

Photograph aesthetical evaluation and classification with deep convolutional neural networks

Yunlan Tan^{a,b,c}, Pengjie Tang^{b,c}, Yimin Zhou^c, Wenlang Luo^{a,b,*}, Yongping Kang^{a,b}, Guangyao Li^c

^a School of Electronics and Information Engineering, Jinggangshan University, Ji'an, Jiangxi, China

^b Key laboratory of watershed ecology and geographical environment monitoring, National Administration of Surveying, Mapping and Geoinformation, Ji'an, Jiangxi, China

^c College of Electronics and Information Engineering, Tongji University, Shanghai, China

ARTICLE INFO

Keywords:

Image aesthetics
Quality assessment
High aesthetics
Low aesthetics
Image classification
Deep convolutional neural network
GoogLeNet
Feature representation

ABSTRACT

In response to the growth of digital photography and its many related applications, researchers have been actively investigating methods for providing automated aesthetical evaluation and classification of photographs. For computational networks to recognize aesthetic qualities, the learning algorithms must be trained using sample sets of characteristics that have known aesthetic values. Traditional methods for developing this training have required manual extraction of aesthetic features for use in the practice datasets. With abundant appearance of convolutional neural networks (CNN), the networks have learned features automatically and have acted as important tools for evaluation and classification. At the time of our research, several existing convolutional neural networks for photograph aesthetical classification only used shallow depth networks, which limit the improvement of performance. In addition, most methods have extracted only one patch as a training sample, such as a down-sized crop from each image. However, a single patch might not represent the entire image accurately, which could cause ambiguity during training. What's more, for existing datasets, the numbers of high quality images of each category are mostly too small to train deep CNN networks. To solve these problems, we introduce a novel photograph aesthetic classifier with a deep and wide CNN for fine-granularity aesthetical quality prediction. First, we download a large number of consumer photographic images from DPChallenge.com (a well-known online photography portal) to construct a dataset suitable for aesthetic quality assessment. Then, we zoom out the images into 256×256 by bilinear interpolation and crop 10 patches (Center+four Corners+Flipping). Once we have associated the set with the image's training labels, we feed the images with the bag of patches into the fine-tuned networks. Our proposed computational method is configured to classify the photographs into high and low aesthetic values. A training pattern specifying an output of (0, 1) indicates that the corresponding image belongs to the “low aesthetic quality” set. Likewise, a training pattern with an output of (1, 0) indicates that the corresponding image belongs to the “high aesthetic quality” set. Experimental results show that the accuracy of classification provided by our method is greater than 87.10%, which is noticeably better than the state-of-the-art methods. In addition, our experiments show that our results are fundamentally consistent with human visual perception and aesthetic judgments.

1. Introduction

With the rapid growth of the imaging and mobile technologies, taking digital photographs has become a daily life activity these days. People want to take, share, and view photographs that have high aesthetic quality. According to statistics provided by DPChallenge.com [1], photographers have uploaded over 587,000 professional photographs. This convenience has stimulated the development of network sharing centers and portals that attract professional or amateur photographers from around the world. Therefore, automatic assessment of aesthetic quality of photographs is a promising technique in

many applications related to production, enhancement, management, retrieval, and recommendation of photographs. It is a challenge to establish the standards for differentiating between high and low quality images. First, visual data are very rich and ambiguous, because aesthetical assessment and prediction of the aesthetic value of images are highly subjective and not universal. Second, when judging photographs, people often apply their own cultural relativity to aesthetic judgments. Finally, even if we could gain agreement that certain degradations of photos (e.g., an image that is out of focus) are indicators of poor quality, it is far more difficult to find consensus about higher level positive visual properties such as color harmonies,

* Corresponding author at: School of Electronics and Information Engineering, Jinggangshan University, Ji'an, Jiangxi, China.

<http://dx.doi.org/10.1016/j.neucom.2016.08.098>

Received 17 February 2016; Received in revised form 16 August 2016; Accepted 31 August 2016

Available online xxxx

0925-2312/ © 2016 Elsevier B.V. All rights reserved.

layout, and lighting. With all these challenges, one might question the possibility of creating generic models that encode photographic preferences and then developing an automatic classifier that can learn and use these preferences accurately. Nevertheless, automatic assessment of photo quality from the perspective of visual aesthetics is of great interest in high-level vision research and has drawn much attention in recent decades [2].

Image features (e.g., low-level, high-level) are adopted by computer vision researchers for the purpose of objective image quality assessment [3], and many successful algorithms have been proposed for this purpose [4–8]. Wang and Datta were the first to realize the quantization of image features, including brightness, color distribution, wavelet, region composition and depth of field. By applying a support vector machine (SVM) or linear regression to distinguish high from low quality photographs, they achieved an accuracy of 70.12%. Luo et al. [5] proposed a photo quality assessment method that first extracted the subject region from a photo, and then formulated a number of high level semantic features for photograph quality classification. Ke et al. [6] designed some high level semantic features based on the spatial distribution of edges, blur, and the histograms of low-level color properties such as brightness and hue. They tested their features on a large, diverse dataset, and their system was able to achieve a classification rate of 72%. Tong et al. [7] tried to classify photos as either professional or snapshots, but they used the Corel image database, which is too homogeneous to separate the two classes. In addition, they simply collected a large set of low level features from the image retrieval literature. Liu et al. [8] presented a Multiple Kernel Learning (MKL) method for aesthetic image classification without using explicit feature selection steps. They produced better results than [6,7] using fewer features, but the results still had a classification rate of only 78.3% on a set of 6000 data. Sun et al. [9] tried to design optimized visual features to mimic human perception in photo assessment. Wong et al. [10] presented a saliency-enhanced image classification method, but this method must be preprocessed to detect the salient region. Zhao and Huang et al. [11] proposed a novel local binary count descriptor for rotation invariant texture classification. Guo et al. [12] used both hand-crafting and semantic features to improve the performance of automatic aesthetic assessment of images. Different from reference images, Zhang et al. [13] assessed the image quality with no reference image by using sparse feature representation. They extracted features from image structure patches and PCMs (pixel correlated matrix). Their method achieved comparable high accuracy.

Apart from image features extraction, it is helpful to develop new image enhancement tools to make images look better with photo aesthetic quality assessment [14–18]. Bhattacharya et al. [14] presented an interactive application that enables users to improve the visual aesthetics of their digital photographs by using spatial re-composition. Rather than prescribing a user-guided image segmentation and inpainting solution, Xu et al. [15] developed a photo-taking interface that provides real-time feedback on how to position the subject of interest according to the photography composition rule-of-thirds. Aydin et al. [16] presented a perceptually calibrated system for automatic aesthetic evaluation of photographic images based on a set of fundamental and meaningful aesthetic attributes such as sharpness, depth, clarity, colorfulness, and tone. They can give “aesthetic signature” of each image automatically and various photo editing applications. Wang et al. [17] built an image aesthetic classification and evaluation system with Hadoop cloud computing. Because of using parallel-processing strategy, the system shorts the processing time of images aesthetic analysis and has an efficient aesthetic evaluation for the application of mobile devices.

In addition, researchers have explored the optimization of image processing algorithms using perceptual measures of visual quality as objective functions [19–23]. Marchesotti et al. [19] proposed a method using BoVW (Bag of Visual Words) and FV (Fisher Vector) to encode the feature vectors using descriptors such as SIFT (Scale Invariant

Feature Transform). Their method outperforms previous methods by a significant margin. Their method learned the linear SVMs with a hinge loss using the primal formulation and a Stochastic Gradient Descent (SGD) algorithm, but they achieved just 68.55% accuracy. Jiang et al. [20] proposed a novel regression method named Diff-RankBoost based on RankBoost and support vector techniques. They predicted coarse-granularity aesthetic categories for more than 450 real consumer photographic images. Wang et al. [21] performed large scale studies analyzing algorithm performance, namely the structural similarity index (SSIM). They developed the ACQUINE [22] aesthetics value measurement system, an aesthetic evaluation and search engine. However, this method is impractical for large datasets due to the limits of SVM. Wu et al. [23] extended their SVM classification method to predict aesthetic adjectives rather than aesthetic scores. They introduced a probabilistic post-processing step that alleviated effects due to misleadingly labeled training data.

Most previous works on aesthetic image analysis have focused on designing appropriate features. Those methods suffer from several drawbacks. First, all attributes that are related to image aesthetic quality assessment cannot be pointed out. In fact, it is hard to discover all the attributes that can affect image aesthetic quality. Researchers just adopt a small number of them which are well-known and easy to be implemented. Second, it is hard to explain how those image features affect the image aesthetic quality. For example, we may agree on that image color will affect the image aesthetic quality a lot. Professional photographers carefully set the image color to gain better visual effect. But there is no fixed discipline that the color scheme of beautiful images must be followed. Both images with simple color palette and images which are colorful can have high aesthetic quality. Third, it is not easy to describe the image features accurately with mathematical model even we are sure that the features can affect the image aesthetic quality.

Since designing handcrafted features has long been regarded as an appropriate method for predicting image aesthetics, none of the above researchers delved further into the classification methods with regard to machine learning. More recently, new research efforts in deep learning methods have brought breakthroughs in many traditional computer vision problems, becoming one of the most powerful learning architectures for many vision tasks including object recognition, image classification, and video classification [24–28]. CNNs are designed for hierarchical feature representation mechanisms from lower level to higher level in which each level consists of a certain number of feature maps. The feature maps of each level are obtained from the maps contained in the previous level by applying several operations such as linear convolution, non-linear activation, and spatial pooling. Researchers who explored these methods [26–28] were able to show some improvement in automatic aesthetical assessment by modifying the network structures slightly (e.g., adding a layer or adding a column) or by adjusting the training strategy (e.g., fine-tuning). Besides useful techniques such as ReLU, dropout, and data augmentation introduced in [24], there are several efficient CNN approaches to improve performance in various classification problems. Wu et al. [29] modeled a weakly supervised, deep multiple instance learning framework, and they achieved convincing performance in vision tasks including classification and image annotation. To improve the geometric invariance of CNN activations, Gong et al. [30] presented a simple but effective scheme called multi-scale orderless pooling (MOP-CNN). Generally, the existing deep convolutional neural networks (DCNN) require a fixed-size input image. To eliminate the above constraint, He et al. [31] equipped the networks with a spatial pyramid pooling strategy. Their new network structure (SPP-net) could boost the accuracy of image classification, allowing their methods to rank third in image classification in ILSVRC 2014.

Multi-column neural networks [32,33] have been demonstrated to be an efficient approach to improving the performance of single-column neural networks in various classification problems. In multi-

column neural networks, one can also constrain the multi-column structures to share weights and aggregate multi-column outputs using max-pooling. Sermanet et al. [32] achieved the best performance compared with other reported results on all major pedestrian detection datasets. The model created by Ciresan et al. [33] reached near-human performance on the MNIST1 dataset. In [34,35], researchers computed ImageNet features (i.e., features extracted by the neural network trained in the ImageNet Challenge [24]) from a multi-scale image pyramid for object recognition, scene recognition, and object detection.

In image aesthetics quality assessment using deep learning, several attempts have been made to apply DCNN into image aesthetic evaluation. Dong et al. [36] implemented a deep convolutional neural network which has eight layers and millions of parameters to “teach” this network enough knowledge about images. The DCNN can “understand” images well and conduct the photo aesthetic quality assessment with a higher accuracy of 83.52%. In addition, Dong et al. [37] also applied image descriptors generated from their deep convolutional neural network with 8 layers of content-based aesthetic features. Further, they fused new aesthetic features to predict multi-level rather than binary photo quality. The MKL method performed quite well, obtaining an accuracy of 78.92%. Guo et al. [38] further proposed a paralleled convolutional neural network (PDCNN), which parallels two DCNNs with three convolutional layers and a DCNN with four convolutional layers, with multi-level structures to automatically adapt to the training dataset. The paralleling architectures of different complexity are able to improve the fitness for different scale datasets. Lu et al. [39] proposed a training approach that used a deep multi-patch aggregation network with two novel layers – a statistics layer and a sorting layer – to enable aggregation of multiple input sources. The proposed deep multi-patch aggregation network improved the performance on the AVA style dataset for image aesthetics categorization and image quality estimation on real-world photos with the accuracy of 75.41% and 89.2% respectively. Meanwhile, Lu et al. [40] developed a double-column deep convolutional neural network with four layers in each column to support input of global and local views to capture both global and local characteristics of images. They employed the style and semantic attributes of images to further boost the aesthetics categorization performance to an accuracy of 75.42%.

The foregoing works constructed the shallow depth CNN that achieved limited aesthetics categorization performance, while recent works found that steadily increasing the depth of networks by adding more very small convolutional layers could achieve a significant improvement. Krizhevsky [24], Lin [41], Szegedy [42], Simonyan [43], and He et al. [44] investigated the effect of the convolutional network depth on its accuracy in the large-scale image recognition setting. Their results inspire us to introduce a deep convolutional network in our work. Our research has been motivated by the feature learning power of deep convolutional neural networks [24,41–44], where feature learning is unified with classifier training using RGB images. To solve the challenge of bridging the “semantic gap” between visual features specific to automatic aesthetic perceptions of images and creating an effective learning model for image aesthetic quality assessment, we propose a novel deep neural network architecture to extract image features automatically rather than designing hand-crafted features. We formulate the learning problem by representing an input image with a small set or bag of patches [39] and associating the set with the image’s training label. We apply the deep convolutional neural networks with the same architecture as introduced in [42] to the image aesthetic quality assessment. This deep network is trained well on part of the ImageNet image database, which contains millions of images with various categories. This deep network can directly compress the image into a relatively lower dimensional feature vector, and meanwhile reserve most information in the image. For existing datasets, the numbers of high quality images of each category are mostly too small to train with deep CNN networks. Therefore, we have chosen to build a large and diverse training and testing database based

on the Web source DPChallenge.com. Our method achieves much better performance compared with state-of-the-art methods. The experimental results obtained on the large and reliable datasets further confirm our assumption.

Our research makes three contributions to this field. First, we highlight the developments in photograph aesthetical quality assessment that have resulted from features extracted manually and automatically. Second, different from those methods that used a ready-made database, we construct a large and diversified benchmark dataset suitable for the research of photograph aesthetical quality assessment. This dataset includes 47,304 photos with manually labeled ground truth. Third, in contrast to those methods that did not learn deep representations, we incorporate deep learning and a multiple instance framework, and we fine-tune the parameters of the classifier to gain a higher level of precision. The classifier can achieve fine-granularity aesthetic quality prediction with an accuracy of 87.1%. In the age of Internet, many websites want to recommend high quality pictures to users automatically. This paper presents an ideal method to solve the problem.

The remainder of this paper provides details of our research. In the next Section, Classification framework and training details of DCNN model training and learning is elaborated. In Section 3, the evaluation and classification dataset and setting are introduced. In Section 4, the experimental results are presented with a discussion. In Section 5, we present our conclusion.

2. Classification framework

2.1. Over architecture

We propose a deep network architecture to support scene details learning of photos, utilizing multiple patches cropped from one image based on GoogLeNet [42]. As depicted in Fig. 1, the input images from database are zoomed out into 256×256 by bilinear interpolation. Motivated by part-based approaches [39], we randomly crop regions of 224×224 from four corners and center of the images. The five cropped regions are flipped horizontally, and then the 10 different sub-crops of size 224×224 are put into the first convolutional layer filters. In the deep network architecture, the first convolutional layer filters the 224×224 input cropping regions with 64 kernels of size 7×7 with a stride of 2 pixels. (This is the distance between the receptive field centers of neighboring neurons in a kernel map.) In turn, the maxpooling layer takes the output of the first convolutional layer as input and filters it with 64 kernels of size 3×3 . After response-normalization layers, the second and third convolutional layers are connected to one another without any intervening pooling or normalization layers. The second convolutional layer has 64 kernels of size 1×1 with a stride of 1 pixel connected to the max pooling layer outputs. The third convolutional layer has 192 kernels of size 3×3 . After response-normalizing and max-pooling layers, there are 9 inception modules (described in Section 2.2.3). With inception modules and deeper layer-by-layer convolutions structure, the DCNN has sharply improved the recognition rate. Two auxiliary classifiers are added to connect to the inception module layers, which can increase the gradient signal that is propagated back and provides additional regularization. After the inception module operation, the fully-connected layers with a final 2-way softmax have 1024 neurons. The neurons in the fully connected layer are connected to all neurons in the previous average pooling layer. The outputs of the last fully-connected layer are fed to a 2-way softmax. The 2-way softmax produces a distribution over the aesthetic quality label (0 or 1), which denotes low or high aesthetic quality respectively.

2.2. Mechanism of DCNN model training and learning

Neural networks have emerged as important tools for classification.

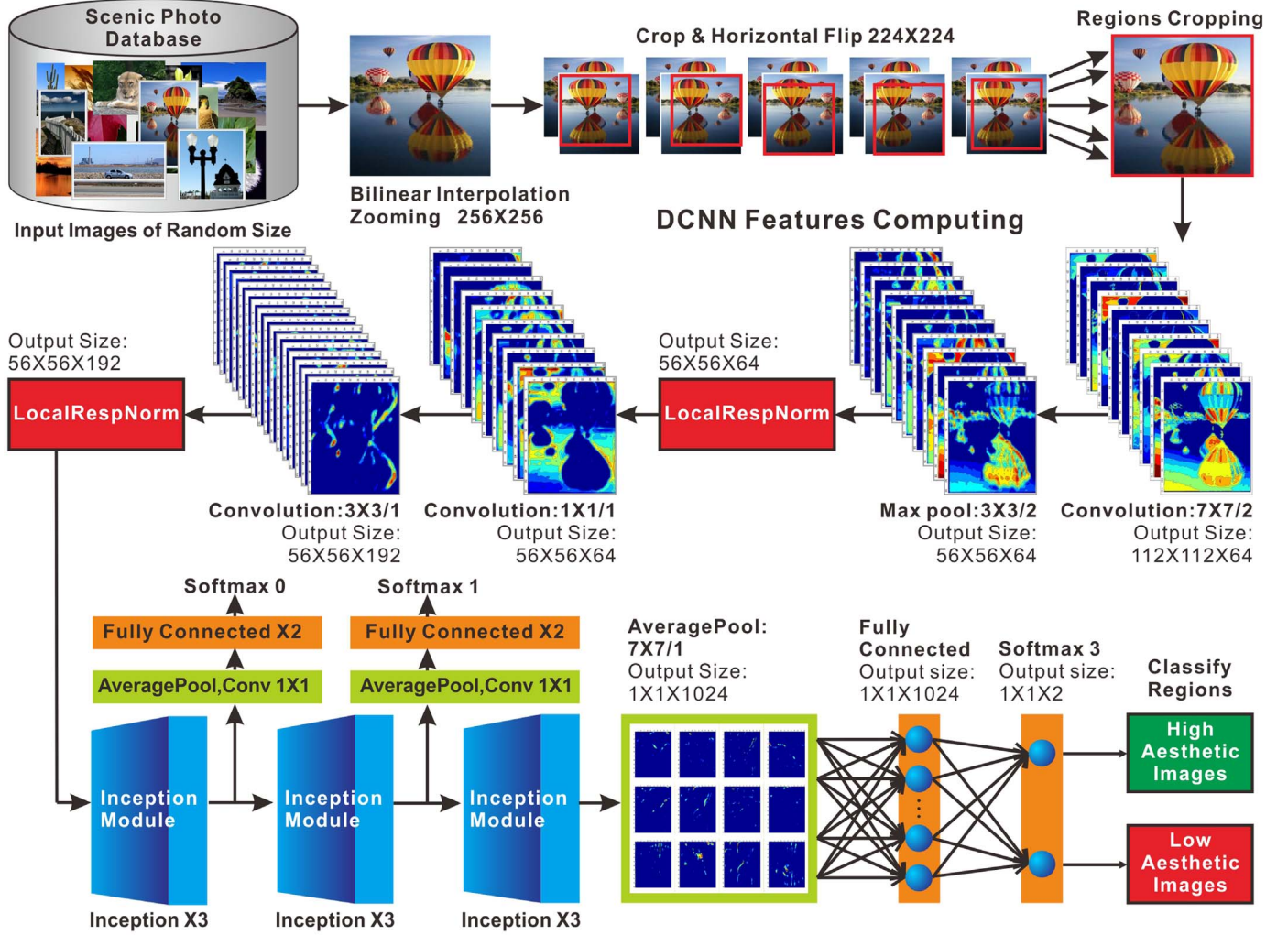


Fig. 1. The pipeline for photograph aesthetical evaluation and classification.

The advantage of neural networks lies in the following theoretical aspects. First, neural networks are data driven self-adaptive methods in that they can adjust themselves to the data without any explicit specification of functional or distributional form for the underlying model [45]. Second, there are universal functional approximators in the neural networks that can approximate any function with arbitrary accuracy [46,47]. We implement a deep convolutional neural network that has the same architecture as [42]. We change the 1000-way softmax into 2-way in the network. Inspired by Lu's bag of patches [39], we formulate the deep neural networks for photograph aesthetics assessment, which is considered as a smooth mapping function $g: X \rightarrow Y$ from dataset C . That is to say, the $224 \times 224 \times 3$ -dim feature vector is mapped into a 2-dim classification vector. Suppose we have the input training vector data and label part, denoted as X_i, Y_i , where $i = \{1, \dots, N\}$, with $X_i \in \mathbb{R}^d$ and $Y_i \in \mathbb{R}^d$. C is a collection of training examples, $C = \{(X_{p_i}, Y_i)\}_{i \in \{1, N\}}$, N is the size of the training samples. In multiple instance learning, data are organized as bags X_{p_i} with K patches cropped from source image X_i . Within each bag there are a number of instances $\{X_{p_{ij}}\}_{j \in \{1, K\}}$. $Y_i \in \{0, 1\}$ donates the associated aesthetics quality label from source image X_i . The patches are cropped from the four corners and the center, as well as their horizontal reflections (hence ten patches in all). Here K is set to 10. The mapping function g performs three steps sequentially: extracting features of individual patches in a bag, aggregating the features, and predicting the label of the bag. There are several key ideas and techniques in the network, including convolution feature map operation, backward

propagation learning mechanism and inception module [42].

2.2.1. Convolution feature map operation

Our network architecture is a hierarchical deep convolutional neural network that extracts local features by convolving input with a group of kernel filters. The convolutional layers generate feature maps by linear convolutional filters followed by nonlinear activation functions (rectifier). Outputs of neurons in the same layer form a plane called a feature map. The obtained convolutional feature maps are then sub-sampled (denoted as pooling) and filtered out to next layer. We can obtain different feature maps by setting different kernel filters for the local receptive field. Given $X_l^i \in \mathbb{R}^{M_l \times M_l}$ represents the i th feature map in the l th layer. The j th kernel filter in the l th layer connected to the i th map in the $(l-1)$ th layer is denoted as $k_l^{ij} \in \mathbb{R}^{K_l \times K_l}$ and index maps set $M_j = \{i\}$, i th in the $(l-1)$ th layer map connected to j th map in the l th layer. We involve the optimization of these bases into the optimization of the network. The convolution operation can be expressed as Rectified Linear Units (ReLUs), which are several times faster than their equivalents using tanh units in deep convolutional neural networks [24]. Then local maximum selection operations and pooling are conducted all over the feature map. Convolution operation and maximum selection operation can be described in Eq. (1).

$$X_l^i = \max \left(0, \sum_{i \in M_j} X_{l-1}^i * k_l^{ij} + b_l^j \right) \quad (1)$$

where $\max(\cdot)$ is ReLU non-linearity activation function, and b_l^i is bias. The feature $\text{map} X_l^i$ is activated after convolution operation $\sum_{i \in M_j} X_{l-1}^i * k_l^i + b_l^i$. We still find that the local response normalization scheme aids generalization. We apply this normalization after applying the ReLU nonlinearity in certain layers. $a_{x,y}^i$ donates the activity of a neuron computed by applying kernel i at position (x, y) . The response-normalized activity $b_{x,y}^i$ is given by the following expression (2):

$$b_{x,y}^i = a_{x,y}^i \left(k + \alpha \sum_{j=\max(0,i-n/2)}^{\min(NK-1,i+n/2)} (a_{x,y}^j)^2 \right)^\beta \quad (2)$$

where the sum runs over n “adjacent” kernel maps at the same spatial position, and NK is the total number of kernels in the layer.

The ordering of the kernel maps is arbitrary, of course, and is determined before training begins. This sort of response normalization implements a form of lateral inhibition modeled on the type found in real neurons. This lateral inhibition creates competition for big activities among neuron outputs computed using different kernels. The constants k, n, α and β are hyper-parameters whose values are determined using a validation set. In our experiment, we let $k = 2, n = 5, \alpha = 10^{-4}$ and $\beta = 0.75$.

Pooling equation can be described in Eq. (3).

$$X_l^i = \text{down}(X_{l-1}^i) \quad (3)$$

where $\text{down}(\cdot)$ is sub-sampling function to computer the max value of each $n \times n$ region in X_{l-1}^i map. In all cases, the convolutional layer and pooling layer appear alternately in the convolutional neural network.

2.2.2. Backward propagation learning mechanism

The output layer is fully connected with its previous layer. It produces a feature vector which can be transferred to a logistic regression layer to accomplish the recognition task. All the weights in the network are learned with the back-propagation method [40]. For back-propagation, the gradient of the l th convolutional layer is computed by Eq. (4):

$$\begin{cases} y_l = \omega_l x_l + b_l \\ x_l = f(y_{l-1}) \\ \Delta y_l = f'(y_l) \Delta x_{l+1} \end{cases} \quad (4)$$

where ω_l represents the weights of a filter, b_l is a vector of biases, and y_l is the response at a pixel of the output map. $f(\cdot)$ is the activation, and f' is the derivative of f . The update rule for weight ω_l is given by the expression (5).

$$\begin{cases} \mu_l^{i+1} = \alpha \mu_l^i - \lambda \cdot \eta \cdot \omega_l^i - \eta \cdot \left\| \frac{\partial \varepsilon}{\partial \omega_l} \right\|_{\omega_l^i} \Delta y_l^i \\ \omega_l^{i+1} = \omega_l^i + \mu_l^{i+1} \end{cases} \quad (5)$$

where i is the iteration index, α is the momentum factor, μ is the momentum variable, λ is the weight decay, η is the learning rate, and $\left\| \frac{\partial \varepsilon}{\partial \omega_l} \right\|_{\omega_l^i}$ is the average over the i th batch D_i of the derivative of the objective function ε with respect to ω_l , evaluated at ω_l^i .

We train the model using stochastic gradient descent with a batch size of 40 examples. We find that the small amount of weight decay is important for the model to learn. Weight decay can reduce the model's training error, which was fine-tuned to 0.0005 in our experiment. We also find that dropout is a technique that prevents overfitting in training neural networks [24]. Typically, dropout and momentum can improve learning [48]. Since applying dropout to all layers causes a significant increase in the time to reach convergence, in our experiments we set dropout for the fully connected layers with ratio of 0.5 or 0.6, $\alpha = 0.9, \lambda = 0.0005$.

As softmax regression problem, we describe the last fully connected layer for multiple instance aesthetics classification. In our network, the

last fully connected layer is followed by a softmax2 layer. Given one training sample X_i , the network advances layer by layer to extract representations, starting from the first convolutional layer to the output of the last fully connected layer $fc_{27} \in \mathbb{R}^m$. The last layer can be viewed as containing high level features of the input image. Donote $X_{p_i} = \{X_{p_{ij}} | j = 1, \dots, K\}$ as a bag of K instances and $t = \{t_i | t_j \in \{0, 1\}, j = 1, \dots, m\}$ as associated aesthetics quality label of X_{p_i} for objects of m categories. A multiple instance convolutional neural network extracts representations of the bag: $h = \{h_{ij}\} \in \mathbb{R}^{m \times K}$, in which each column is the representation of an instance. The aggregated representation of the bag for DCNN is given by the expression (6).

$$\hat{h} = f(h_{i1}, h_{i2}, \dots, h_{iK}) \quad (6)$$

where function f is $\text{avg}_j(h_{ij})$. fc_{27} is transformed into a probability distribution $P \in \mathbb{R}^m$. Cross entropy is used to measure the prediction loss L of the network. Specifically, we have the expression (7).

$$P_i = \frac{\exp(\hat{h}_i)}{\sum_i \exp(\hat{h}_i)}, \text{ and } L = - \sum_i t_i \log(P_i) \quad (7)$$

where L is the loss of cross-entropy. The loss function of the DCNN can be minimized by using a stochastic gradient descent (SGD) algorithm during the training process of DCNN. The gradients of the deep convolutional neural network are calculated via back-propagation [40], which is given by the expression (8):

$$\frac{\partial L}{\partial \hat{h}_i} = P_i - t_i \text{ and } \frac{\partial \hat{h}_i}{\partial h_{ij}} = \begin{cases} 1, & h_{ij} = \hat{h}_i \\ 0, & \text{else} \end{cases} \quad (8)$$

2.2.3. Inception module

The new local inception module is another characteristic of GoogLeNet [42], which introduces 9 inception modules. Each Inception module is made up of 1×1 convolutions, 3×3 convolutions, 5×5 convolutions, and 3×3 max pooling. Moreover, 1×1 convolutions are applied to computer reductions for involving with less number of parameters as well as rectified activation before the expensive 3×3 and 5×5 convolutions. The structure is shown as Fig. 2.

The basic idea of the inception module is to find the optimal local construction and to repeat it spatially. This architecture aims to increase significantly the number of units at each stage without causing an uncontrolled blow-up in computational complexity [49]. Thanks to the inception modules, the DCNN can be designed not only very deeply but also efficiently. With 9 inception modules, we can extract local feature representation using flexible convolutional kernel filter sizes with layer-by-layer structure, which is proved to be robust and effective for the large scale high-resolution images. Furthermore, owing to the padding strategy and precise designs, after the inception module operation, we can obtain a number of feature maps of the same sizes, though by means of different scale convolutions as well as pooling and the feature maps are concatenated together by a concat-layer following by each inception module.

2.3. Complexity analysis of DCNN model

The total complexity of DCNN model is comprised of three parts: parameters scale, time complexity and number of neurons. The first part is the parameters scale, which almost designed in fully-connected layers and convolutional layers. The second part is time complexity, which also depends on convolutional computation and fully-connected computation. Given N_{paras} represents the number of parameters scale and T_{Comp} represents the model computational complexity [50], the parameters scale and computational complexity can be described by Eqs. (9) and (10) respectively.

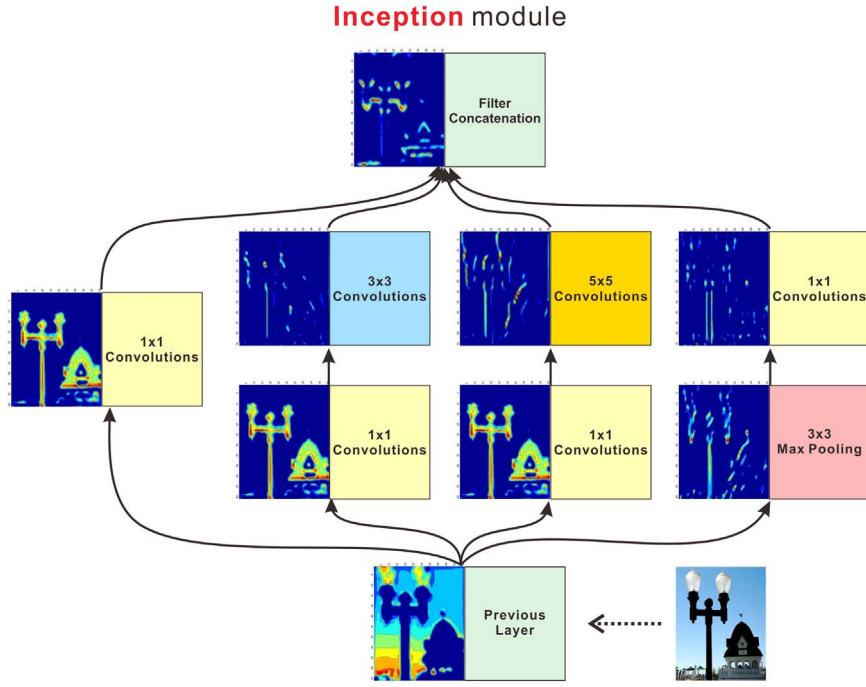


Fig. 2. The architecture of one inception model.

$$N_{paras.} = \sum_{l=1}^d (k_l^h \times k_l^w \times m_{l-1}^{nc} \times m_l^{nc}) \quad (9)$$

$$TComp. = \sum_{l=1}^d (m_{l-1}^{nc} \times k_l^h \times k_l^w \times m_l^{nc} \times m_l^h \times m_l^w) \quad (10)$$

Where d is depth of the DCNN (layers of convolution), k_l^h and k_l^w are the height and width of convolutional kernels (generally, $k_l^h = k_l^w$), m_l^h and m_l^w represent the height and width of l th feature map in the l th layer respectively. m_l^{nc} is the number of the former channel and m_{l-1}^{nc} is the number of the next channel. When $l=1$ and RGB color channels, m_{l-1}^n is set to 3. p_l^h and p_l^w represent the number of horizontal and vertical spatial padding of feature maps with convolution stride and s_l^w respectively. When $l > 3$, we compute m_l^h and m_l^w in Eq. (11):

$$\begin{cases} m_l^h = \lceil (m_{l-1}^h - k_l^h + 2p_l^h) / s_l^h \rceil + 1 \\ m_l^w = \lceil (m_{l-1}^w - k_l^w + 2p_l^w) / s_l^w \rceil + 1 \end{cases} \quad (11)$$

The third part is number of neurons. With d layers put together, the number of neurons $N_{neus.}$ can be described by Eq. (12).

$$N_{neus.} = \sum_{l=1}^d (m_l^n \times m_l^h \times m_l^w) \quad (12)$$

In this paper, the GoogLeNet model is 22 layers deep when counting only layers with parameters, i.e., $d = 22$. The overall number of layers (independent building blocks) used for the construction of the network is about 100. The model input images with random size are resized to 256×256 . Those cropped regions are 224×224 , that is $m_0^h = m_0^w = 224$. k_l^h and k_l^w are set to the same size at same layer (generally, $k_l^h = k_l^w = 1, 3, 5 \text{ or } 7$). All input images are RGB color ones, that is $m_0^n = 3$. In addition, several complexities of CNN models are counted. Table 1 shows the neurons number and parameters scale of different deep models. GoogLeNet model has the scale of 6.8 M Parameters, which is 4.7% of VGG19 model and the number of neurons of the model is also 37% of VGG19 model. In spite of having a large depth, the number of weights in the GoogLeNet model is less than the number of weights in other shallower networks. This is why we select GoogLeNet model as the generic image descriptors for aesthetic quality assessment.

Table 1

The neurons number and parameters scale of different CNN models.

Model type	Neurons number (K)	Parameters scale (M)	Layers
Alex-Net [24]	650	62.5	9
GoogLeNet [42]	5500	6.8	27
ZFCNN [25]	1470	63.2	8
VGG19 [43]	14,860	144	19

3. Evaluation and classification dataset constructing and setting

3.1. Dpchallenge dataset constructing

The image datasets use to study aesthetics typically consist of photographic images shared on social networks. Online communities such as photo.net and DPChallenge.com gather a large number of expert and amateur photographers who share, view, and judge photos online. These photographers also agree on the most appropriate annotation policy to score the images. Such policies can include a scale of numerical values (ratings). From these annotations, images can be labeled as being visually appealing or not. This type of rating scale allows for a fair, quantitative evaluation of the different methods of automated aesthetical evaluation. Since CUHK only contains 12,000 images [6], the number of high quality images of each category is too small to train CNN networks. The number of votes per image in AVA database [51] ranges over 78 and the images are not latest. We believe that the visual aesthetics of images with higher votes can be understood better by hobbyists and professionals. Therefore, it is necessary to construct a dataset with high quality images to train deep CNN networks. For our evaluation, we construct database in our experiments. Images mined by our Web crawlers from DPChallenge.com have been voted on by a large number of users and can be considered as users' favorites. Therefore, they can be regarded reasonably as professional quality photos. From a statistical point of view, these images must convey some information and common structures that reflect the characteristics of a professional level image. Because high quality photos with differing subject matter can still have underlying aesthetic properties in common, we focus on mining and characterizing the

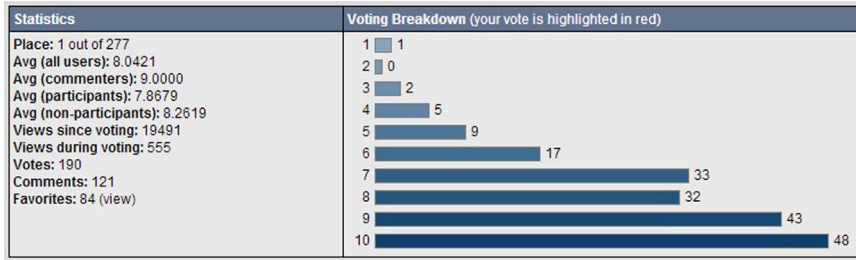


Fig. 3. The screenshot of statistics attributes for ID. 473436.

underlying aesthetic rules for both landscape and nature photos. The galleries contain user scores of 587,136 images with a grade ranging from 1 to 10. We download 22,104 photographs from the category of landscapes and 28,913 photographs from the nature gallery. All photos are voted on by at least 100 users. Then we delete those gray images by manual. There are 20,114 photographs of landscape category and 27,190 photographs of nature gallery left. Fig. 3 is a screenshot showing the statistical attributes provided for each photo; in this case for sample photo ID.473436 (http://www.dpchallenge.com/image.php?IMAGE_ID=473436). As shown in this example, each image is associated with a distribution of scores that corresponds to individual votes. We select one of statistical attributes, Avg (all users), as the ground-truth for all the aesthetic values of photographs used in our experiments.

Since aesthetics scores above 8.0 or below 3.0 rarely appear in this dataset, we select the observed median value 5.5 to be the median aesthetics rating. We set the gap $\delta = 1.0$, different from Datta's $\delta = 1.6$ [4]. δ is used to create a gap artificially between high and low quality images because pictures lying in this gap are likely to represent noisy data in the peer-rating process. In all our classification experiments, we use a scale in which ratings $\geq 5.5 + \delta/2$ are very good and ratings $\leq 5.5 - \delta/2$ are very bad. Accordingly, we consider images with ratings ≥ 6.0 as high aesthetic value and images with ratings ≤ 5.0 as low aesthetic value. In addition, we consider ratings > 5.5 & ≤ 6.0 as common good aesthetics, and ratings ≥ 5.0 & ≤ 5.5 as common bad aesthetics. Images with ratings > 5.5 are counted as good images and those with ratings ≤ 5.5 counted as bad ones in all our classification experiments. In turn, these ratings are associated with a corresponding aesthetic quality label: 0 refers to a low-quality (bad) image, and 1 refers to a high-quality (good) image. With these standards in mind, we construct a set of 11,089 high aesthetics images and 9025 low ones in the category of landscape, and a set of 12,874 high aesthetics images and 14,316 low ones in the category of nature. The number of different aesthetics value interval is shown in Fig. 4. There are 47,304 images in all in our database. The positive samples 23,963 and negative examples 23,341 are roughly equivalent, which is suitable for applying in aesthetic quality assessment.

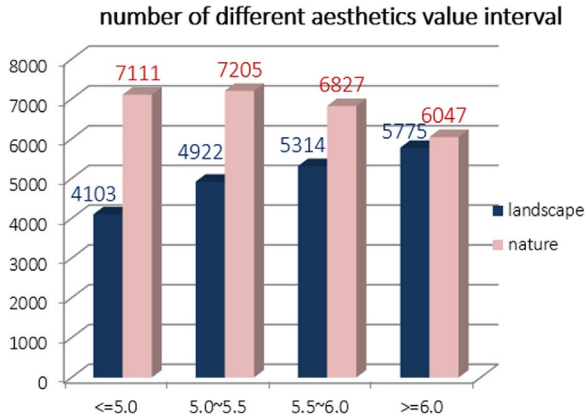


Fig. 4. Number of different aesthetics value interval on landscape and nature datasets.

3.2. Classification challenge setup

To reduce time of constructing DCNN, we utilize a pre-trained model file, googlenet_place205_train_iter_2400000.caffemodel, for individual photograph aesthetic evaluation and classification. To train and test the GoogLeNet models, we follow the multi-view classification method [31,52]. Before training, the high-resolution images with random size are resized to 256×256 . Specifically, we randomly crop regions of 224×224 from four corners and center of the image. After that, the cropped regions are horizontally flipped. Therefore, we obtain 10 views, each of which is fed into GoogLeNet models for prediction. We start with training the 27-layer GoogLeNet [42], where network weights are well-initialized. We fine-tune the DCNN on photograph aesthetics identification training set using stochastic gradient descent with a batch size of 40 to update the parameters. For better classification performance, the learning rate is set initially as 0.01. The momentum is set to 0.9. To reduce the effect of over-fitting, the training is regularized by weight decay of 0.0005 and dropout for the fully connected layers with ratio of 0.5 or 0.6. The training stops after 20,000 or 30,000 iterations. We test the model on tens of thousands of color images depicting typical scenes of natural landscapes by 2 GTX Titan-X GPUs. The experiments are scheduled in ubuntu14.04 and caffe platform. The code is written in MATLAB2013a and clean, efficient C++ with opencv3.0 used for GPU computation.

4. Experimental results and discussions

4.1. Comparison with different policy

In order to compare visual features that show correlation with community-based aesthetics scores, we experiment with a variety of ratings classifications for training and testing on the dataset. The experiments are showed in Table 2. The results of four experiments are convincing. In Classification experiment1, when images ratings about $score > 6.0$ & $score < 5.0$ are trained and the rest ratings about $score > 6.0$ & $score < 5.0$ are tested, the accuracy is 87.10%. In Classification experiment2, when images ratings about $5.0 \leq score \leq 6.0$ are added to test at the same condition of Classification experiment3, the accuracy is decrease to 74.13%. In Classification experiment3, when images ratings about $score > 6.0$ & $score < 5.0$ & $5.0 \leq score \leq 6.0$ are trained and the rest ratings about $score > 6.0$ & $score < 5.0$ are tested, the accuracy is 86.05%. In Classification experiment4, when adding images ratings about $5.0 \leq score \leq 6.0$ at the same condition on Classification experiment3, the accuracy is decreased to 75.12%. As expected, when images ratings about $score > 6.0$ & $score < 5.0$ are trained and tested, the best accuracy obtained on our dataset is even above 87%.

The DCNN has a 7×7 effective receptive field over the input. It incorporates three non-linear rectification layers instead of a single one, which makes the decision function more discriminative. A unique structure called the Inception module allows for increases in both the depth and the width of our networks without getting into computational difficulties. In some research such as [9,39,41], the authors often picked out the top 10% of the images from their respective datasets as

Table 2
Comparison with Different Policy by proposed method on DPChallenge collection.

Experiments	Policy	Accuracy
Classification experiment 1	training: score > 6.0 & score < 5.0, samples:21,036 testing: score > 6.0 & score < 5.0, samples:2000	87.10%
Classification experiment2	training: score > 6.0 & score < 5.0, samples:21,036 testing: score > 6.0 & score < 5.0 & 5.0 ≤ score ≤ 6.0, samples:4000	74.13%
Classification experiment3	training: score > 6.0 & score < 5.0 & 5.0 ≤ score ≤ 6.0, samples: 43,308 testing: score > 6.0 & score < 5.0, samples:2000	86.05%
Classification experiment4	training: score > 6.0 & score < 5.0 & 5.0 ≤ score ≤ 6.0, samples: 43,308 testing: score > 6.0 & score < 5.0 & 5.0 ≤ score ≤ 6.0, samples:4000	75.12%

Table 3
Comparison between our method and state-of-the-art methods.

Methods	Accuracy (%)
VGG19 [43]	71.34
Marchesotti's [19]	68.55
Dong's [37]	78.92
Lu's [39]	75.41
Lu's [40]	75.42
Dong's [36]	83.52
The proposed method	87.10

the positive class and the bottom 10% of the images as the negative class. As shown on Table 2 and Table 3, our Classification experiment 1 achieves accuracy of 87.10%, which is better than the state-of-the-art methods. This result suggests that the classifier would benefit strongly from large scale datasets with high discriminable images inputs.

4.2. Accuracies with different dropout under different iterations

Dropout is a technique that prevents overfitting in training neural networks. In [42], a dropout with 70% ratio of dropped outputs is applied in the last fully-connected layer. In our work, we adopt dropout and shuffle the training data in each epoch to alleviate overfitting. In

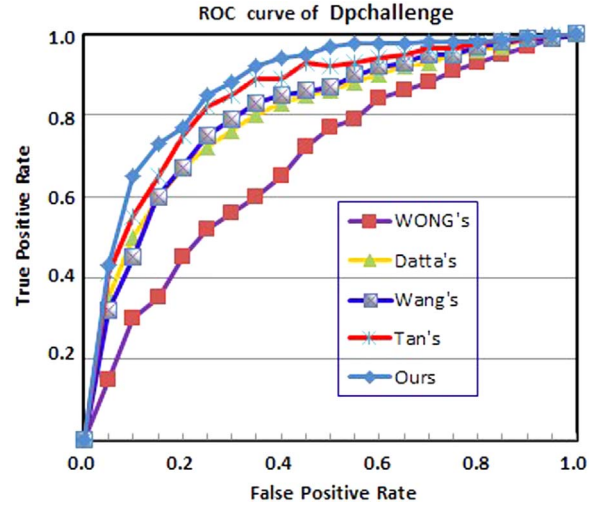


Fig. 6. Photo quality assessment performance comparisons with different methods. False Positive Rate and True Positive Rate denote positive sample classification accuracy and error rates respectively.

this type of research, typically the outputs of neurons are set to zero with a probability of 0.5 in the training stage and divided by 2 in the test stage. In our experiments, we fine-tune the dropout rate as 0.5 or 0.6 for the binary classification problem with a batch size of 40. As we increase dropout, training takes slightly longer to converge. The training has mostly converged by 20,000 iterations Or 30,000 iterations. The classifier based on DCNN is robust enough to produce good accuracy in separating high and low rated photographs after 20,000 and 30,000 iterations. Fig. 5 shows the accuracies with different dropouts under different iterations. These experiments demonstrate that the best accuracy obtains on the DPChallenge.com collection is above 88%.

4.3. Evaluation with ROC curve

To further compare the classification performance of our proposed method with other methods, the receiver operating characteristic (ROC) curves are plotted by adjusting the performance of the classifier. The ROC curve is a direct representation of the performance of the classifier: the larger area beneath of the ROC curve, the better performance of the classifier. In Fig. 6, we show the ROC curves resulting from the research of Wong [10], Datta [4], Wang [17] and our own. As can be seen from the figure, the area under our ROC curve is much larger than the area under Wong's curve and slightly larger than the area under the curve obtained from Wang's and Datta's methods. This result shows that our proposed method outperformed the state-of-

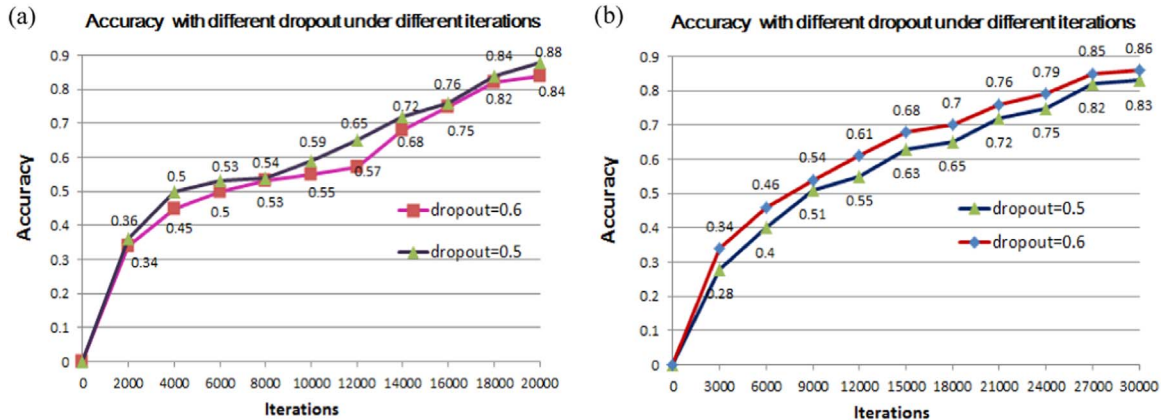


Fig. 5. (a) Accuracy on Classification experiment1 under 20,000 iterations (b) Accuracy on Classification experiment3 under 30,000 iterations.

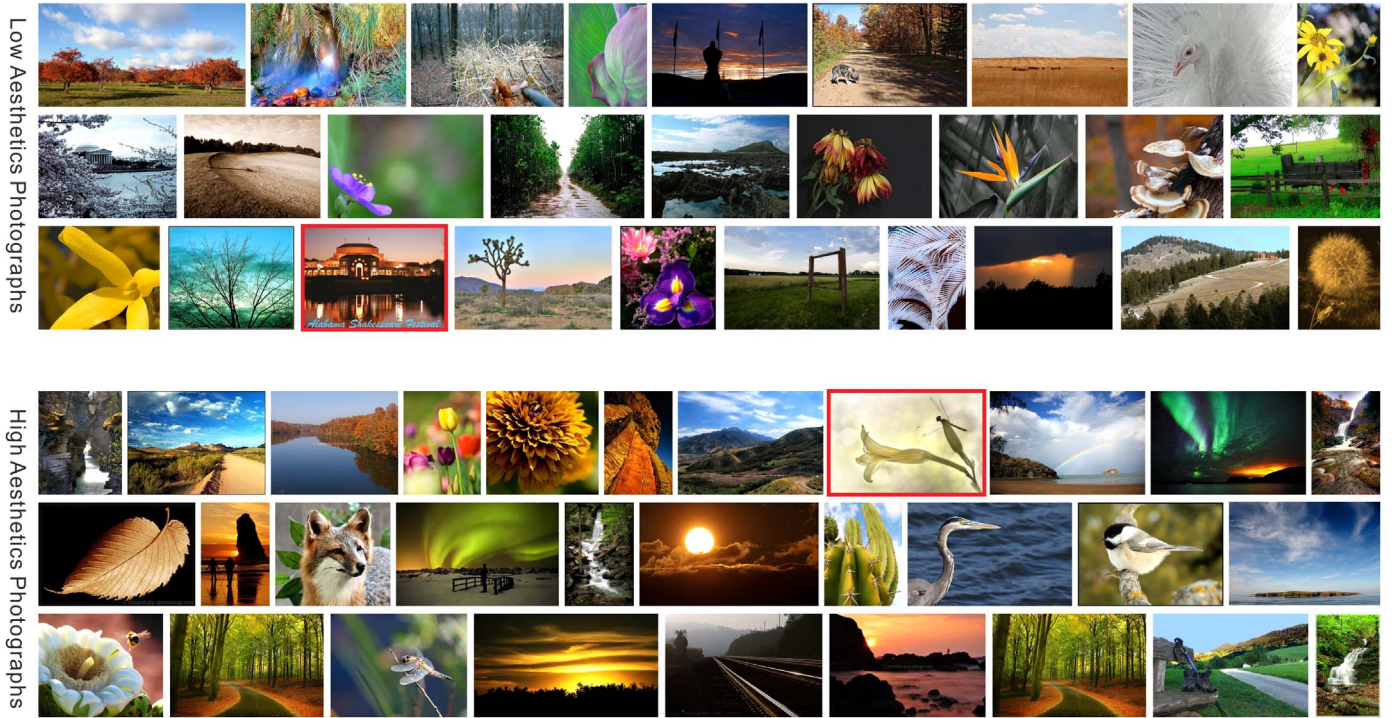


Fig. 7. Photographs of high aesthetics quality and low esthetics quality.

the-art methods significantly.

4.4. Comparison between our method and state-of-the-art methods

We compare the proposed method with VGG19's [43], Marchesotti's [19], Dong's [37], Xin's [39], Xin's [40] and Dong's [36] on the AVA dataset or part categories of it. Table 3 illustrates the comparison with different methods.

To prove our method's better performance, we make an experiment using VGG19 [43] (another DCNN method) on the same constructing dataset. The classifier VGG19 achieves an accuracy of 71.34%. Clearly, our method based on a DCNN scheme provides the best performance, which is 15.76% higher than VGG19 [43]. This comparison demonstrates that the representation depth of our method is beneficial for the classification accuracy. Dong's [36] method using deep convolutional neural network achieved the second place with an accuracy of 83.52%, which can "understands" images well. They [37] also constructed another deep convolutional neural network to predict multi-level photo quality with an accuracy of 78.92%. Lu's method [40], which used a double-column deep convolutional neural network, achieved the aesthetics categorization performance accuracy of 75.42%. Lu's DMA-Net approach [39] alone has improved the better performance on the AVA style dataset for image style classification, image aesthetics categorization with the accuracy of 75.41%. Marchesotti's method [19], which used a general image descriptor, achieved an accuracy of 68.55%. Based on Table 3, we can see that our method achieves the best performance, scoring 3.58% higher than Marchesotti's [19], 3.58% higher than Dong's [37], and 11.68% higher than Lu's [40]. It should be noted that our dataset and AVA [51] are different from those used by these other researchers, but the data are all collected from the same Web source. Images in our constructed database are downloaded directly from the website with careful selection. In compare with AVA, our labels inevitably have less noise inside. These results show that our constructed dataset makes it possible for our method to learn features automatically and assess image aesthetics using our deep learning approach, thereby providing significantly improved performance.

4.5. Quality assessment of high aesthetics quality and low aesthetics quality

Fig. 7 shows the 30 images that receive high aesthetic quality ratings and the 28 images that receive low aesthetic quality ratings under our proposed method. All high and low aesthetics images are from the Dpchallenge.com dataset that we have constructed. Our proposed method can classify more than 87.1% of images into the right category. The image descriptor directly learns discriminant features from normalized multiple instance image pixels and directly compress the image into a relatively lower dimensional feature vector. Meanwhile, the deep network can reserve most information in the image. The experiment results show that the accuracies of aesthetic quality prediction all can greatly be improved. As show by Fig. 7, the high esthetic photographs as rated by our method are more beautiful than the low esthetic photographs, a result mostly consistent with human visual perception.

The ground-truth labels are displayed in the form of colored frames (white frames for images assessed as correct, and red frames for images assessed as incorrect). In Fig. 7, note that a false positive image is considered to have high aesthetics quality in the column of Low Aesthetics Photographs. The reason for that incorrect categorization is that some regions from the image might not be considered as beautiful or have as bright colors as other regions. The second false positive image is considered to have low aesthetics quality in the column of High Aesthetics Photographs. The reason for that incorrect result is that the image might have been considered as having low contrast or as lacking in highlights, depth, or texture.

5. Conclusion

In this paper, we introduce a deep and wide neural network architecture to solve the challenge of bridging the "semantic gap" between the emotion-related concepts and visual features specific to image aesthetic quality assessment. The deep convolutional neural network with the same architecture as GoogLeNet [42] is trained to "understand" images well using part of the ImageNet image database

containing millions of images in various subject categories. This deep network can compress the image directly into a relatively lower dimensional feature vector, and meanwhile reserve most information in the image. The input images are provided with a bag of patches [39]. Once we have associated the set with the image's training labels, we feed the images with the bag of patches into the fine-tuned networks. The experimental results obtained on the large and reliable datasets further confirm our assumption. Our method achieves much better performance compared with state-of-the-art methods. In the future, we will further apply multi-level methods to combine our DCNN features with other hand-crafted aesthetic features for achieving a human-level performance.

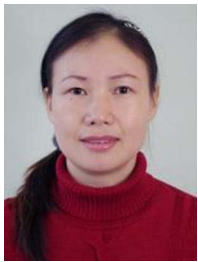
Acknowledgements

This work is supported by the Scientific Research Foundation of the Education Bureau of Jiangxi Province (No. GJJ150788), the Key Laboratory of Watershed Ecology and Geographical Environment Monitoring, NASG (WE2016015 and WE2016013), the Art Science Program of Jiangxi Province (No. YG2015081), the Natural Science Foundation of Jiangxi (20151BAB207016) and the National Key Technology Research and Development Program of the Ministry of Science and Technology (No. 2012BAC11B01).

References

- [1] Dpchallenge, (<http://www.dpchallenge.com>).
- [2] W. Luo, X.G. Wang, X.O. Tang, Content-based photo quality assessment, in: Proceedings of the 13th IEEE International Conference on Computer Vision (ICCV2011), Barcelona, Spain, November 2011, pp. 2206–2213.
- [3] Y.L. Tan, Y.M. Zhou, G.Y. Li, A.M. Huang, Computational aesthetics of photos quality assessment based on improved artificial neural network combined with an autoencoder technique, *Neurocomputing* 188 (2016) 50–62.
- [4] R. Datta, D. Joshi, J. Li, J.Z. Wang, Studying aesthetics in photographic images using a computational approach, in: Proceedings of the 9th European Conference on Computer Vision (ECCV2006), Graz, Austria, May 2006, pp. 288–301.
- [5] Y.W. Luo, X.O. Tang, Photo and video quality evaluation: Focusing on the subject, in: Proceedings of the 10th European Conference on Computer Vision (ECCV2008), Marseille, France, October 2008, pp. 386–399.
- [6] Y. Ke, X.O. Tang, F. Jing, The design of high-level features for photo quality assessment, in: Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR2006), June 2006, pp. 419–426.
- [7] H.H. Tong, M.J. Li, H.J. Zhang, J.R. He, C.S. Zhang, Classification of digital photos taken by photographers or home users, in: Proceedings of the 5th Pacific Rim Conference on Multimedia (MM2004), Tokyo, Japan, November 2004, pp. 198–205.
- [8] N.N. Liu, X. Jin, H. Lin, D. Zhang, Aesthetic image classification based on multiple kernel learning, in: Proceedings of CCF Chinese Conference (CCCV 2015), Xi'an, China, September 2015, pp. 229–236.
- [9] X.H. Sun, H.X. Yao, R.R. Ji, S.H. Liu, Photo assessment based on computational visual attention model, in: Proceedings of the 17th ACM international conference on Multimedia (MM2009), New York, USA, October 2009, pp. 541–544.
- [10] L.K. Wong, K.L. Low, Saliency-enhanced image aesthetics class prediction, in: Proceedings of the 16th IEEE International Conference on Image Processing (ICIP2009), Cairo, Egypt, November 2009, pp. 997–1000.
- [11] Y. Zhao, D.S. Huang, J. Jia, Completed local binary count for rotation invariant texture classification, *IEEE Trans. Image Process.* 21 (10) (2012) 4492–4497.
- [12] L.H. Guo, Y.C. Xiong, Q.H. Huang, X.L. Li, Image esthetic assessment using both hand-crafting and semantic features, *Neurocomputing* 143 (2014) 14–26.
- [13] C. Zhang, J. Pan, S.T. Chen, T.T. Wang, D.J. Sun, No reference image quality assessment using sparse feature representation in two dimensions spatial correlation, *Neurocomputing* 173 (2016) 462–470.
- [14] S. Bhattacharya, R. Sukthankar, M. Shah, A framework for photo-quality assessment and enhancement based on visual aesthetics, in: Proceedings of the 18th ACM International Conference on Multimedia (2010MM), Firenze, Italy, October 2010, pp. 271–280.
- [15] Y. Xu, J. Ratcliff, J. Scovell, G. Speiginer, R. Azuma, Real-time guidance camera interface to enhance photo aesthetic quality, in: Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems (CHI 2015), Seoul, Korea, December 2015, pp. 1183–1186.
- [16] T.O. Aydin, A. Smolic, M. Gross, Automated aesthetic analysis of photographic images, *IEEE Trans. Vis. Comput. Graph.* 21 (1) (2015) 31–42.
- [17] W.N. Wang, W.J. Zhao, C.J. Cai, J.X. Huang, X.M. Xu, L. Li, An efficient image aesthetic analysis system using Hadoop, *Signal Process.: Image Commun.* 39 (2015) 499–508.
- [18] W.H. Kim, J.H. Choi, J.S. Lee, Subjectivity in aesthetic quality assessment of digital photographs: analysis of user comments, in: Proceedings of the 23rd Annual ACM Conference on Multimedia Conference (2015MM), Brisbane, Australia, October 2015, pp. 983–986.
- [19] L. Marchesotti, F. Perronnin, D. Larlus, G. Csurka, Assessing the aesthetic quality of photographs using generic image descriptors, in: Proceedings of the 13th International Conference on Computer Vision (ICCV2011), Barcelona, Spain, November 2011, pp. 1784–1791.
- [20] W. Jiang, A.C. Loui, C.D. Cerosaletti, Automatic aesthetic value assessment in photographic images, in: Proceedings of the IEEE International Conference on Multimedia and Expo (ICME2010), Suntec, Singapore, July 2010, pp. 920–925.
- [21] Z. Wang, A.C. Bovik, H.R. Sheikh, E.P. Simoncelli, Image quality assessment: from error visibility to structural similarity, *IEEE Trans. Image Process.* 13 (4) (2004) 600–612.
- [22] R. Datta, J.Z. Wang, ACQUINE: aesthetic quality inference engine-real-time automatic rating of photo aesthetics, in: Proceedings of the International Conference on Multimedia Information Retrieval (MIR2010), Philadelphia, Pennsylvania, March 2010, pp. 421–424.
- [23] Y. Wu, C. Bauckhage, C. Thureau, The good, the bad, and the ugly: Predicting aesthetic image labels, in: Proceedings of the 20th International Conference on Pattern Recognition (ICPR2010), Istanbul, Turkey, August 2010, pp. 1586–1589.
- [24] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, in: Proceedings of Advances in neural information processing systems (NIPS 2012), Nevada, USA, December 2012, pp. 1097–1105.
- [25] M.D. Zeiler, R. Fergus, Visualizing and understanding convolutional networks, in: Proceedings of the 13th European Conference on Computer Vision (ECCV2014), Zurich, Italy, September 2014, pp. 818–833.
- [26] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR2014), Columbus, USA, June 2014, pp. 580–587.
- [27] L. Kang, P. Ye, Y. Li, D. Doermann, Convolutional neural networks for no-reference image quality assessment, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR2014), Columbus, USA, June 2014, pp. 1733–1740.
- [28] Y. Taigman, M. Yang, M. Ranzato, L. Wolf, DeepFace: closing the gap to human-level performance in face verification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR2014), Columbus, USA, June 2014, pp. 1701–1708.
- [29] J.J. Wu, Y.N. Yu, C. Huang, Y. Kai, Deep multiple instance learning for image classification and auto-annotation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR2015), Boston, USA, June 2015, pp. 3460–3469.
- [30] Y. Gong, L. Wang, R. Guo, S. Lazebnik, Multi-scale orderless pooling of deep convolutional activation features, in: Proceedings of the 13th European Conference on Computer Vision (ECCV2014), Zurich, Italy, September 2014, pp. 392–407.
- [31] K. He, X. Zhang, S. Ren, J. Sun, Spatial pyramid pooling in deep convolutional networks for visual recognition, in: Proceedings of the 13th European Conference on Computer Vision (ECCV2014), Zurich, Italy, September 2014, pp. 346–361.
- [32] P. Sermanet, K. Kavukcuoglu, S. Chintala, Y. LeCun, Pedestrian detection with unsupervised multi-stage features learning, in: Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR2013), Portland, USA, June 2013, pp. 3626–3633.
- [33] D. Ciresan, U. Meier, J. Schmidhuber, Multi-column deep neural networks for image classification, in: Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR2012), Providence, USA, June 2012, pp. 3642–3649.
- [34] F. Iandola, M. Moskewicz, S. Karayev, R. Girshick, T. Darrell, K. Keutzer, DenseNet: Implementing efficient Convnet descriptor pyramids. Technical report, University of California, Berkeley, arXiv preprint [arXiv:1404.1869](https://arxiv.org/abs/1404.1869)[cs.CV], 2014.
- [35] M. Koskela, J. Laaksonen, Convolutional network features for scene recognition, in: Proceedings of the 22nd ACM International Conference on Multimedia (MM2014), Orlando, USA, November 2014, pp. 1169–1172.
- [36] Z. Dong, X. Shen, H.Q. Li, X.M. Tian, Photo quality assessment with dcnn that understands image well, in: Proceedings of the 21st International Conference on Multimedia Modeling (MMM2015), Sydney, Australia, January 2015, pp. 524–535.
- [37] Z. Dong, X.M. Tian, Multi-level photo quality assessment with Multi-view features, *Neurocomputing* 168 (2015) 308–319.
- [38] L.H. Guo, F.D. Li, Image Aesthetic Evaluation Using Paralleled Deep Convolution Neural Network, arXiv preprint [arXiv:1505.05225](https://arxiv.org/abs/1505.05225)[cs.CV], 2015.
- [39] X. Lu, Z. Lin, X.H. Shen, R. Mech, J.Z. Wang, Deep multi-patch aggregation network for image style, aesthetics, and quality estimation, in: Proceedings of the 15th IEEE International Conference on Computer Vision (ICCV2015), Santiago, Chile, December 2015, pp. 990–998.
- [40] X. Lu, Z. Lin, H.L. Jin, J.C. Yang, Rating image aesthetics using deep learning, *IEEE Trans. Multimed.* 17 (11) (2015) 2021–2034.
- [41] M. Lin, Q. Chen, S.C. Yan, Network in network, arXiv preprint [arXiv:1312.4400](https://arxiv.org/abs/1312.4400) [cs.NE], 2013.
- [42] C. Szegedy, W. Liu, Y. Q. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, Going deeper with convolutions, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR2015) Boston, USA, June 2015, pp. 1–9.
- [43] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-scale Image Recognition, arXiv preprint [arXiv:1409.1556](https://arxiv.org/abs/1409.1556), 2014.
- [44] K.M. He, X.Y. Zhang, S.Q. Ren, J. Sun, Delving deep into rectifiers: Surpassing human-level performance on imagenet classification, in: Proceedings of the 15th IEEE International Conference on Computer Vision (ICCV2015), Santiago, Chile, December 2015, pp. 1026–1034.
- [45] D.S. Huang, *Systematic Theory of Neural Networks for Pattern Recognition* 28, Publishing House of Electronic Industry of China, Beijing, China, 1996, pp. 323–332.

- [46] D.S. Huang, Radial basis probabilistic neural networks: model and application, *Int. J. Pattern Recognit. Artif. Intell.* 13 (07) (1999) 1083–1101.
- [47] D.S. Huang, J.X. Du, A constructive hybrid structure optimization methodology for radial basis probabilistic neural networks, *IEEE Trans. Neural Netw.* 19 (12) (2008) 2099–2115.
- [48] G.E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, R.R. Salakhutdinov, Improving Neural Networks by Preventing co-adaptation of Feature Detectors, *arXiv preprint arXiv:1207.0580*, 2012.
- [49] Z. Zhong, L. Jin, Z. Xie, High performance offline handwritten chinese character recognition using GoogLeNet and directional feature maps, in: *IEEE Proceedings of the 13th International Conference on Document Analysis and Recognition (ICDAR2015)*, Tunis, August 2015, pp. 846–850.
- [50] K. M. He, J. Sun, Convolutional neural networks at constrained time cost, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR2015)* Boston, USA, June 2015, pp. 5353–5360.
- [51] N. Murray, L. Marchesotti, F. Perronnin, AVA: a large-scale database for aesthetic visual analysis, in: *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR2012)*, Providence, USA, June 2012, pp. 2408–2415.
- [52] W.F. Liu, D.C. Tao, Multiview hessian regularization for image annotation, *IEEE Trans. Image Process.* 22 (7) (2013) 2676–2687.



Yunlan Tan received her B.Sc. degree in Computer Science from Jiangxi Normal University and her M.Sc. degree in Computer Application Technology from East China Normal University, China, in 1996 and 2004 respectively. Now she is an associate professor in the School of Electrical and Information Engineering, Jinggangshan University, Jiangxi, China and she is pursuing her Ph.D. degree in the College of Electrical and Information Engineering, Tongji University, Shanghai, China. Her current research interests include image processing and machine learning.



Pengjie Tang received his B.S. degree in Computer Science and Technology from Jinggangshan College and his M.S. degree in Computer Software and Theory from Nanchang University, Jiangxi, China, in 2006 and 2009 respectively. Now he is working at Jinggangshan University, Jiangxi, China, as a lecturer, and he is a Ph.D. candidate at the Department of Computer Science and Technology, Tongji University, Shanghai, China. His current research interests include multimedia intelligent computing, deep learning and computer vision.



Yimin Zhou received his B.S. degree from Shanghai Jiaotong University, Shanghai, China, in 2002. He received his M.S. degree in Tongji University, Shanghai, China, in 2005. Now he is pursuing his Ph.D. degree in the College of Electrical and Information Engineering, Tongji University, Shanghai, China. His research interests concentrate on image processing, machine learning and computer graphics.



Wenlang Luo received his B.Sc. degree in Physics from Beijing Normal University and his Ph.D. degree in Atomic and Molecular Physics from Sichuan University, China, in 1998 and 2009 respectively. Now he is a professor in the School of Electrical and Information Engineering, Jinggangshan University, Jiangxi, China. His current research interests include scientific computation and numerical simulation.



Yongping Kang received his M.Sc. degree of engineering in Computer Technology from Tongji University in 2010. Beginning with 2008, he is the director of network information center of Jinggangshan University. His current research interests include image processing, panorama, Virtualization and information technology.



Guangyao Li received the B.Sc., M.Sc. degrees and the Ph.D. degree from Nanjing University of Aeronautics and Astronautics in 1986, 1989 and 1997, respectively. Currently, he is a professor and Ph.D. supervisor in the School of Electrical and Information Engineering, Tongji University, Shanghai, China. His main research interests include graphics and images processing, and virtual reality.