provides the image segmentation.

60

## 61     1.2      Dynamic architecture image models

62     Models whose architecture can be dynamically altered provide a powerful means of
63     addressing the image translation issue. The distinction between these and the fixed
64     architecture models is that dynamic models don't merely instantiate hidden variables in a
65     fixed structure conditionally on the supplied data, but they seek also to find relationships
66     between substructures of these variables and are able to modify these relationships on the fly
67     should subsequent data demand it.

68     Such issues in images have been considered by von der Malsburg [23] who proposed a model
69     of deformable templates called the Dynamic Link Architecture (DLA). In this model the
70     input image triggers feature detector cells which are then elastically matched to templates
71     residing in cells in the layer above, by exploring dynamic links between the two layers. The
72     templates are labeled graphs of objects expected to be found in the image, and the task is
73     essentially one of labeled graph matching. Though dynamic in architecture this model is
74     non-hierarchical with all object templates residing in a single layer.

75     Geman and Geman [20] introduce line processes as an extension of the basic MRF approach.
76     These line processes (which also form an MRF) occupy the boundaries between pixels, and
77     the connection of a number of these segments together produced the regions. Line processes
78     are dynamically combined as edge elements which describe the boundaries between regions
79     in the image. They perform this within a Bayesian framework and apply the model to an
80     image restoration problem. This is an interesting model, but it still suffers from the
81     disadvantages of MRF approaches in that inference is NP-hard.

82     A promising alternative to the above is to allow certain flexibility to the network architecture
83     such that the sub-tree elements can be constructed in a way that their boundaries correspond
84     directly to the natural boundaries in the image, producing unbalanced TSBNs. It is
85     anticipated that such models would have the flexibility required to alleviate the "blockiness"
86     experienced by balanced TSBNs and be invariant to image translation. Unbalanced tree
87     structures have already been used by other areas of image modeling such as Meer and
88     Connelly with an algorithmic approach. Modeling such structures within a Bayesian
89     framework seems very attractive. Maintaining a tree structure without cross-connections
90     would further allow the use of the attractive linear-time inference algorithms of Pearl.

91     Some work applying a hierarchical Bayesian model to images has already been undertaken
92     by Utans and Gindi [25]. Their motivation was object recognition and, like von der
93     Malsburg, sought a means of identifying particular objects in an image independent of
94     position, scale and rotation. These may appear anywhere in the image and the task was to
95     identify and match them against a previously constructed model of a single instance of the
96     object. The hierarchy consists of three layers, with a single object match neuron at the top
97     level connected to the second level object subparts whose labeling is based on the stored

object structure. The lowest level neurons identify the individual components of the object in the given image, and the algorithm then tries to automatically label and connect these neurons to parents in the layer above. As in von der Malsburg [23], prior knowledge of the object is required and generalizing the approach to scenes containing multiple objects, including unknown and occluded types is a very difficult task.

## 2    Probabilistic modeling of images with tree

The main limitation of MRF models is that the inference is not computationally efficient. They also lack a hierarchical structure and do not provide a multi-scale interpretation of the image, which is highly desirable.

We concentrate on a newly emerging and promising image model called the Tree- Structured Belief Network (TSBN). Tree-structured belief networks provide a natural way of modeling images within a probabilistic framework. By this method a balanced tree-structured belief network is constructed with a single root node and the image is presented at the leaves. Inference can then be conducted by an efficient linear-time algorithm [10]. Fixed-structure TSBNs have been used by a number of authors as models of images such as Bouman and Shapiro [11]; Irving et al.[12]. They have an attractive multi-scale architecture, but suffer from problems due to the fixed tree structure, which can lead to very blocky segmentations.

Consider for instance the four level binary tree of Figure 1.1(a). The image applied in the example is a 1-d image of a black bar on a white background, and initially all the other nodes in the tree are uninstantiated. The structure and size of the TSBN directly determines the images it can deal with. The binary tree shown handles 1-d images, but more commonly quad-trees are used which better model the real 2-d images that are of interest.



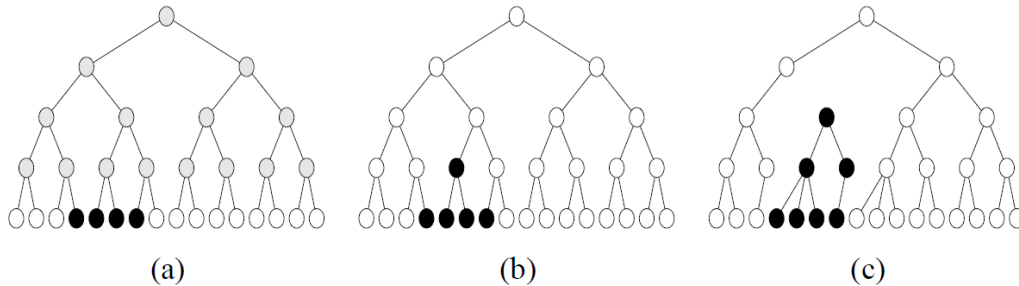(a)                      (b)                      (c)

Figure 1: (a) balanced tree with image applied (b) resulting segmentation (c) an example of dynamic tree

The hierarchical structure of TSBNs naturally leads to coarser-scale representations of the image at successive levels. As well as providing a natural mechanism for all regions in the image to have some influence over each other and thus exert global consistency, there is also potential for using the segmentations given at higher levels in image coding applications.

However some problems arise when the natural boundaries in the image do not coincide with those of sub-trees in the TSBN. This effect is illustrated by Figure 1.1(b), where the black bar spans two sub-trees with roots at the third level. The resulting segmented images exhibit an undesirable blockiness as a consequence. The aim of this work is to attempt to find models which produce good representations of natural images and to use these models to improve on current image segmentation techniques. One strategy to overcome the fixed structure of TSBNs is to break away from the tree structure, and use belief networks with cross connections. However, this means losing the linear-time belief-propagation algorithms that can be used in trees (Pearl, 1988a) and using approximate inference algorithms.

TSBNs have many attractive properties and we believe that models based upon them, but having a dynamically adjustable tree structure to enable their boundaries to better reflect those of the image, provide a very promising starting point. Such models we have named dynamic trees (DTs) and one such tree produced for the toy image data is shown in the Figure 1.1(c).

141 Dynamic trees are a generalization of the fixed-architecture tree structured belief networks.
142 Belief networks can be used in image analysis by grouping its nodes into visible units $\mathbf{X}v$, to
143 which the data is applied, and having hidden units $\mathbf{X}h$, connected by some suitable
144 configuration, to conduct a consistent inference. DTs set a prior $P(\mathbf{Z})$. Over tree structures $\mathbf{Z}$
145 which allows each node to choose its parent so as to find one which gives a higher posterior
146 probability $P(\mathbf{Z},\mathbf{X}h/\mathbf{X}v)$ for a given image $\mathbf{X}v$. This effectively produces forests of TSBNs.

147 There are two essential components that make up tree architectures and the nodes and
148 conditional probability tables (CPTs) in the given tree. There are a very large number of
149 possible structures; in fact for a set of nodes created from a balanced tree with branching
150 factor b and depth D (with the top level indexed by 1) there are $\prod_{d=2}^{D}(b^{(d-2)} +$
151 $1)^{b^{d-1}}$ possible forest structures.
152 Our objective will be to obtain the maximum a posteriori (MAP) state from the posterior
153 $P(Z|Xv) \propto P(Z)P(Xv|Z)$. For any Z it is possible to compute $P(Xv|Z)$ (using Pearl message
154 passing) and P(Z). However, since the number of configurations of Z is typically very large
155 it will usually be intractable to enumerate over them all and other approaches need to be
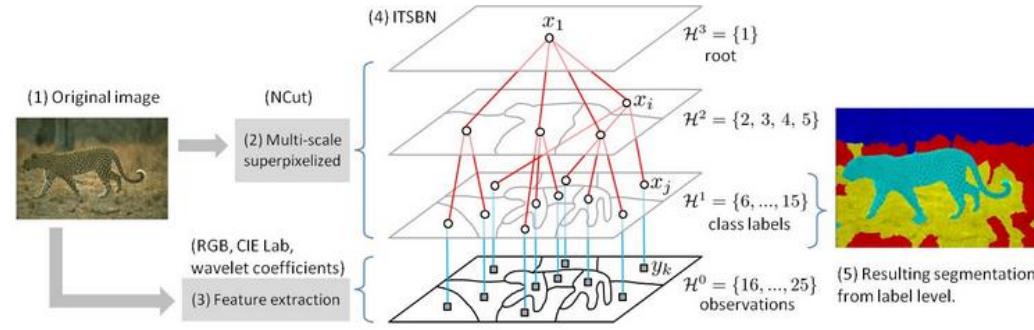156 adopted.
157



159 Figure 2. The overview of Irregular Tree-Structured Bayesian Network (ITSBN) framework

160 Sampling by Markov Chain Monte Carlo (MCMC) techniques, or search using Simulated
161 Annealing are possibilities which have been considered in [28], but the drawback is that they
162 are slow. An alternative to sampling from the posterior $P(Z, Xh|Xv)$ is to use approximate
163 inference. One possibility is to use a mean-field-type approximation to the posterior of the
164 form $Q_z(Z)Q_{X_h}(X_h)$ and this is which have been considered in [14].
165 One another alternative is to find the best structural connectivity for Z and then using it to
166 find the best possible labeling configuration for the tree-structure of the image segments. In
167 this way, we will have an approximation of the best structure connectivity between nodes of
168 the tree. Suppose that the structure connectivity is characterized by an $N \times N$ matrix
169 adjacency matrix Z, where $z_{ij}$ takes a values 1 if node j is the parent of node i. there is only
170 one constraint which each node could be connected to a parent node in the directly adjacent
171 layer. Therefore, a realization of structure matrix Z can have at most one entry equal to 1 in
172 each row.
173 Suppose we have an image which is segmented by our method, see figure 2. Each node is a
174 discrete random variable $x_i$ of an image section i. This random variable takes one of C
175 possible labels/classes. The probability which defines nodes j and its parents i has been
176 labeled v and u respectively is $\theta_{ij}^{uv} = p(x_j = v | x_i = u)$
177 We can rewrite this equation by defining two indicator variables $x_j, x_i$ which pick out the
178 correct probabilities.

$$p(x_j = v | x_i = u) = \prod_{v=1}^{C} \prod_{u=1}^{C} [\theta_{ij}^{uv}]^{x_j^v x_i^u}$$

180 Note that this framework is unsupervised, hence the probabilistic model of a class c is not
181 provided a priori, and the number of labels allowed for each input image must be defined
182 before executing the algorithm. Estimating the appropriate number of classes for each image

183   (the model selection problem) is explained in more detail in Section 3.

184   We introduce observed variables represented by shaded square-shaped nodes in the structure
185   as illustrated in Fig. 1. Each observed random variable $y_e \in R_d$ of an image site $e \in V$
186   represents the relevant image features such as color or texture which take on continuous
187   values. Extensive details on our choice of features will be discussed in the experimental
188   results section. We model the feature vector $y_e$ using a multivariate Gaussian distribution
189   given as:

190
$$p(y_e|x_i) = \prod_{c=1}^{C} \mathcal{N}(y_e|\mu_c, \Lambda_c^{-1})^{x_{ic}}$$

191   Where $\mathcal{N}(y_e|\mu_c, \Lambda_c^{-1})^{x_{ic}}$ is the multivariate Gaussian distribution and $\mu_c, \Lambda_c^{-1}$ are mean and
192   covariance matrices for class c respectively. Using the notation described above, we can now
193   write the hidden labels collectively as $X = \{x_i\}_{i \in H}$ which H denotes the first hidden layers
194   and the observed image features as $Y = \{y_e\}_{e \in V}$. The log-likelihood of the complete data can
195   be expressed as:

196
$$\log p(X,Y|Z,\theta) = \sum_{e \in V} \sum_{i \in H^1} z_{ei} \prod_c x_{ic} \log \mathcal{N}(y_e|\mu_c, \Lambda_c^{-1})^{x_{ic}}$$

197
$$+ \sum_l (\sum_{j \in H^1} \sum_{i \in H} z_{ji} \sum_{v=1}^{c} \sum_{u=1}^{c} x_j^v x_i^u \log \theta_{ij}^{uv})$$

198   Note that the connectivity structure matrix Z is assumed known from each input image,
199   hence is excluded from the parameter set. In the next section, we derive an efficient inference
200   and parameter estimation algorithm for ITSBN.

201

## 3       Inference and parameter Learning

203   In this section, we present a maximum likelihood estimation algorithm for ITSBN model
204   parameter $\theta$. For models with hidden variables, Expectation-Maximization (EM) [13] is a
205   principled framework that alternates between inferring the posterior over hidden variables in
206   the E-step, and maximizing the expected value of the complete data likelihood with respect
207   to the model parameter in the M-step. More specifically, we infer the posterior probability
208   over the hidden labels in the E-step and compute relevant sufficient statistics needed in the
209   M-step for maximizing $\theta$. Since exact inference on a tree structure graph is tractable, the
210   objective function can be written as:

211
$$\mathcal{F}(\theta, \theta^{t-1}) \equiv \langle \log p(X,Y|Z,\theta) \rangle_{p(X|Y,Z,\theta^{t-1})}$$

212   Where $\langle f(x) \rangle_{q(x)}$ denotes the expectation of function f(x) respect to the distribution q(x). In
213   M-step, by maximizing $\mathcal{F}$ with respect to $\theta$, we derive close-form update equations for
214   $\mu_c, \Lambda_c^{-1}$ and $\theta_l^{uv}$ as follows:

215
$$\mu_c = \frac{\sum_{e \in V} \sum_{i \in H^1} z_{ei} \langle x_i^c \rangle_{p(X|Y,Z,\theta^{t-1})} y_e}{\sum_{e \in V} \sum_{i \in H^1} z_{ei} \langle x_i^c \rangle_{p(X|Y,Z,\theta^{t-1})}}$$

216
$$\Lambda_c^{-1} = \frac{\sum_{e \in V} \sum_{i \in H^1} z_{ei} \langle x_i^c \rangle_{p(X|Y,Z,\theta^{t-1})} (y_e - \mu_c)(y_e - \mu_c)^T}{\sum_{e \in V} \sum_{i \in H^1} z_{ei} \langle x_i^c \rangle_{p(X|Y,Z,\theta^{t-1})}}$$

217

218
$$\theta_l^{uv} = \frac{\hat{\theta}_l^{uv}}{\sum_{\acute{u}} \hat{\theta}_l^{\acute{u}v}}$$

219   Where $\hat{\theta}_l^{uv} = \sum_{j \in H^l} \sum_{i \in H} z_{ji} \langle x_i^u x_j^v \rangle_{p(X|Y,Z,\theta^{t-1})}$ denotes the unnormalized class-transition CPT.
220   To compactly explain our inference step, we express ITSBN in terms of a factor graph [4]
221   using 2 types of nodes: 1) a variable node for each random variable $x_i$ and 2) a factor node
222   for each local function, namely CPT in this case. In the operation of sum-product algorithm,
223   we compute a message from variable node x to factor node $f: m_{x \to f}(x)$, and message

224   from factor node f to variable node $f: m_{f \to x}(x)$, which can be expressed as:

225
$$m_{x \to f}(x) = \prod_{h \in n(x) \backslash \{f\}} m_{h \to f}(x)$$

226

227
$$m_{f \to x}(x) = \sum_{n(f) \backslash x} \left( f(n(x)) \prod_{y \in n(x) \backslash \{x\}} m_{h \to f}(x) \right)$$

228
229  where $n(x)$ denotes the set of neighbors of a given node x, and $n(x)\{f\}$ denotes the
230  remaining after f is removed from the set $n(x)$.

231

## 4    Building the tree

233  In this section, we discuss the methodology applied to each input image in order to create an
234  irregular tree structure, where each level in the hierarchy can be regarded as a scale of
235  visualization (coarse-to-fine representation). The details of over-segmentation algorithm and
236  connecting superpixels across the adjacent level are discussed as followed.

237

### 4.1    Superpixel segmentation

239  There are many approaches to generate superpixels, each with its own advantages and
240  drawbacks that may be better suited to a particular application. For example, if adherence to
241  image boundaries is of paramount importance, the graph-based method of [4] may be an ideal
242  choice. However, if superpixels are to be used to build a graph, a method that produces a
243  more regular lattice, such as [1], is probably a better choice.

244  In this paper, we use the implementation in [15] due to its computational efficiency and
245  homogeneity in the resulting superpixel output e.g. pixels with similar texture, color, and
246  residing within the same object boundary are grouped into the same superpixel.

247  In our experiment, the number of hierarchical level L is fixed to 6 for every image, and the
248  number of superpixels from level $l = L - 2$ to leaf $l = 0$ are approximately 5, 20, 50, 200
249  and 200 respectively. As mentioned earlier the root level contains only 1 node while the leaf
250  level has the same number of nodes as level $l = 1$. The same setting is applied to every
251  image in the dataset.

252

### 4.2    Structure of tree

254  We describe a method to build an irregular tree structure from each input image and its
255  corresponding multi-scale superpixels derived from the output of [15]. Two constraints are
256  required in order to build such a tree in our experiment. First, each observed node $y_e$ in the leaf
257  level can connect to its parent in level $l = 1$ in a 1-to-1 transparent manner as shown in Fig. 1
258  since both levels are indeed derived from the same over-segmenting process. Second, each child
259  node in level $l \in \{1 \ldots L - 2\}$ can only be connected to a parent in level $l + 1$. This regularization
260  enables smoothness of the information transfer across the scale of image. The constraint is
261  formulated formally as follows. For all $j \in H^l$, $i \in H^{l+1}$ and $l \in \{1 \ldots L - 2\}$ the connectivity
262  $z_{ji} = 1$ when $i = i^*$ and otherwise $z_{ji} = 0$; when $i^* = argmax |S_i \cap S_j|$.
263  Where $|S_i \cap S_j|$ means the number of overlapped pixels when intersecting superpixel Si and Sj
264  together. When regarding Sj as an arbitrary set, this idea can be extended to build a datadriven tree
265  structure in other settings beyond superpixels, for instance block pixels, voxels, abstract image
266  regions. Furthermore, the proposed method still works in the case where Sj is not a strict subset of
267  $S_{i^*}$ which makes the construction process quite robust and general.

268

## 5    Extracting features

270  In the bottom of the tree, we have a set of visible nodes which corresponds to image
271  segments. We will use feature learning to extract image features from segments. In this
272  section we discuss about feature learning and its advantage to our work.

273  Images are highly variable (because of viewpoint changes, shape variation, etc.) and

274 high-dimensional (often hundreds of thousands of pixels). It is difficult for computer vision
275 algorithms to run on raw image data and to generalize well from training data to unseen data.
276 Therefore, it is desirable to extract features to make computer vision tasks easier to solve and
277 finding "good" feature representations is vital for this kind of high-level computer vision
278 task. Features are defined here as attributes that can be extracted from the input data, and
279 these feature representations could then facilitate subsequent high-level computer vision
280 tasks.

281 Both color and texture features are extracted from an input image. Color feature is obtained
282 by averaging RGB and CIE Lab color space within a superpixel. Texture feature is acquired
283 from applying 3-level complex wavelet transform (CWT) to the input image, then averaging
284 the magnitude of wavelet coefficients within a superpixel. Totally, we have 15-dimensional
285 feature vector from each superpixel; 6 from color and 9 from texture.

286 State-of-the-art feature extractors usually consist of a filter bank, a non-linear transformation,
287 and some kind of feature pooling method [15], which collect nearby values through an
288 average, maximum, or histogramming computation. The filter bank could be oriented edge
289 detection filters such as Scale Invariant Feature Transform (SIFT) [16]. Hand-crafted feature
290 methods, such as SIFT, are still the most common ones. Typically, these hand-crafted
291 features use low-level concepts (e.g. edges) rather than higher-level concepts such as parallel
292 lines or basic geometric elements [17].

293 It has been shown that dictionary learning can be used in computer vision to learn features
294 for modeling the local appearance of natural images, leading to state-of-the-art performance
295 [17]. Thus, learning hierarchical feature representations automatically is a more flexible and
296 more general strategy. In this project, we adapt Deconvolutional Networks [17] to form the
297 hierarchical sparse convolutional feature representations in a purely unsupervised fashion.
298 However, the drawback of Deconvolutional Networks architecture is its fix geometry. As
299 the features are aggregated according to a predefined pattern, the higher
300 level features represent data with poor spatial accuracy [29]. Therefore, we only use a 2-layer
301 architecture in this segmentation project.

302

303 ## 6      Evaluation and experiments

304 While evaluation methodology for supervised image segmentation result is obvious, there is
305 no clear standard to evaluate unsupervised image segmentation. We adopt the methodology
306 mentioned in [26] which regards the border between each pair of adjacent resulting
307 segmentation regions as the object boundary which can be evaluated using boundary-based
308 precision-recall (PR) framework of [27]. In the framework, each pixel in the resulting
309 boundary image can take a real value ranging from 0 to 1; 1 when the pixel is believed to be
310 a boundary and 0 when is not. However, the framework is still vulnerable to boundary
311 misalignment occurred naturally between the boundaries produced by segmentation
312 algorithms and ones produced by several human subjects.

313 We report our results from experiment on segmentation and matching of boundary images
314 from Berkeley Segmentation Data Set and Benchmarks 500 (BSDS500), reported in [26].
315 The dataset contains 300 training and 200 testing images, each of which has ground-truth
316 boundary images done manually by five different human subjects on average. Since this is
317 unsupervised framework, we only run our algorithm on the 200 testing images directly
318 without training the model beforehand, and thus the 300 training images are discarded.

319 We used two kinds of feature vectors in this research. In the original paper, authors have used
320 a simple feature vectors which extracted from color and texture information. We also used
321 feature learning to use a feature vector for each pixel and then each superpixels. It shows
322 promising result and has better performance in terms of precision and accuracy.

323 In the first experiments, both color and texture features are extracted from an input image.
324 Color feature is obtained by averaging RGB and CIE Lab color space within a superpixel.
325 Texture feature is acquired from applying 3-level complex wavelet transform (CWT) to the
326 input image, then averaging the magnitude of wavelet coefficients within a superpixel.
327 Totally, we have 15-dimensional feature vector from each superpixel; 6 from color and 9
328 from texture. Note that we have not used significant number of features in our experiment

compare to others in the literature. That is because our focus is on the performance of the ITSBN alone, not on the image feature.

We extract a 45-dimensional learned feature vectors using 2-layer Deconvolutional Networks [17]. Features have been extracted from a 9*9 kernel around each pixel. As we use supetpixels in the leaves of the tree, we construct a feature vector for each superpixels by averaging the feature vectors of pixels belongs to. So for each superpixels we have a 45-dimension feature vector.

At the end, the values of precision, recall, F-measure and computational run-time are averaged across all the test images in the dataset as shown in Table 1. The values summarizes the performance of both methods applied to the dataset. At the same parameter setting, the results show that ITSBN slightly outperforms supGMM. The greater precision value indicates less noisy image segmentation results in ITSBN than supGMM. This makes a lot of sense because ITSBN has tree-structured prior which is equivalent to adding the label smoothness regularization to the maximization of the log-likelihood. However, supGMM seems to slightly does a better job on recalling the boundary pixel.

Table 1. The across-dataset average of precision, recall and F-measure of supGMM, ITSBN and FL

| Method | Precision | Recall | F-measure | Run-time |
|---|---|---|---|---|
| supGMM | 0.2597 | 0.5711 | 0.3470 | 30sec/image |
| ITSBN | 0.2636 | 0.5697 | 0.3501 | 2min/image |
| FL+ITSBN | 0.2154 | 0.4127 | 0.2831 | 2min/image+ |

## 7    Conclusion

In this project, we have described and implemented an efficient unsupervised image segmentation algorithm using probabilistic tree-structure Bayesian network. We implemented Irregular Tree Structure Bayesian Networks (ITSBNs). The tree structure of ITSBN can be tailored to fit natural boundaries in each input image, thus eliminates blocky segmentation results. Based on the factor graph representation, we derived a sum-product algorithm for the inference step. In order to avoid adapting the tree structure at each iteration, we presented a fast and simple method to build an irregular tree from a set of superpixels in different scales extracted from each input image. This flexibility indeed is very much desired as a great number of such algorithms have been made publicly available, each with different properties. We also used feature learning method to extract features from each superpixels. By integrating non-parametric density estimation technique with the traditional precision-recall framework, the proposed method is more robust to boundary inconsistency due to human subjects. We experimentally show the improvement of ITSBN over the baseline method which motivates us to further investigate the model of similar type. Although, feature learning method did not extract effective features for segmentation, we can improve it by combining these learned feature and simple features.

**References**

[1] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.

[2] T. Cour, F. Benezit, and J. Shi. Spectral segmentation with multiscale graph decomposition. In *CVPR*, 2005.

[3] A. Levinshtein, A. Stere, K. Kutulakos, D. Fleet, S. Dickinson, and K. Siddiqi. Turbopixels: Fast superpixels using geometric flows. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2009.

[4] Pedro Felzenszwalb and Daniel Huttenlocher. Efficient graph-based image segmentation. In *CVPR*, 2004.

[5] Alastair Moore, Simon Prince, Jonathan Warrell, Umar Mohammed, and Graham Jones. Superpixel Lattices. In *CVPR*, 2008

[6]   O. Veksler, Y. Boykov, and P. Mehrani. Superpixels and supervoxels in an energy optimization framework. In *ECCV*, 2010

[7]   D. Comaniciu and P. Meer. Mean shift: a robust approach toward feature space analysis. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, May 2002.

[8]   A. Vedaldi and S. Soatto. Quick shift and kernel methods for mode seeking. In *ECCV*, 2008.

[9]   Luc Vincent and Pierre Soille. Watersheds in digital spaces: An efficient algorithm based on immersion simulations. In *IEEE Transactions on Pattern Analalysis and Machine Intelligence*, 13(6):583–598, 1991.

[10]  J. Pearl. *Probabilistic reasoning in Intelligent Systems*. Morgan Kaufman Publishers, 1998

[11]  C. A. Bouman and M. Shapiro. A Multiscale Random Field Model for Bayesian Image Segmentation. In *IEEE Transaction on Image Processing*, 3(2), 1994.

[12]  W. W. Irving, P. W. Fieguth and A. S. Willsky, An Orverlapping Tree Approach On Multiscale Stochastic Modeling and Estimation. In *IEEE Transaction on Image Processing*, 6(11), 1997.

[13]  G. Mori, Guiding Model Search Using Segmentation, In *ICCV*, 2005.

[14]  N. Adams, C.K.I. Williams, Dynamic Trees for Image modeling, In *Image and Vision Computing*, 2002.

[15]  Kavukcuoglu, K, Ranzato, M.A, LeCun, Y. What is the best multi-stage architecture for object recognition?. In *CVPR*, 2009.

[16]  David G Lowe, Distinctive image features from scale-invariant keypoints, In *Internatonal Journal of Computer Visopm*, 60(1), 2004

[17]  M. Zeiler, D. Krishnan, G. Taylor, R. Fergus. Deconvolutional Networks. In *CVPR,* 2010.

[18]  Chellappa, R. and Chatterie, S. Classification of Textures using Guassian Markov Random Fields. In *IEEE Trans. Accoust., Speech and Signal Processing*, volume 33, pages 959–963, 1985.

[19]  Chellappa, R. and Jain, A. *Markov Random Fields - Theory and Applications*. Academic Press Ltd, London, UK, 1993.

[20]  Geman, S. and Geman, D.. Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, volume 6, no. 6, pages 721–741, 1984.

[21]  Luettgen, M. R. and Willsky, A. S. Likelihood Calculation for a Class of Multiscale Stochastic Models, with Application to Texture Discrimination. In *IEEE Transactions on Image Processing*, 4(2), 194–207, 1995.

[22]  Dayan, P., Hinton, G. E., Neal, R. M., and Zemel, R. S. The Helmholtz Machine. In *Neural Computation*, 7(5), 1995.

[23]  von der Malsburg, C. Pattern Recognition by Labelled Graph Matching. In *Neural Networks*, volume 1, pages 141–148, 1988.

[24]  Montanvert, A., Meer, P., and Rosenfeld, A. Hierarchical Image Analysis Using Irregular Tessellations. In *IEEE Trans. Pattern Analysis and Machine Intelligence*, 13(4), 307–316, 1991.

[25]  Utans, J. and Gindi, G. Improving Convergence in Hierarchical Matching Networks for Object Recognition. In S. J. Hanson, J. D. Cowan, and C. L. Giles, editors, In *NIPS*, 1993.

[26]  Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik, Contour detection and hierarchical image segmentation. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, pp. 898–916, 2011.

[27]  D.R. Martin, C.C. Fowlkes, and J. Malik, Learning to detect natural image boundaries using local brightness, color, and texture cues. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 5, pp. 530 –549, 2004.

[28]  R. Jenssen, D. Erdogmus, K. Hild, J. Principe, and T. Eltoft, Optimizing the Cauchy-Schwarz PDF distance for information theoretic, non-parametric clustering. In *Energy Minimization Methods in Computer Vision and Pattern Recognition*. Springer, pp. 34–45, 2005.

[29]  Leon Bottou. From Machine learning to Machine Reasoning, *arXiv preprint*, 2011