

Saliency Map for Visual Attention Region Prediction: A Comparison of Two Methods

Mao Wang

School of Mechanical and Electrical Engineering, Beijing Information Science & Technology

University, Beijing, 100192, China

wangmao775@163.com

Keywords: Human-Computer Interaction, Saliency Map, Fuzzy Inference, Fuzzy Neural Network, Intention Recognition.

Abstract. Visual attention region prediction has been paid much attention by researchers in intelligent systems recent years because it can make the interaction between human and intelligent agents to be more convenient. Saliency determines the capability of an image detail to attract visual attention and thus guide eye movements in a bottom-up way. A lot of models for saliency map combining color, intensity and orientation feature maps by simple normalization and linear summation, which can not reflect the importance of each feature in saliency map well. Therefore, in this paper, the prediction method of the visual attention region inferred by using fuzzy inference and fuzzy neural network (FNN) after extracting and computing of images feature maps and saliency maps were proposed and compared. A method for training FNN is also proposed. A user experiment was conducted to evaluate and compare the prediction effect of proposed methods by making surveys for the prediction results. Furthermore, t test results shown that there are significant difference between the results got by two different methods. This also indicated that prediction method based on FNN proposed by us has a better performance in the level of attention regions' position prediction according to different images.

Introduction

Recently, the intention recognition, recognizing the intention of a user or an agent by analyzing their actions or changes of state, is becoming an important issue in various research fields of intelligent systems. Much of early work is in the context of speech understanding and response automatically[1]. For example, Pynadath et al. achieved the plan recognition on a problem in traffic monitoring through exploited the context by using a general Bayesian framework[2]. More recently, Pereira et al.[3] described an approach to tackle intention recognition by combining dynamically configurable and situation-sensitive Causal Bayes Networks plus plan generation techniques[4,5]. Mao et al. have presented a utility-based approach to solve the recognition of intention, which is realized by incrementally using plan knowledge and observations to change state probabilities [6].

In their researches, the probability is the main factor which used to infer the human intention. Another method is use an automatic capture of bottom-up salient stimuli and volitional shifts guided by and top-down context factors[7,8], where bottom-up salient stimuli are the external factors to user and top-down contexts are the internal factors. The characteristic of all the methods mentioned above is that one or several characteristic values of the image showing before user have been extracted, calculated and combined as their basis for inferring user's intention. But as mentioned in [7,8], the saliency map which based on color, intensity and orientation feature maps of images, is obtained by simple normalization and linear summation of the three features.

In this paper, two methods for visual attention region prediction system inspired on saliency map are described. The aim of this work is to present a new approach that improves the performance of attention prediction based on saliency map by comparing the effects of fuzzy inference and fuzzy neural network (FNN) methods. In this paper, both the fuzzy inference and FNN employing features of image as input allows us to

combine features and infer with great flexibility some intuitive decision rules based on the visual perception principles.

Overview of Proposed Approach

We propose two new approaches to predict visual attention region based on image's saliency map. The overall procedural flow of proposed approach is summarized in Fig. \ref{fig_fmarchitecture}.

Firstly, intensity, color (red, green, blue and yellow) and orientation (degrees of $0, \pi/4, \pi/2, 3\pi/4$) are extracted in multi-resolution from the Gaussian pyramid by linear filtering. Then a saliency map is generated for each of them by computing the difference between the layers of the pyramids, which imitates the center-surround type receptive field. Finally, the saliency map is generated by fuzzy inference or a trained FNN using the three saliency maps as inputs. The training method of FNN will be explaining in the following section.

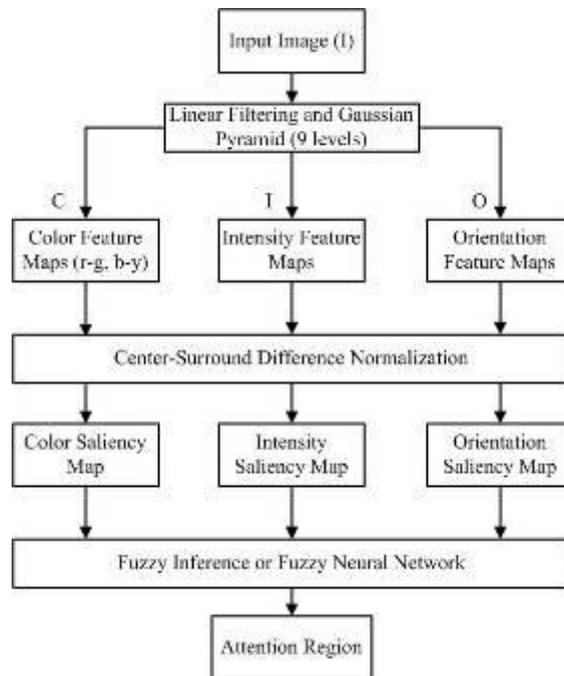


Fig. 1 Overall Procedure for Proposed System

Feature Maps

We propose two new approaches to predict visual attention region based on image's saliency map. The overall procedural flow of proposed approach is summarized in Fig. 1. In Fig. 1, the linear filter is used in order to compute center-surround different of various features at 9 scales. In this paper, the input image is sub-sampled into a dyadic Gaussian pyramid by convolution with a linearly separable Gaussian filter and decimation by a factor of two[10], which means that for example, the third level has a resolution of 1/8 of the input image's.

After filtering, the three features of images have their values at each position according to the input image, which are divided into 9 levels of pyramid ready to be calculated. Then, in this paper, the color feature is reflected by two values defined by us, which are red-green and blue-yellow opponencies. If r, g, b are the red, green, and blue values of the input color image respectively, then the color map of one level can be calculated according to the following equations.

$$M_{r-g} = \frac{r - g}{\max(r, g, b)} \tag{1}$$

$$M_{b-y} = \frac{b - \min(r, g)}{\max(r, g, b)} \quad (2)$$

where M_{r-g} , M_{b-y} stand for red-green and blue-yellow opponencies. Red-green and blue-yellow opponencies are central to modeling the contribution of color to saliency because of these two opponency axes can cover the entire visible light[11]. And note that the definitions deviate from the original model by[12].

The intensity map M_i of one level is calculated as:

$$M_i = \frac{r + b + g}{3} \quad (3)$$

These operations are repeated for each level of the input to obtain an intensity pyramid with also 9 levels. Local orientation map M_o is obtained by applying steerable filters to the intensity pyramid levels M_i [13]. After getting M_{r-g} , M_{b-y} , M_i and M_o , in order to yield the feature maps, we simulate the center-surround receptive fields by subtraction between two maps at the center (c) and the surround (s) levels in these pyramids. They can be calculated as

$$F_{l,c,s} = N(|M_l(c) - M_l(s)|) \quad (4)$$

$$l \in L = L_c \cup L_l \cup L_o$$

where

$$L_c = \{I\},$$

$$L_l = \{r - g, b - y\},$$

$$L_o = \{0, \frac{\pi}{4}, \frac{\pi}{2}, \frac{3\pi}{4}\}.$$

Note that N is an iterative operator for local nonlinear iterative competition between salient locations within each feature map. And F_l is the feature map summed over the center-surround combinations using across-scale addition while F_c and F_s stand for feature maps at the center and the surround levels in these pyramids, respectively.

Finally, by summing over the center-surround combinations and normalizing again according the results obtained in Eq. (4), the feature maps of color, intensity and orientation can be obtained according to Eq.(5) as C_c , C_i , C_o , respectively. Here, all the feature maps we need for building region saliency map are already obtained. Fig. 2 shows an example of three feature maps mentioned above.

$$C_c = F_l,$$

$$C_i = N\left(\sum_{l \in L_c} F_c\right), \quad (5)$$

$$C_o = N\left(\sum_{l \in L_o} F_o\right).$$

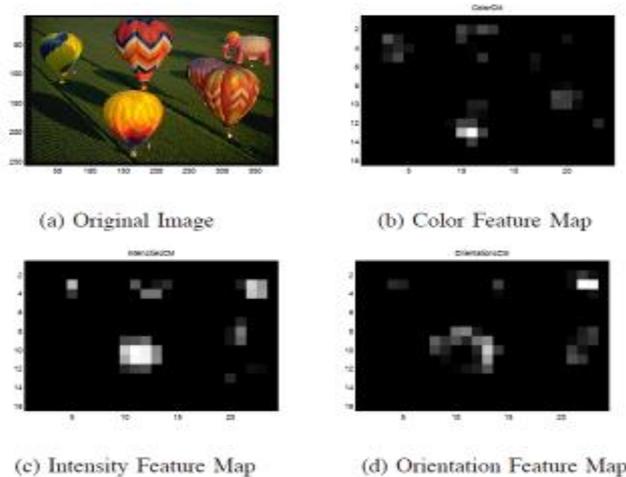


Fig. 2 Example of Feature Maps Got from Image

Saliency Map by Fuzzy Inference

Most attention models are based on a saliency map and a dynamical process for visiting saliency maxima. Itti et al.[7] introduced a model for the bottom-up selective attention based on serially scanning a saliency map, which is computed from local feature contrasts, for salient locations in the order of decreasing saliency[8].

But this method also has its weakness which is the method of simply sum of them gives the same importance at the same time. But based on experiments[9], most people do not pay equal attention to all of them. In order to solve this problem, we proposed a method to compute saliency map by fuzzy inference based on the features of graphics. In this way, the importance of all features can be reflected in fuzzy rule respectively.

In this study, we use the feature variables from color feature map(C_c), intensity feature map(C_i) and orientation feature map (C_o) in the *IF* part while the output value in *THEN* part is value of region saliency map(S_m).

Every value of region saliency map is decided by fuzzy inference rules as shown in Table 1, where C, I, O stand for color, intensity, and orientation respectively. Fig. 3 shows the membership functions and singletons.

Table 1 Fuzzy Inference Rules

C	O			
	I	OL	OM	OH
CL	IL	SVL	SL	SLL
	IM	SL	SLL	SLL
	IH	SLL	SLL	SM
CM	IL	SL	SLL	SM
	IM	SLL	SM	SM
	IH	SM	SM	SLH
CH	IL	SM	SLH	SH
	IM	SLH	SH	SH
	IH	SH	SH	SVH

In Table 1, CL, CM and CH are fuzzy labels that represent each state of value from color feature map is low, medium and high, respectively. And so on to IL, IM, IH, OL, OM and OH.

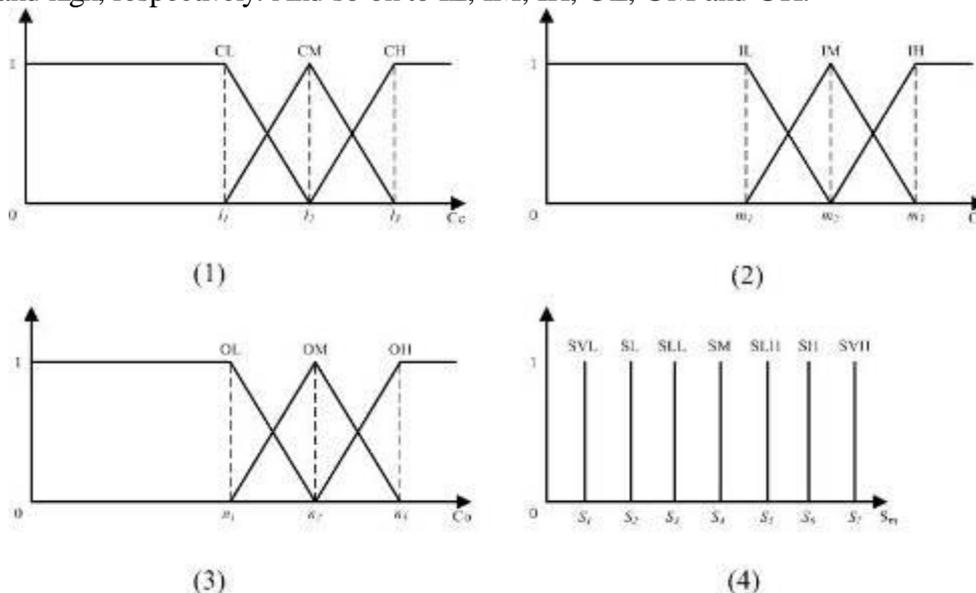


Fig. 3 Membership Functions and Singletons((1)~(3) are Membership Functions in IF Part, (4) is Singletons in THEN Part)

Saliency Map by Fuzzy Neural Network

It is worth to note that the fuzzy inference method for saliency map is only suitable for specific images, which means is not universal. therefore, we proposed a method by using FNN to solve this problem. In this way, the importance of all features can be reflected in fuzzy rule with the human decision making model by the conceptual framework of fuzzy logic.

A FNN system is a learning machine that finds the parameters of fuzzy rules by exploiting approximation techniques from neural networks[14].

We have pointed out the defect of ordinary combination method of feature maps in the beginning of this section. It has not an importance distinction between various features, especially when a feature is more important comparing with others.

In this study, we use the feature variables from color feature map (C_c), intensity feature map (C_i) and orientation feature map (C_o) as input while the output is a value of region saliency map (S_m).

Every value of region saliency map is decided by FNN as shown in Fig. 5 where G stands for Gaussian Function[15].

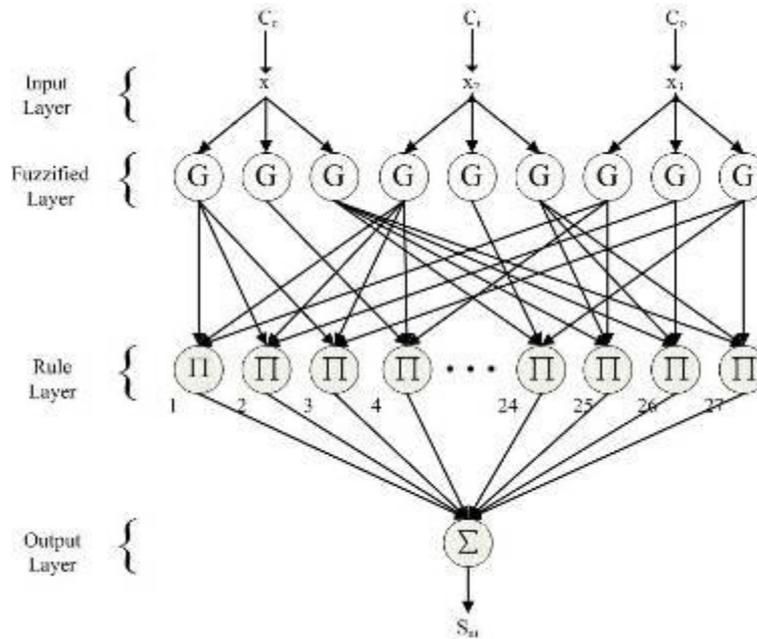


Fig. 4 The Structure of the FNN

As shown in Fig. 4, the FNN structure with three input variables, three input nodes of fuzzified layer for each input variable, 27 rule nodes of hidden layer, and one output node of output layer.

A typical format for a fuzzy rule base consists of a collection of fuzzy IF-THEN rules in the following form:

$$IF x_1 \text{ is } A_{111}^j, \dots; \text{ and } x_n \text{ is } A_{nml}^j, \text{ THEN } S_m^j \text{ is } \beta^j \quad (6)$$

where A_{111}^j and β^j are fuzzy sets and x_i , S_m^j are the input and output of the fuzzy inference rule, respectively.

Fuzzified Layer

This layer uses a Gaussian function as a membership function, so the output of the i th term node associated with x_i is:

$$\mu A_{ijk} = \exp\left(-\frac{(x_i - m_{ijk})^2}{\sigma_{ijk}}\right) \quad (7)$$

where m_{ijk} and σ_{ijk} denote the mean (center) and variance (width) of A_{ijk} , respectively. And i, j, k have the similar meaning as n, m, l in Eq. 6.

Rule Layer

This layer implements the links relating preconditions (fuzzified layer) to consequences (output layer).

The connection criterion is that each rule node has only one antecedent link from a fuzzified node of a linguistic variable. Hence there are 27 rule nodes in the initial form of FNN structure. We mention that there is still no weight adjustment in this layer. The output of the j th rule node is:

$$out_j^3 = \prod_{i=1}^n \mu A_{ikl}(x_i) \quad (8)$$

where the superscript 3 of out stands for the input number is 3, and i, k, l have the similar meaning as n, m, l in Eq. 6. Only noted that l is determined by the connection criterion.

Output Layer

All consequence links are fully connected to the output nodes and interpreted directly as the strength of the output action. This layer performs defuzzification to obtain the numerical output:

$$S_m = \sum_{j=1}^m \beta^j \prod_{i=1}^n \mu A_{ikl}(x_i) \quad (9)$$

where m is the number of fuzzy IF-THEN rules and n is inputs number.

Supervised Learning of Fuzzy Neural Network

The adjustment of the parameters in the proposed FNN can be divided into two tasks, corresponding to the IF (antecedent) part and THEN (consequent) part of the fuzzy inference rules. A simple and intuitive method of initializing the center and width for Gaussian functions is to use normal fuzzy sets to fully cover the input space. In this paper, we initialized these singletons based on the fuzzy inference method mentioned above.

A gradient-descent-based BP algorithm is employed to adjust FNN's parameters[15,16]. The goal is to minimize the error function:

$$E = \frac{1}{2} (d - S_m)^2 \quad (10)$$

where S_m is the output of the FNN and d is the desired output for the input pattern. If w_{ijk} is the adjusted parameter, then the learning rule is:

$$w_{ijk}(t+1) = w_{ijk}(t) - \eta \frac{\partial E}{\partial w_{ijk}} + \alpha \Delta w_{ijk}(t) \quad (11)$$

and

$$\Delta w_{ijk}(t) = w_{ijk}(t) - w_{ijk}(t-1) \quad (12)$$

where η is the learning rate and α ($0 < \alpha < 1$) is the momentum parameter.

In this paper, the sample data for training is a McGill calibrated color Image Database. The data-base provides a large number of color images of natural scenes, calibrated, for use in biological and computer vision research. The feature maps of image are calculated as explained above as input data. And the output data for training is the saliency value of the image calculated based on Itti's model[8] but adjusted according to the actual attention region given by user who look over the sample images.

Experimental Results

The initial structure of the FNN uses three input nodes for x_1, x_2 and x_3 , which stand for C_o, C_i and C_o , respectively. So in this case we have 27 initial rules. Suppose one epoch of learning takes 384 points. The supervised learning is continued for 500 epochs of training. The fuzzy sets for these linguistic term nodes are normally and uniformly initialized. We choose $\eta = 0.02$ and $\alpha = 0.85$ for supervised learning. The desired error d is got from the adjusted saliency map value calculated by Itti's method. Finally, the mean squared error (MSE) is 0.000497. The learning curve is illustrated in Fig. 5. From the figure we can see that the learning speed is very fast. This is because there are only 20 groups sample data used as inputs in this time.

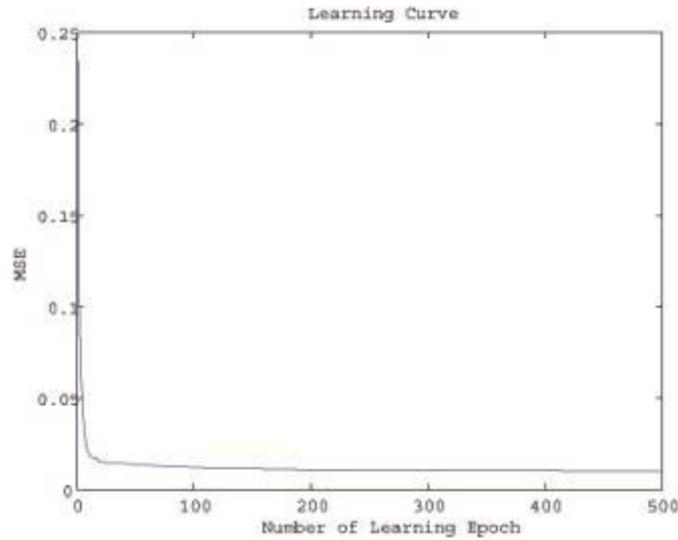
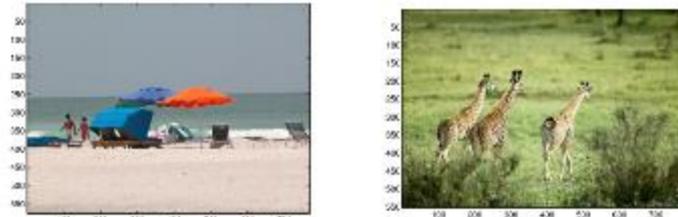


Fig. 5 The Learning Curve of the FNN

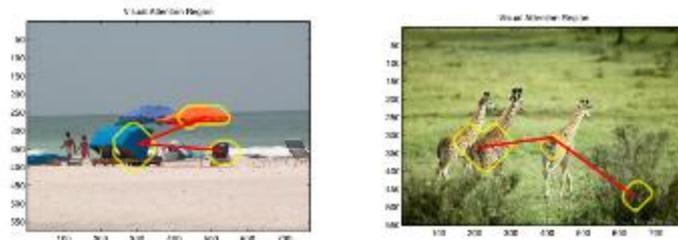
After the training of FNN, We conduct several experiments to demonstrate the inference result of the proposed methods and also compare with the performance of the proposed methods with Itti's model[8]. As mentioned in this paper, we get the various feature saliency maps at first. Here, two different images are used as the input images. The input image is processed for low-level features at multiple scales, and center-surround differences are computed according to Eq.(4). Then, the resulting feature maps are combined into feature saliency maps according to Eq.(5), which is shown in Fig. 6.



(a) Original Image



(b) Attention Region by Fuzzy Inference Method



(c) Attention Region by FNN Method

Fig. 6 Two Examples of Saliency Maps and Attention Regions

After getting the feature saliency maps, the region locations in the saliency map compete for the highest saliency value by the two methods proposed by us. After segmentation around the most salient region location,

this saliency map is used for obtaining a smooth object mask at image resolution and for object-based inhibition of return. The results of saliency map and attention region by both sum feature maps method and FNN method are shown in Fig. 6. The attention regions are marked by yellow lines while red lines express the order easy to be paid attention of them.

As we can see in the figure, there are only little differences between the saliency maps of two methods. And for the first example, the approximate locations of attention regions and the orders are basically the same between two methods while the little difference is the shape and size of region. This is because the color, intensity and orientation feature are all reflect obviously in regions marked of it compared with the rest regions, which also means that the differences of importance for the three features are small. So the method we proposed has not functioned very efficiently. But from the result of the second example we can see that the attention regions of our proposed method have better result. But only from these results we cannot yet say whether our proposed method is better than Itti's or not definitely. So we conduct another experiment to verify it.

The following experiment is asked 5 males who are between 20 to 30 years old to look over 20 images. After all the experiments they will be showed the results of attention region got by the two methods and asked to compare with the ones they actually attending and looking at in the experiment process. Finally, an evaluation of user's attitude to the results is carried on and shown in Fig. 7.

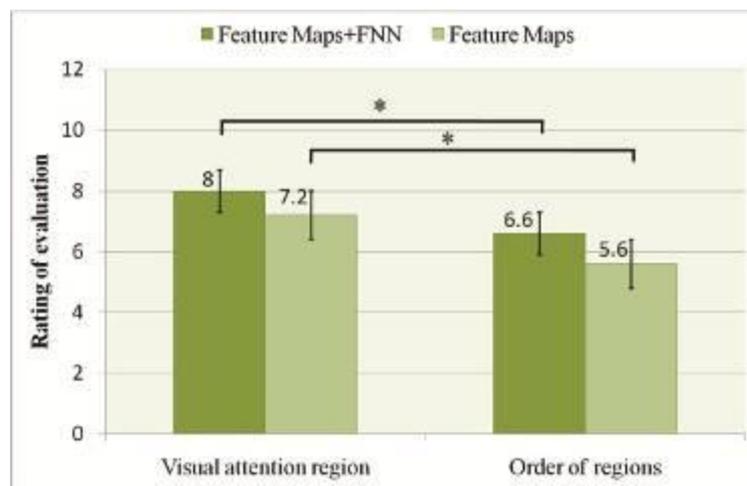


Fig. 7 Evaluation Value and Standard Deviation of Attention Regions Obtained by Two Methods

Every factor of them has five ranks and represented by 1~10 from worst to best in these figures. We can see from the evaluation results that the performance of our proposed method is higher at the order of visual attention region while a little lower at the accuracy of region detection. This also illustrated that proposed FNN method can improve the performance of attention region prediction at some aspect.

Conclusions

In this paper, we proposed a FNN method based on color, intensity and orientation feature maps of images to predict the visual attention regions and compared with the fuzzy inference method. We also conducted a series of attention region predict experiments. The prediction accuracy of our proposed approach was evaluated in experiments, and the results confirmed the effectiveness of our method in visual attention region prediction.

The problem still existing is the input image data for training of FNN is still less and the advantages of FNN have not reflected very well. And another problem is the method of getting the training data is not very improvement yet. In the future work, we will work on proposing a method to get more sample data for training FNN effectively and verify the method.

References

- [1] F. Sadri, Logic-Based Approaches to Intention Recognition, Handbook of Research on Ambient Intelligence: Trends and Perspectives, 2010.
- [2] D. V. Pynadath and M. P. Wellman, Accounting for Context in Plan Recognition, with Application to Traffic Monitoring, Proceedings of the Eleventh International Conference on Uncertainty in Artificial Intelligence. (1995) pp.472-481.
- [3] L. M. Pereira and H. T. Anh, Intention Recognition via Causal Bayes Networks Plus Plan Generation, Progress in Artificial Intelligence. (2009) pp. 138-149.
- [4] K. A. Tahboub, Intelligent Human-Machine Interaction Based on Dynamic Bayesian Networks Probabilistic Intention Recognition, Journal of Intelligent and Robotic Systems. Vol.45(2006), pp.31-52.
- [5] J. W. Harris and H. Stocker, Handbook of Mathematics and Computational Science, Springer-Verlag New York, 1998.
- [6] W. Mao and J. Gratch, A Utility-Based Approach to Intention Recognition, Proceedings of the AAMAS 2004 Workshop on Agent Tracking: Modeling Other Agents from Observations (2004).
- [7] L. Itti, C. Koch and E. Niebur, A Model of Saliency-Based Visual Attention for Rapid Scene Analysis, IEEE Trans. Pattern Analysis and Machine Intelligence. Vol.20(1998), No.11, pp.1254-1259.
- [8] L. Itti and C. Koch, A saliency-based search mechanism for overt and covert shifts of visual attention, Vision Research. Vol.40(2000), pp.1489-1506.
- [9] D. Walther, U. Rutishauser, C. Koch and P. Perona, On the usefulness of attention for object recognition, Workshop on Attention and Performance in Computational Vision. (2004) pp.96-103.
- [10] L. Itti, Models of bottom-up and top-down visual attention, PhD thesis, California Institute of Technology, 2000.
- [11] L. M. Hurvich and D. Jameson, An opponent-process theory of color vision, Psychological Review. vol.63, pp.384-404, 1957.
- [12] R. Manduchi, P. Perona and D. Shy, Efficient deformable filter banks, IEEE Transactions on Signal Processing. Vol.46(1998), No.4, pp.1168-1173.
- [13] W. O. Lee, J. W. Lee, K. R. Park, E. C. Lee and M. Whang, Object recognition and selection method by gaze tracking and SURF algorithm, 2011 International Conference on Multimedia and Signal Processing. (2011) pp.261-265.
- [14] G. B. Huang, Q. Y. Zhu and C. K. Siew, Extreme learning machine: a new learning scheme of feedforward neural networks, Neural Networks, 2004. Proceedings. 2004 IEEE International Joint Conference. Vol.2(2004) pp.985-990.
- [15] C. M. Lin, C. F. Hsu, Supervisory recurrent fuzzy neural network control of wing rock for slender delta wings, IEEE TRANSACTIONS ON FUZZY SYSTEMS. Vol.12(2004), No.5, pp.733-742.
- [16] A. Olmos, F. A. A. Kingdom, A biologically inspired algorithm for the recovery of shading and reflectance images, Perception. Vol.33(2004), No.12, pp.1463-1473.

Acknowledgment

The work is supported by Advanced Funds for Overseas Students Sci-tech Programs of Beijing Municipal Human Resources and Social Security Bureau.