

Using Convolutional Neural Networks to demystify aesthetic works of art

Pujun Bhatnagar
Stanford University
353 Serra Mall, Stanford, CA, 94305
pujun@cs.stanford.edu

Abstract

When one visits a museum and looks at some piece of art, she/he not only understands the subject matter that is being depicted, but is also able to analyze other ‘subjective’ qualities such as the aesthetic quality, emotional intensity, creativity etc. associated with the piece. And after 1-2 short glimpses, the person is able to come up with a description that aptly describes it. All this is possible because humans have a notion of aesthetic appeal of art works. In this paper, these abilities are transferred to a computer. Different types of Convolutional Neural Networks architectures and Siamese Networks are used to train models that hierarchically understand and intelligently predict the degree of aesthetics of an image. I visualize the learned weights in order to understand the underlying features that makes some art works so appealing and am able to perform better than the state of the art techniques that have been previously used to tackle this problem.

1. Introduction

The field of Computer Vision often deals with problem related to classification or localization of objects in images. However, when humans look at images, they not only just see objects, but are able to analyze and make comments about the subjective quality of an image, talk about its aesthetic quality, its emotional intensity and comment about artist’s creativity. Unlike computers, which with the advent of Deep Learning architectures, often treat images as digital dark matter that can be mined to learn hidden structure and representations and recognize, classify and localize images and objects, humans go further and are able to easily describe these images with words.

1.1. Subjective Qualities of Images

Many psychological case studies have been conducted to provide insight into the human brain to find how humans evaluate a piece of art. These show that for most humans, looking at a photo or a piece of art results in psychological

responses in various parts of the brain. [12]. These have been found to be correlated to ‘quality’ of an image, which can be split into objective and subjective quality of images.

The objective qualities of images can be related to the photo properties like exposure, color, contrast, background and sharpness. For instance, case studies show that a photo is fully appreciated when it is well exposed and not-blurry. Blurry shots and over or under-exposed photos are generally regarded as bad shots [8].

The remaining subjective qualities of the images, as the name suggests are harder to quantify and these are the attributes which make this regression problem challenging. For instance, some case studies show that people, when asked about any image, often attribute aesthetics to questions like ‘Is the composition balanced?’ or ‘How is the color distribution in the image?’ More often than not, there are no nice mathematical ways to capture these subjective qualities as they vary from person to person [9].

Therefore, it is easy to see that even though the objective qualities can be easily quantified in order to predict which photos are aesthetically pleasing, we still haven’t found ways to capture the subjective qualities of images. There are some heuristics, like *The Rule of the Third* and *The Golden Angle*, that attempt to capture these notions with varying degree of success[12]. Finally, photos are also associated to people’s memory, which may contribute to their appeal. Both extrinsic and intrinsic image features for making an image memorable have been studied [9].

1.2. Problem Statement

This paper aims at providing insights to the understanding of human being’s appreciation of art and beauty; aims at answering the pivotal question ‘*What is it that makes an image/art a magical mystical work that blows people away when they look at it?*’ Using the hierarchical learning capabilities of a Convolutional Neural Networks, the models learn these unique set of features that contribute to the works’ popularity and, when given an input image, are able to assign a score to the image based on its aesthetic appeal. Towards the end, the developed model is used to predict the

scores of renaissance artworks to see how do these learned features do against world renowned artworks.

1.3. Challenges

One of the biggest challenges in estimating image aesthetics is that most current techniques only utilize low level features, that are unable to capture the higher level representations that may contribute to aesthetics. Also, traditional techniques, being very imperative and non-hierarchical in nature, are unable to figure out the best ways to combine and integrate the lower level features to assign correct score. Finally, before 2012, when AVA dataset was released, there was no proper dataset, with enough detailed meta-data about aesthetics, that could be used to train such models.

2. Related Work

Some studies, aiming to quantify the aesthetic quality of images, have been done in the past. Most of the early approaches tried to approach these problems by using hand-crafted features based on photographic intuitions [8]. For instance, the use of low level features such as spatial distribution of edges, color frequency distributions, contrast in addition to using some high level features like object placement, orientation templates depth of view etc. have been used to train classifiers to solve this problem. Most of these approaches work well with their own dataset but fail to deliver similar results on more generalized data. [12]

More recent works in this domain have started using Convolutional Neural Networks. For instance, Lu et al. showed state of the art aesthetic classification performance using convolutional Neural Network [9]. In 2012, Murray et al. introduced the AVA which a large scale dataset of images with aesthetic and style rating. In their original work, they formulated a binary classification problem and trained SVMs with Fisher Vector signatures to achieve 67% accuracy [11]. Lu et al. also used the same AVA dataset and same experimental settings as [11] and were able to obtain 71% accuracy on this dataset by using single column deep convolutional networks. The convolutional layer contained 64 kernels, size of 11, 5, 3 and 3 respectively and the fully connected layer contained 1000, and 256 neurons respectively. In 2014, the RAPIDNet, was able to outperform existing models and techniques in this domain [10]. Their most accurate architecture contained 4 convolutional layers and 2 fully connected layers. They achieved this by using a specialized dual column network where one column is the same as their original aesthetics classification network and the second column is used to classify image style which then combine into a aesthetic prediction [9].

3. Technical Approach

3.1. Problem formulation

Based on the AVA dataset, the problem statement turns itself into a regression problem where each of the input image would be given a score from 1 to 10 based on its aesthetic evaluation. However, since CNNs rarely work well on regression problems, I decided to discretize the scores into 10 buckets, 1...10, thus transforming this into a classification problem. Additionally, in order to compare my results with past results, these scores were also transformed into two labels: "HIGH" or "LOW" (i.e. as a binary classification problem) based on a threshold δ , where, any $score \leq \delta$ was labeled as LOW while $score > \delta$ was treated as HIGH. The initial value of δ was chosen to be 5.

3.2. Dataset

For this project, AVA dataset [11], containing 250 thousand labeled images was used. This data is an amalgamation of photos that were submitted to a worldwide photography competition, where each contestant submitted entries in various competitions and the members of the community voted on each of the entries. Each image is labeled on a scale of 1 to 10 based on overall appeal for the given contest theme by an average of 200 site members [11].

It was decided to use this dataset for a variety of reasons. First, it has enough images that I did not need to worry about having too little data. Second, various research studies have already been done using this dataset which can serve as good baselines and comparisons for the project. Finally, while starting to work on this problem, this was the only dataset that offered so much meta-data in a well-defined format.

The dataset was obtained by scraping from **DPchallenge.com** using a python script. Here is the distribution of the resulting data and the associated statistics:

For more information about the dataset, one can refer to [11]. It contains a detail comparison of AVA with other datasets and documents good reasons why AVA is a suitable choice when training a model for image/art aesthetics.

4. Experiment

4.1. Data Pre-processing

As highlighted in [12][8], there are objective and subjective characteristics that lead to an overall aesthetic effect. Since the input is not of any default size or aspect ratio, data pre-processing was done for the model to properly pick up these features. Firstly based on 231N lectures, the mean was subtracted from all the samples. Also, Based on [9] suggestions, the images were re-sized in two ways: first, by warping the images to be 256 X 256 without considering the aspect ratio (i.e. simply down-sampling to 256



Figure 1: Preview of AVA data set

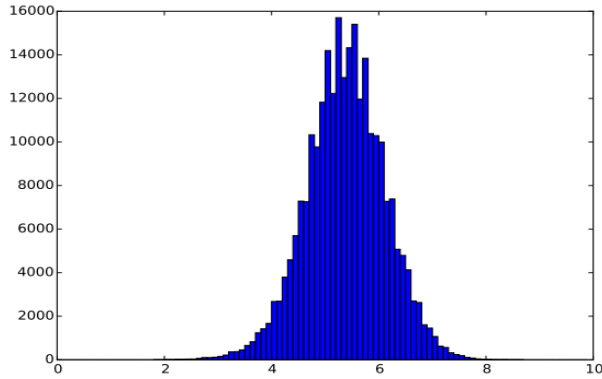


Figure 2: Histogram showing rating meta-data of AVA dataset

X 256). Second: by doing data augmentation. 10 random crops of size 256 X 256 from Top-Left, Top-Right, Bottom-Left, Bottom-Right and Center of the image and taking a mirror image of these while preserving the associated label were also taken as an input.

4.2. Models

4.2.1 Naive Baseline: A Non-CNN model

Using WEKA data mining and machine learning software [4], a 3 layer Neural net model and a support vector regressor were trained to perform regression task on the aesthetics of images. These models served as baselines for the more sophisticated CNN models. PCA was performed on the input images to do away with dimensions with little variation and train these models faster.

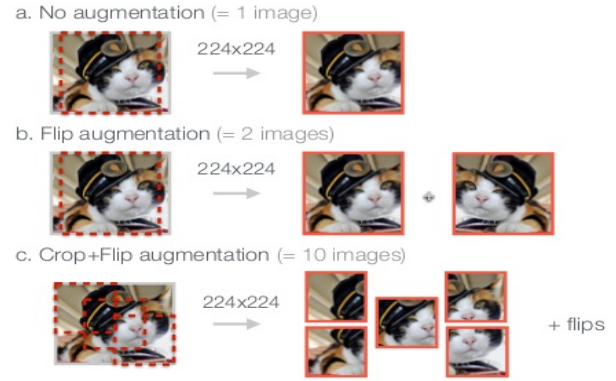


Figure 3: Data Augmentation used to increase performance

4.2.2 Better Baseline- Machine Learning applied on specific feature

Studies in the past have used different properties of images to come up with heuristics for image aesthetics [3]. Therefore, a stronger baseline was developed using these features, which included image sharpness (calculated using Gradient and Gradient Entrophy methods described in [8]), color distribution, SIFT BOW descriptor and salient maps.

4.2.3 CaffeNet

CaffeNet model was the first CNN model that was used to do the classification task. It is available in Caffe modelzoo [7] and is similar to the AlexNet architecture. It contains 5 convolutional layers and 3 fully connected layers. To perform this task, the last fully connected layer was replaced with a FC layer, initialized with random weights, and fine tuned on the AVA dataset.

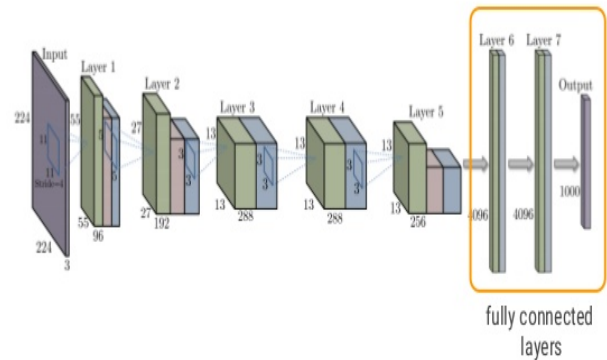


Figure 4: CaffeNet

Additionally, to compare my experiments with the previously published results, the scores were converted into HIGH and LOW labels, i.e. binary classification task. The

deep network was initialized with ImageNet weights to and was trained over 120k epochs with batch size of 128.

4.2.4 VGGNet

After playing with CaffeNet, it was decided to move to a deeper architecture in hope of seeing performance gains. Therefore, I downloaded VGGNet and repeated the same steps as CaffeNet. VGGNet is a 16-layer (13 convolutions and 3 fully connected layers) network from Simonyan et al. [13]. Even with this, the last fully connected layer was replaced with a Soft max loss layer, initialized with random weights, and fine tuned on AVA dataset. Again, as done with the CaffeNet, a separate copy of VGGNet was trained to do classification task.

Training this model was more cumbersome than imagined and the first couple of trails failed in achieving any accuracy (when compared to CaffeNet). After investigation, it was found that this network architecture was much more sensitive to learning rate and the initial learning rate was too high which made the loss blow up. Therefore, I went with a learning rate of 10^{-6} . Also, since this model is more computationally intensive than CaffeNet, I trained it for 100k epochs.

NOTE: I also tried using PReLU instead of ReLU for both the CaffeNet and VGGNet, while keeping all other hyper-parameters constant to study its affect on accuracy.

4.3. Convolutional Neural Network with ‘Deeper’ channels

This was inspired [3] and by the intuition that the current industry standard is to use features like color distribution, composition, edge distribution. Therefore, what if, for every image, these features were computed and augmented to make them more deeper. After looking into various features and based on [5], the following features were augmented as additional channels to every image:

1. Salient Map as channel 4
2. Color distribution Map as channel 5
3. sharpness Map as channel 6

I decided to make a convolutional network with deeper 8 conv layer architecture followed by 2 fully connected layer. Also, since I was using $256 * 256$ warped images with 6 channels (augmented information), I came up with the following architecture:

This was trained in Python using numPy (see the end of the paper for the source code).

4.3.1 Siemesse Networks

This was inspired by [9]. The generic structure of Convolutional Neural Network, containing M **Conv-ReLU-Pool** followed by N **Affine Layers + regression** layer, remained the same but two CNNs were used instead of one. The intuition

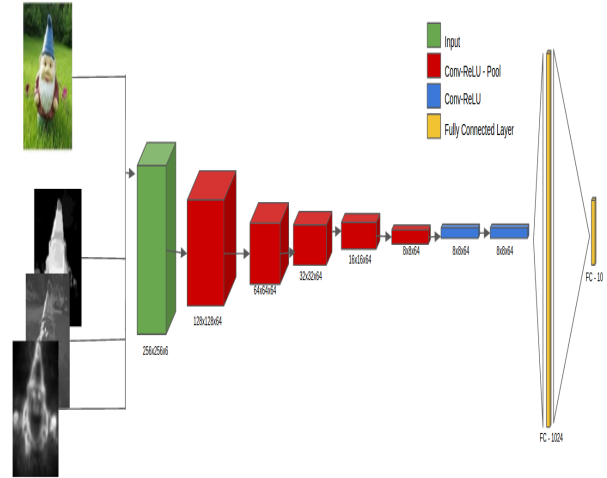


Figure 5: Integrating global and local information into deeper channels and running the CNN

Layer	Parameters	Channels	ReLU	POOL	Layer*
Data	256x256	6			
CONV1	3x3x6	64	Yes	Yes	128x128x64
CONV2	3x3x64	64	Yes	Yes	64x64x64
CONV3	3x3x64	64	Yes	Yes	32x32x64
CONV4	3x3x64	64	Yes	Yes	16x16x64
CONV5	3x3x64	64	Yes	Yes	8x8x64
CONV6	3x3x64	64	-	-	8x8x64
CONV7	3x3x64	64	-	-	8x8x64
CONV8	3x3x64	64	-	-	8x8x64
FC-1	FCx1024				
FC-2	FCx10				

Table 1: Custom CNN network trained using numPy and Python

was that since the convolutional network was used to solve the task of predicting image aesthetics, which was modeled as returning a number between 1-10, a L-2 loss was used in the last layer while partially training (only a few layers weights’ were allowed to be updated) and fine-tuning the model.

Research and techniques described in these areas, especially those covered in CS231N lectures and by Lu.et al [10], were significantly used while playing with data augmentation and to fine tune the parameters for different CNNs. For instance, inspired by [10], I used the idea of using two CNNs together. The first was a normal CNN that was taking the entire image as the input while the second CNN tried to integrate the global and the local view, while simultaneously accounting for style information.

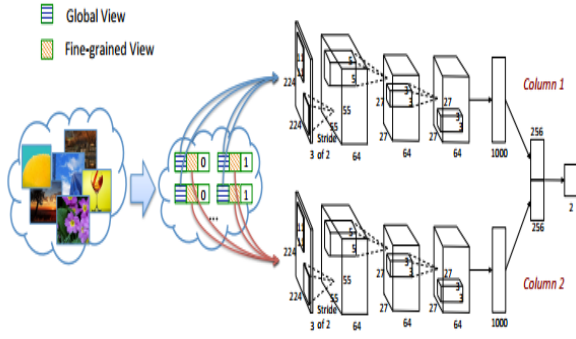


Figure 6: RAPIDNet [9]

5. Results

5.1. Baseline Results

In [9], Lu et al. achieved 64% baseline accuracy. [11], Murray et al. were able to achieve a 67% accuracy without using CNNs. The following were my baseline results (* means that in addition to image values, custom features were also used):

SVR	NN	SVM*	NN*	[9]	[11]
66.7%	67.6%	68.5%	69.01%	64%	67 %

Table 2: Comparing baseline accuracies for predicting image aesthetics

Therefore, I decided to choose the highest value, 69%, as the baseline for further experiments.

5.2. CaffeNet experiment

The first experiment was training the CaffeNet with randomly cropped images. After training it for 100k epochs, I got 74.3% accuracy. After reading [14], I decided to try warping the images. The intuition was that CNNs initialized with ImageNet weights, which are trained on down sampled whole images, do not inherently work well with non-down sampled random crops [3]. I was able to achieve 77.02% accuracy. This was counter intuitive, since I was initially assuming that random crops would perform better since the power of CNN comes from the local features.

My results		Results from [14]	
Random Crops	Warps	Random Crops	Warps
74.3%	77.02%	74.1%	76.82%

Table 3: Comparing CaffeNet performance

My results were better than previous studies. This could be attributed to having more data, more computing power and more number of epochs. Also, my results were consistent with previous studies which concluded that warps are better than random crops for CaffeNet because CaffeNet, originally was trained on downsampled whole trained.

5.3. VGGNet experiment

After training CaffeNet, I moved to VGG16 in hope to see more improvements. Here is how my results compare with the experiments done in [14] and [9]:

My VGGNet	[14] VGGNet	[9] ConvNet
78.6%	77.07%	74.46%

Table 4: Comparison of VGGNet experiments

This shows that the fine-tuned version of VGGNet was a little better than [14] and much better than [9]. I believe that this increase in accuracy happened because of training for more epochs while keeping low learning_rate.

5.4. Results of Custom CNN with deeper channels

This was inspired by [3] and I implemented this by using Python and numPy. It took longer than CaffeNet and VGGNet to train and the model was able to achieve 77.8% accuracy.

[3] ConvNet1	[3] ConvNet2	ConvNet w/ deeper channels
70.6%	75.2%	77.8%

Table 5: Comparison of CNN with deeper channels with others

5.5. Comparing all one-column CNNs tried so far

After training CaffeNet, VGGNet and custom CNN, I compared my experiment results with previous studies [14] [9] [11] as a midpoint check. I see that overall my model achieves better results than past studies. This could be attributed to two main reasons: I initialized the CNNs with ImageNet weights whereas past studies had trained things from scratch. Also, my CNN models were deeper which may have allowed the network to have more expressive power.

When experimenting with ReLU and PReLU, it turned out that since the original CNN weights had used ReLU during training, replacing with PReLU led to a lower performance. Also, while doing this, I had to add another parameter for PReLU and this lead to the exploding loss

phenomena in some cases. After some failed attempts, the learning rate was set to 10^{-6} to make sure that the weights did not blow up.

5.6. Weight Visualizations

Here are the weight visualizations of the convolutional network with deeper channels.

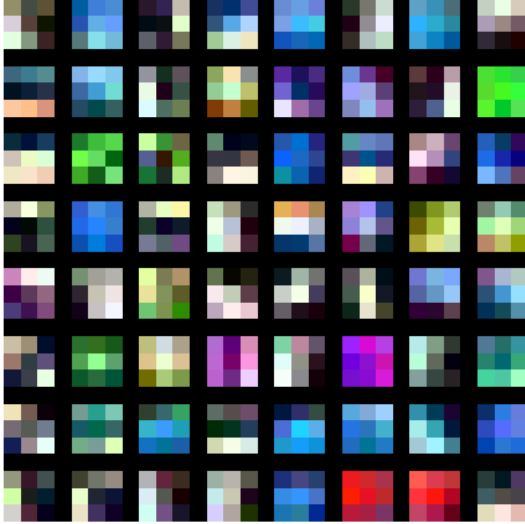


Figure 7: weight visualizations of trained CNN with deeper channels: 1st convolutional layer

In the first layer, most of them are lower level features that capture a lot of the edges and color information. As seen in the figure, there is a lot of blue and green and some red. After seeing these visualizations, I revisited the data to confirm if AVA indeed had a lot more blue and green colors. It turns out that most of the highly aesthetic images of landscape either contain an ocean, sky and/or a tree in the background. This explains why the first layer had abundance of those colors.

When I looked at the weights of the fully connected layer, most of them captured the higher level features. However, as seen in the figure above, the particular set of weights looked most visually interesting. The neuron in the 1st row 1st column and 2nd row, 2nd col (zero indexed) were activated for most of *Caravaggio's work*, which are most famous for *chiaroscuro*: contrasting treatment given to light and shade in drawings and paintings. This shows that the network may be able to hierarchically combine features into more abstract and broad concepts like shades.

Therefore, after doing these experiments, it was confirmed that pre-trained and fine-tuned CNNs could indeed be used for predicting picture aesthetics.

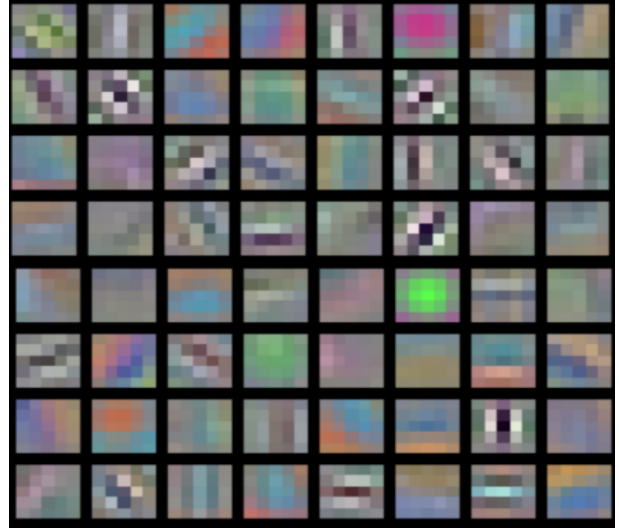


Figure 8: Weight Visualizations of FC-1024

5.7. Results of using Siamese Networks

Inspired by [9], I decided that the following would be my baseline while trying Siamese networks. Caffe has a

δ	[24]	SCNN	AVG_SCNN	DCNN	RDCNN
0	66.7%	71.20%	69.91%	73.25%	74.46%
1	67%	68.63%	71.26%	73.05%	73.70%

Figure 9: RAPIDNet results from [9]

Siamese model, [1], which was used to replicate RAPIDNet. Like [3], I decided to use a L-2 loss for the last layer and train the model. I was able to achieve 71.45% accuracy which is lower than the original paper. I tried playing with the parameters but was not able to achieve a much better result.

5.8. Choosing the best network to classify popular art forms from Renaissance Era

6. Classification and mis-classifications

6.1. From the AVADataset

Most of the image that correctly labeled had very explicit attributes that allowed the model to assign the correct label. For instance, most of the HIGH aesthetics images had good contrast, great placement of objects and overall sense of aesthetics. Similarly, most of the images that were correctly labeled as LOW had attributes, like most of them were dull and had random colors, spots and no homogeneity, that made it easy for the models to classify them correctly.

I noticed an interesting trend with the images that were mis-classified. First, most of them, from a broader perspec-

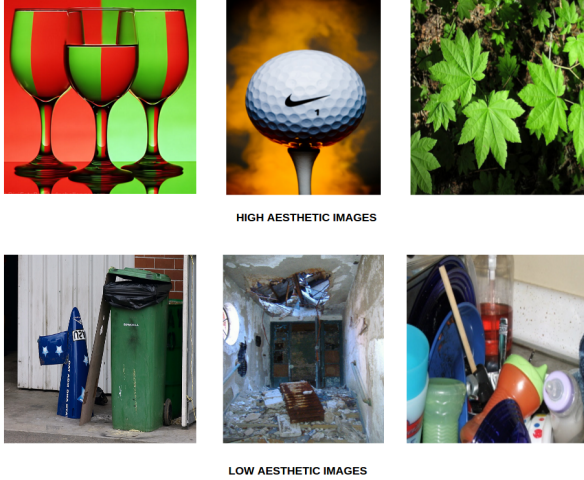


Figure 10: Correctly classified examples from AVA dataset



Figure 11: Incorrectly classified examples from AVA dataset

tive, looked aesthetically appealing, but if one were to pay attention to the detail, one could point of elements of the image that was responsible for the label. For instance, in the above figure, the image on the top left, looks relatively aesthetic if one looks at it. There are trees in the background and some part of the sky is visible. However, if one were to pay a closer attention, one would notice that the flag hanging from the tree is covered with mud, which is the reason why, even though the model thinks that this has HIGH aesthetics, it ought to be classified as LOW (in accordance to the comments on DPchallenge). Also, in some cases like the middle column, it was really hard to come to a decision on whether a particular image should be classified as having HIGH or LOW aesthetics. In such cases, I tried using different δ values to get some performance boost.

Finally, it is worthy to note that hand sketches, most of which had HIGH ratings, were mis-classified as images with LOW aesthetics.

6.2. Images from the Renaissance Era

After training and fine-tuning the model on the AVA dataset, I tried to test my model with renaissance art images. The idea was to test if the model was able to predict the correct label of these world famous pieces of art. Here are the results: Most of the portraits, specially with the dark



Figure 12: Best model classifying renaissance images

background were classified correctly by the model. For instance, works of *Caravaggio*, who is world renowned for this mastery of experimenting with light and shadows, were correctly classified. Also, works of *Leonardo Da Vinci*, were classified correctly to a great degree.

However, for most of the images of artworks that were sketches or with really light colored background, the model had real difficulty in classifying correctly. As seen in the above figure, one of the greatest works, *The Renaissance Man*, was mis-classified as LOW. After seeing this, I proposed that this was due to the speckles that were part of the original artwork. I found a replica of the same image, which did not have the speckles, and the model was able to classify it correctly. I concluded that some of the works, were mislabeled because of the granular nature of the art.

7. Failed experiments and things learned

7.1. During the Training Process

Initially, while playing with different architectures, I found the accuracy of my models to be low. After digging deeper, I realized that I had to resolve a potential issue with my loss function. Since, I had modeled the original regression problem as a classification problem, I had to introduce

a L-2 loss function which accounted for mis-classification appropriately. For instance, if my model predicted an image as a 8 when the true score was a 7, the penalty ought to be smaller than when the true score was a 7 and the predicted score was a 5. Therefore, I ended up using L-2 loss:

$$\frac{\sum_{n=1}^n (x_i - \bar{x})^2}{2n}$$

and got better results.

After training a bunch of models in Caffe and getting low accuracy, I realized that monitoring the effect of learning rate was quintessential to getting good accuracies, especially in the case of VGGNet whose weights would easily explode with high learning rates.

After getting a good model, I spent some time looking at what were the areas where the model was not performing well. The following figure shows different types of mis-classifications: For the AVA dataset, the model had most

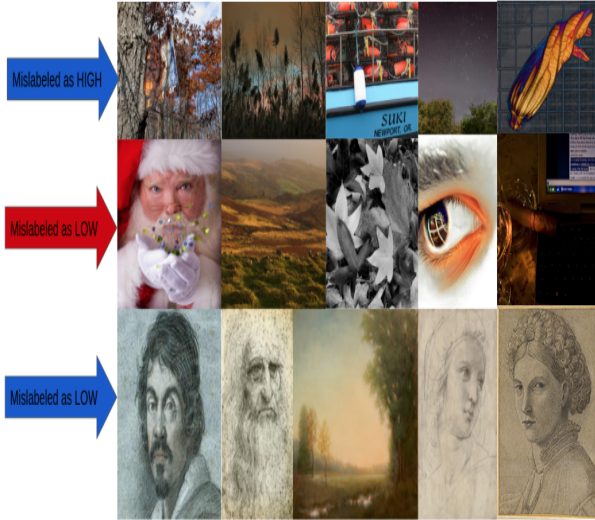


Figure 13: Mis-classifications by the best model

issues with images that broadly looked aesthetically appealing but had certain strong components that made them less appealing. In those cases, the model was not able to pick up those attributes and ended up mis-classifying them. Also, some of the images, like the one with Santa Claus, that had inherent graininess in the depicted subject, led to mis-classification. Finally for most of the renaissance images that were sketches were classified incorrectly.

There are a couple of reasons why I think this happened. Most of the images that were rated as LOW in the AVA dataset, either were of construction site or had bad exposure. In both these cases, the associated images were really grainy and in some of them, one could see dust like particles. I believe that the model learned to associate the graininess with less aesthetic beauty and thus mis-classified these

images.

Also, AVA dataset did not have any sketch images. Thus, I hypothesized when the model encountered all the renaissance era sketches, it mis-classified them because of the inherent graininess associated with sketching. In order to test out this, I found a clear image of the *Renaissance Man* and fed it to the network and model gave it a higher score.

8. Conclusion

In this paper, I experimented with predefined Convolutional Network architectures: CaffeNet, VGGNet and Siamese Net and showed that CNNs can be successfully trained to predict the aesthetics of an image. Then, by augmenting additional information into extra channel of images, I was able implement a deeper channel architecture which was able to perform better than the existing state of the art models trained on AVA dataset. I compared my best model with other existing models and finally, I evaluated it on renaissance art to measure how well it generalized on predicting image aesthetics for these world renowned pieces.

9. Future Work

From the experiments covered in this and other related papers (see references), there should be little doubt in CNNs potential to be used for the task of rating images based on their aesthetics. This is an interesting problem domain where one can spend more time to achieve better results. I would love to look further into different loss functions, better data augmentation and regularized and transfer learning to improve accuracy.

While starting to work on this problem, I spent some time looking for appropriate data sets that could be used for training and fine tuning. However, I could only find AVA that fit this criteria. I believe that, with a more specialized dataset, for instance specific to art in Renaissance Era, I could build better more specific models and see the features that make those world class paintings so amazing.

Finally, while testing the model, I tried gathering human data on how the model was doing. In the limited time I had, I created a small user interface which showed all the images, along with their bucket (score) labels, that were classified as HIGH by the model. I noticed, though the model captures the broad notion of art, it could do even better if I were able to incorporate people's personal preferences. Given more time, I would like to explore this domain. I hope that using similar trained models, one could use a computer to evaluate the aesthetics of any image and provide editing recommendations for artists, designers and photographers.

References

- [1] C. community. Siamese Networks in caffe.
- [2] C. cumminity. CaffeNet and VGGnet in caffe.
- [3] E. Gong. Estimating photo aesthetic rating using convolutional neural networks. 2015.
- [4] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. ; the WEKA data mining software: An update; SIGKDD explorations. 11, 2009.
- [5] M.-M. C. N. J. M. X. H. P. H. S. T. S.-M. Hu. Global contrast based salient region detection. 2011.
- [6] Y. Jia. Caffe: An open source convolutional architecture for fast feature embedding. 2013.
- [7] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [8] Y. Ke, X. Tang, and F. Jing. The design of high-level features for photo quality assessment. In *in IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2006, vol.1*, page 419426. 2006.
- [9] X. Lu, Z. Lin, H. Jin, J. Yang, and J. Z. Wang. Rapid: Rating pictorial aesthetics using deep learning. In *Proceedings of the ACM International Conference on Multimedia, MM 14*, page 457466. 2014.
- [10] X. Lu, Z. Lin, H. Jin, J. Yang, and J. Z. Wang. Rapid: Rating pictorial aesthetics using deep learning. In *Proceedings of the ACM International Conference on Multimedia*, page 457466. 2014.
- [11] N. Murray, L. Marchesotti, and F. Perronnin. Ava: A largescale database for aesthetic visual analysis. In *Computer Vision and Pattern Recognition (CVPR)*, page 24082415. June 2012.
- [12] V. O. S. Dhar and T. Berg. High level describable attributes for predicting aesthetics and interestingness. *Computer Vision and Pattern Recognition (CVPR)*, pages 1657–1664, 2011.
- [13] K. Simonyan and A. Zisserman. *Very deep convolutional networks for large-scale image recognition*. CoRR, abs/1409.1556, 2014.
- [14] P. Veerina. Learning good taste: Classifying aesthetic images. 2015.

10. Supplementary Material



Figure 14: Correctly classified HIGH aesthetic images



Figure 15: Correctly classified LOW aesthetic images



Figure 16: Correctly classified Renaissance images

10.1. Source Code

10.1.1 Check if AVA dataset links are valid

```
1
2
3
4 #
5 # While downloading AVA images, you may need to
6 # check the un-downloaded images.
7 # This script will display information of missing
8 # images in console.
9 #
10 # Note: This script should be placed under
11 #   AVA_dataset/script/ and images are saved at
12 #   AVA_dataset/image/ folder.
13
14 import os
15 import sys
16
17 savePath = r'../image'
18 AVAtxt = r'../AVA.txt'
19
20 # arg check
21 if len(sys.argv) < 3:
22     print 'arg failer! # python checkWholeness.py
23         beginIndex stopIndex'
24     exit()
25
26 # must >= 1
27 beginIndex = int(sys.argv[1])
28 # must >= 255530
29 stopIndex = int(sys.argv[2])
30
31 f = open(AVAtxt)
32 for line in f:
33     line = line.strip().split(' ')
34     imageIndex = line[0]
35     imageID = line[1]
36
37     # use begin and stop index constrain
38     if int(imageIndex) < beginIndex:
39         continue
40     elif int(imageIndex) >= stopIndex:
41         break
42
43     # if the image does not exist, display its '
44     # index ID' pair
45     if os.path.isfile(os.path.join(savePath,
46     imageIndex + '.jpg')) == False:
47         print imageIndex, imageID
```

10.1.2 Download AVA data-set images

```
1 #
2 # Use 'AVA.txt' provided by database author to
3 # download all images from 'dpchallenge.com'.
4 # This script should be placed under
5 #   AVA_dataset/script/ and images are saved at
6 #   AVA_dataset/image/ folder.
7 #
8 # Usage:
9 #   python downloadImage.py beginIndex stopIndex
10 # So you can download in multi-process.
11 #
```

```
11 # Note: few images were deleted from the website,
12 # please use checkWholeness.py to check
13 # missed images and manually delete their
14 # items from AVA.txt
15
16 import time
17 import os
18 import urllib
19 import sys
20 from HTMLParser import HTMLParser
21
22 ##
23 # Heritate from HTMLParser to parse main image's
24 # url
25
26 class dpchallengeImageParser(HTMLParser):
27     def __init__(self):
28         HTMLParser.__init__(self)
29         self.name = None
30
31     # fetch url and loop out when both image's
32     # width and height exceed 200 pix.
33     def handle_starttag(self, tag, attrs):
34         if self.name is not None:
35             return
36         if tag == 'img':
37             tmpWidth = 0
38             tmpHeight = 0
39             for key, value in attrs:
40                 if key == 'src':
41                     tmpName = value
42                     tmpWidth = 0
43                     tmpHeight = 0
44                 elif key == 'width':
45                     tmpWidth = float(value)
46                 elif key == 'height':
47                     tmpHeight = float(value)
48
49     # extract main (big) image only
50     if (tmpWidth > 150) and (tmpHeight > 150):
51         self.name = tmpName
52         break
53
54 savePath = r'../image'
55 URLprefix = r'http://www.dpchallenge.com/image.
56 php?IMAGE_ID='
57 AVAtxt = r'../AVA.txt'
58
59 # arg check
60 if len(sys.argv) < 3:
61     print 'arg failer! # python downloadImage.py
62         beginIndex stopIndex'
63     exit()
64
65 # must >= 1
66 beginIndex = int(sys.argv[1])
67 # must >= 255530
68 stopIndex = int(sys.argv[2])
69
70 f = open(AVAtxt)
71 for line in f:
72     line = line.strip().split(' ')
73     imageIndex = line[0]
74
75     # use begin and stop index constrain
76     if int(imageIndex) < beginIndex:
77         continue
```

```

72     elif int(imageIndex) >= stopIndex:
73         break
74
75         # skip exiting images
76         if os.path.isfile(os.path.join(savePath,
77                                         imageIndex + '.jpg')) == True:
78             continue;
79
80         # get display webpage url
81         imageID = line[1]
82         URL = URLprefix + imageID
83
84         # url request and write image
85         tic = time.time()
86         # get 'dpchallenge' display page
87         urlopen=urllib.URLopener()
88         fp = urlopen.open(URL)
89         data = fp.read()
90         fp.close()
91         # parse out image url in 'self.name'
92         urlParser = dpchallengeImageParser()
93         # ignore illegal characters (almost
94         # impossible appear in image's url)
95         data = data.decode('ascii', 'ignore')
96         urlParser.feed(data)
97
98         # some main images are even smaller than
99         # advertisements
100        # or have been deleted from website
101        # they need manual downloading
102        if urlParser.name is None:
103            continue
104
105        # get iamge
106        fimg = urlopen.open(urlParser.name)
107        data = fimg.read()
108        fimg.close()
109        # write image
110        fout = open(os.path.join(savePath, imageIndex
111                                  + '.jpg'), 'wb')
112        fout.write(data)
113        fout.close()
114        toc = time.time()
115
116        # display progress
117        print 'Processing image: {} \ttime elapse: {} \
118              tstop at: {}'.format(int(imageIndex), toc -
119                                    tic, stopIndex)

```

10.1.3 transform original AVA dataset to make it work with Caffe

```

1 import numpy as np
2
3 def get_avg_ratings(ratings):
4     ratings = [float(x) for x in ratings]
5     assert len(ratings) == 10, "something fishy yo!"
6
7     count = 0.0
8     for indx, r in enumerate(ratings):
9         count += (indx + 1) * float(r)
10
11     return count / np.sum(ratings)
12
13 def generate_AVA_transformed():
14     delta = 5.0 # threshold value for
15                 classification

```

```

14 f = open("AVA_dataset/AVA.67140.txt")
15 classification = open("AVA_dataset/
16                       AVA_classification.txt", 'w')
17 regression = open("AVA_dataset/AVA_regression.
18                   txt", 'w')
19 for i in f:
20     arr = i.split(" ")
21     file_name = "%s.jpg" % arr[0]
22     ratings = arr[2:12]
23     avg_rating = get_avg_ratings(ratings)
24     classification_rating = 1 if avg_rating >
25                             delta else 0
26     regression.write("%s %s\n" %
27                     (file_name, avg_rating))
28     classification.write("%s %s\n" %
29                          (file_name, classification_rating))
30
31 if __name__ == "__main__":
32     print "Your task has started..."
33     generate_AVA_transformed()
34     print "Your task is complete!"

```

10.1.4 Transform images into Caffe supported LMDB format

```

1 #!/usr/bin/env sh
2 # Create the imagenet lmdb inputs
3 # N.B. set the path to the imagenet train + val
4   data dirs
5
6 EXAMPLE=/mnt/AVA_dataset
7 DATA=/mnt/AVA_dataset
8 TOOLS=build/tools
9
10 TRAIN_DATA_ROOT=/mnt/AVA_dataset/image/
11 VAL_DATA_ROOT=/mnt/AVA_dataset/image/
12
13 # Set RESIZE=true to resize the images to 256x256
14   . Leave as false if images have
15   # already been resized using another tool.
16 RESIZE=true
17 if $RESIZE; then
18     RESIZE_HEIGHT=256
19     RESIZE_WIDTH=256
20 else
21     RESIZE_HEIGHT=0
22     RESIZE_WIDTH=0
23 fi
24
25 if [ ! -d "$TRAIN_DATA_ROOT" ]; then
26     echo "Error: TRAIN_DATA_ROOT is not a path to a
27     directory: $TRAIN_DATA_ROOT"
28     echo "Set the TRAIN_DATA_ROOT variable in
29     create_imagenet.sh to the path" \
30     "where the ImageNet training data is
31     stored."
32     exit 1
33 fi
34
35 if [ ! -d "$VAL_DATA_ROOT" ]; then
36     echo "Error: VAL_DATA_ROOT is not a path to a
37     directory: $VAL_DATA_ROOT"
38     echo "Set the VAL_DATA_ROOT variable in
39     create_imagenet.sh to the path" \
40     "where the ImageNet validation data is
41     stored."

```



```

34     exit 1
35 fi
36
37 echo "Creating train lmdb..."
38
39 GLOG_logtostderr=1 $TOOLS/convert_imageset \
40     --resize_height=$RESIZE_HEIGHT \
41     --resize_width=$RESIZE_WIDTH \
42     $TRAIN_DATA_ROOT \
43     $DATA/train.txt \
44     $EXAMPLE/AVANet_train_lmdb
45
46 echo "Creating val lmdb..."
47
48 GLOG_logtostderr=1 $TOOLS/convert_imageset \
49     --resize_height=$RESIZE_HEIGHT \
50     --resize_width=$RESIZE_WIDTH \
51     $VAL_DATA_ROOT \
52     $DATA/val.txt \
53     $EXAMPLE/AVANet_val_lmdb
54
55 echo "Done."

```

10.1.5 Custom convolutional layer for images with Deeper channels

```

1 # NOTE: the following file could be written by
  # using loops, thus decreasing the number of
  # lines of code.
2 # however the author decided to write this in
  # the unrolled fashion for debugging ease
3
4 import numpy as np
5
6 from cs231n.layers import *
7 from cs231n.fast_layers import *
8 from cs231n.layer_utils import *
9
10
11 class DeepChannelNet(object):
12     # conv - relu - 2x2 max pool
13     # conv - relu - 2x2 max pool
14     # conv - relu - 2x2 max pool
15     # conv - relu - 2x2 max pool
16     # conv - relu - 2x2 max pool
17     # conv - relu - 2x2 max pool
18     # conv - relu - 2x2 max pool
19     # - affine - relu
20     # - affine
21     # - softmax
22
23
24     def __init__(self, input_dim=(3, 32, 32),
25                 num_filters=32, filter_size=7,
26                 hidden_dim=100, num_classes=10,
27                 weight_scale=1e-3, reg=0.0,
28                 dtype=np.float32):
29         self.params = {}
30         self.reg = reg
31         self.dtype = dtype
32         # extracting various parameters from the
33         # input dimensions
34         C, H, W = input_dim
35         print "C: %s H: %s W: %s" % input_dim
36
37 # all my convolutional layers
38 self.params['W1'] = weight_scale * np.random.

```

```

39     randn(num_filters, C, filter_size,
40           filter_size)
41     self.params['b1'] = np.zeros(num_filters)
42
43     self.params['W2'] = weight_scale * np.random.
44     randn(num_filters, num_filters, filter_size,
45           filter_size)
46     self.params['b2'] = np.zeros(num_filters)
47
48     self.params['W3'] = weight_scale * np.random.
49     randn(num_filters, num_filters, filter_size,
50           filter_size)
51     self.params['b3'] = np.zeros(num_filters)
52
53     self.params['W4'] = weight_scale * np.random.
54     randn(num_filters, num_filters, filter_size,
55           filter_size)
56     self.params['b4'] = np.zeros(num_filters)
57
58     self.params['W5'] = weight_scale * np.random.
59     randn(num_filters, num_filters, filter_size,
60           filter_size)
61     self.params['b5'] = np.zeros(num_filters)
62
63     self.params['W6'] = weight_scale * np.random.
64     randn(num_filters, num_filters, filter_size,
65           filter_size)
66     self.params['b6'] = np.zeros(num_filters)
67
68     self.params['W7'] = weight_scale * np.random.
69     randn(num_filters, num_filters, filter_size,
70           filter_size)
71     self.params['b7'] = np.zeros(num_filters)
72
73 # end of the convolutional era
74 self.params['W8'] = weight_scale * np.random.
75     randn(H * W * num_filters / pow(4, 2), H * W *
76           num_filters / pow(4, 3))
77     self.params['b8'] = np.zeros(H * W *
78           num_filters / pow(4, 3))
79
80 # for the output affine layer
81 self.params['W9'] = weight_scale * np.random.
82     randn(H * W * num_filters / pow(4, 3),
83           num_classes)
84     self.params['b9'] = np.zeros(num_classes)
85
86 for k, v in self.params.iteritems():
87     self.params[k] = v.astype(dtype)
88
89 def loss(self, X, y=None):
90     W1, b1 = self.params['W1'], self.params['b1']
91     W2, b2 = self.params['W2'], self.params['b2']
92     W3, b3 = self.params['W3'], self.params['b3']
93     W4, b4 = self.params['W4'], self.params['b4']
94     W5, b5 = self.params['W5'], self.params['b5']
95     W6, b6 = self.params['W6'], self.params['b6']
96     W7, b7 = self.params['W7'], self.params['b7']
97     W8, b8 = self.params['W8'], self.params['b8']
98     W9, b9 = self.params['W9'], self.params['b9']
99
100 # pass conv_param to the forward pass for the
101 # convolutional layer
102 filter_size = W1.shape[2]
103 conv_param = {'stride': 1, 'pad': (
104     filter_size - 1) / 2}

```

```

82
83 # pass pool_param to the forward pass for the
84   max-pooling layer
85 pool_param = {'pool_height': 2, 'pool_width':
86   2, 'stride': 2}
87
88 scores = None
89 # print "statistics for X, W1", X.shape, W1.
90   shape
91 a1, c1 = conv_relu_pool_forward(X, W1, b1,
92   conv_param, pool_param)
93 a2, c2 = conv_relu_pool_forward(a1, W2, b2,
94   conv_param, pool_param)
95 a3, c3 = conv_relu_pool_forward(a2, W3, b3,
96   conv_param, pool_param)
97 a4, c4 = conv_relu_pool_forward(a3, W4, b4,
98   conv_param, pool_param)
99 a5, c5 = conv_relu_pool_forward(a4, W5, b5,
100   conv_param, pool_param)
101 a6, c6 = conv_relu_pool_forward(a5, W6, b6,
102   conv_param)
103 a7, c7 = conv_relu_pool_forward(a6, W7, b7,
104   conv_param)
105
106 a8, c8 = affine_relu_forward(a7, W8, b8)
107 scores, c9 = affine_forward(a8, W9, b9)
108
109 # for i, a in enumerate([a1, a2, a3, a4, a5,
110   a6, a7, a8]):
111   # print "A%s shape is %s" % ((i+1), a.
112     shape)
113
114   # print "W8 shape and bias", W8.shape, b8.
115     shape
116   # print "W9 shape and bias", W9.shape, b9.
117     shape
118   if y is None:
119     return scores
120
121 loss, grads = 0, {}
122 dloss, dscores = softmax_loss(scores, y)
123
124 da8, dw9, db9 = affine_backward(dscores, c9)
125 da7, dw8, db8 = affine_relu_backward(da8, c8)
126 da6, dw7, db7 = conv_relu_backward(da7, c7)
127 da5, dw6, db6 = conv_relu_backward(da6, c6)
128 da4, dw5, db5 = conv_relu_pool_backward(da5,
129   c5)
130 da3, dw4, db4 = conv_relu_pool_backward(da4,
131   c4)
132 da2, dw3, db3 = conv_relu_pool_backward(da3,
133   c3)
134 da1, dw2, db2 = conv_relu_pool_backward(da2,
135   c2)
136 dX, dw1, db1 = conv_relu_pool_backward(da1,
137   c1)
138
139 grads['W1'] = dw1
140 grads['W2'] = dw2
141 grads['W3'] = dw3
142 grads['W4'] = dw4
143 grads['W5'] = dw5
144 grads['W6'] = dw6
145 grads['W7'] = dw7
146 grads['W8'] = dw8
147 grads['W9'] = dw9

```

```

130 grads['b1'] = db1
131 grads['b2'] = db2
132 grads['b3'] = db3
133 grads['b4'] = db4
134 grads['b5'] = db5
135 grads['b6'] = db6
136 grads['b7'] = db7
137 grads['b8'] = db8
138 grads['b9'] = db9
139
140 reg_loss = 0.0
141 for i in range(1,10):
142     layer_name = str(i)
143     # print "working on layer %s" %
144     layer_name
145     reg_loss += 0.5 * self.reg * np.sum(self.
146     params['W' + layer_name] * self.params['W' +
147     layer_name])
148     grads['W' + layer_name] += self.reg *
149     self.params['W' + layer_name]
150
151 loss = dloss + reg_loss
152 return loss, grads

```

10.1.6 CaffeNet, VGGNet, Siamese Network implementation

For training CaffeNet, VGGNet and Siamese network, Caffe[6], along with native scripts available in the model zoo[2] and [1] were used.