

Blind Image Quality Assessment using Semi-supervised Rectifier Networks

Huixuan Tang
University of Toronto
hxtang@dgp.toronto.edu

Neel Joshi
Microsoft Research
neel@microsoft.com

Ashish Kapoor
Microsoft Research
akapoor@microsoft.com

Abstract

It is often desirable to evaluate images quality with a perceptually relevant measure that does not require a reference image. Recent approaches to this problem use human provided quality scores with machine learning to learn a measure. The biggest hurdles to these efforts are: 1) the difficulty of generalizing across diverse types of distortions and 2) collecting the enormity of human scored training data that is needed to learn the measure. We present a new blind image quality measure that addresses these difficulties by learning a robust, nonlinear kernel regression function using a rectifier neural network. The method is pre-trained with unlabeled data and fine-tuned with labeled data. It generalizes across a large set of images and distortion types without the need for a large amount of labeled data. We evaluate our approach on two benchmark datasets and show that it not only outperforms the current state of the art in blind image quality estimation, but also outperforms the state of the art in non-blind measures. Furthermore, we show that our semi-supervised approach is robust to using varying amounts of labeled data.

1. Introduction

A simple, scalar measure of perceptual image quality can inform a number of algorithms in computer vision and computer graphics. Such a measure can be used for evaluating image processing methods, driving image compression techniques, or filtering out low quality images before image recognitions tasks. Computing perceptual image quality is challenging due to variations in image content and the underlying image distortion process.

Blind measures of image quality, i.e., those that do not require groundtruth reference images, are challenging to create but are much more desirable than those that require a reference image. Recent approaches have used labeled data and machine learning to model perceptual image quality. Such methods first extract hand-crafted features from images and then learn a mapping of these features to subjective quality scores by kernel or nearest neighbor regression.

For the regression model to work well, the kernel function needs to be highly informative with respect to the image distortions alone and not be affected by other aspects such as the image content. Previous methods often define this

function in the original or a linearly compressed space (e.g. PCA) of raw features and therefore critically rely on the design of the features to de-correlate image distortion from image content. In practice, however, the features often capture a combination of similarity in distortion and image content, thus making discrimination on distortion alone quite challenging. Though previous methods show good performance on datasets of small numbers of distortion types [13], they usually perform poorly on datasets with more distortion types [10], as images with different distortions overlap in the feature space in a way that is not separable with the linear models used in previous work.

The performance of previous methods also degrades significantly when labels are sparse, as is the case with many machine learning methods. Using more training data helps, but the collection of large datasets for image quality assessment (IQA) is non-trivial and expensive [13], as the degraded images need to be collected across a wide range of image quality, content, and type of degradations, thus resulting in millions of trials. Furthermore, each degraded image requires the evaluation of many subjects to eliminate bias across subjects and content. Meanwhile, all experiments need to be conducted under a controlled environment and therefore cannot easily be migrated to a crowd-sourcing platform. As a result, the current standard blind image quality assessment datasets [13, 10] are all generated from a relatively small number of natural images.

In this paper, we present a neural network approach to alleviate these problems. Specifically we define the kernel function as a simple radial basis function on the output of a deep belief network of rectified linear hidden units [9, 2]. We first pre-train the rectifier networks in an unsupervised manner and then fine-tune them with labeled data. Finally we model the quality of images with Gaussian Process regression. Our approach outperforms both blind and non-blind methods (Fig. 1) – this is the first *blind* measure we are aware of that outperforms *non-blind* measures on perceptual image quality datasets.

The key advantage of our model is that its success mostly relies on an *unsupervised* pre-training stage. Unlabeled data is easy to generate en masse and thus our approach benefits from using a large amount of data that much more densely samples the space of distortion types and image content. This enhances generalization power in both di-

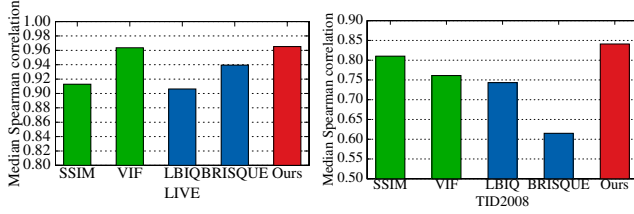


Figure 1: Best performance of our method and state-of-arts IQA measures on the LIVE and TID2008 dataset. All methods using 20% of the labeled data set for training and learning based method use the remaining 80% for testing. Higher Spearman correlation is better.

mensions, and we show that our resulting measure performs well even with fairly small labeled training sets across many distortion types, which is a significant advantage over previous work. Specifically, our contributions include: 1) a state of the art image quality assessment method that significantly outperforms existing measures on a wide range of distortion types, and 2) a regression method that is robust to using a small amount of labeled training data.

2. Related work

Blind image quality assessment Much effort in blind image quality assessment has been devoted to crafting features to de-correlate distortion information from image content [14, 6]. Yet the learning models of these approaches are strikingly simple. As a result, there is usually an image content dependent bias among distortion types. Moorthy and Bovik [8] proposed to address the problem by identifying the distortion types first, but the classification problem becomes harder with many distortion types. Tang et al. [14] suggests these biases can be eliminated to a great extent by correlating quality scores with thousands of features. However, this measure still exhibits some bias on image content. While previous work has focused mostly on the features used for image quality assessment, we focus on using more sophisticated learning methods and overcoming the burdens associated with training data.

The difficulty of acquiring many labels is increasingly drawing the attention of the research community. One way to overcome this is by supervised learning. Mittal et al. [7] train a Gaussian model of carefully designed features. However, the simplicity of the Gaussian model makes it prone to fail as the number of distortion types increase. Xue et al. [16] train a regression model, but replace the human scores with a non-blind quality measure. Our use of unlabeled images in the pre-training stage is similar in spirit, but our goal is to *improve* the performance when labels are limited rather than design a method that uses no human scores at the cost of *worse* performance. Ye et al. [17] also uses unlabeled distorted images but only use them to generate

distortion-specific codewords and not to learn the mapping to image quality. Therefore, their method does not address the case where labeled image quality data is sparse.

Deep neural networks Restricted Boltzmann machines, deep belief networks and their variations are proven to be compact universal approximators [4, 12] and achieve impressive performance in various field of applications as a way to model, visualize, and infer complex nonlinear data. A very important finding we exploit in this work is that the training of such neural networks can be decomposed into unsupervised pre-training and a supervised refinement stage [3]. Pre-training builds a robust probabilistic model for the input data and usually achieves great performance because the model has good generalization power. This property motivated us to perform semi-supervised learning [18] using such models, so that we can exploit unlabeled images to deal with the limited availability of human scores for image quality assessment.

We are specifically inspired by the work of Salakhutdinov and Hinton [11], which uses a binary deep belief network to formulate a Gaussian process kernel and recognizes face orientations and digits from labeled data. Our learning model extends this work to semi-supervised rectifier networks, and we apply this technique to overcome the challenges of blind image quality assessment.

3. Overview

The goal of the proposed framework is to provide a measure of image quality from the relevant features extracted from images. Specifically, for the purpose of this paper we extracted the same set of image features as the LBIQ measure [14]. These features include univariate and cross-scale histograms and statistics of complex wavelet transform of images (the real part, absolute value, and phase) as well as a few direct blur and noise measures. Finally, the input data is whitened via a discrete cosine transform before it is used in our proposed system.

We choose the LBIQ features over other features as LBIQ has the best average performance across both the LIVE and TID datasets, as seen in Fig 1. Particularly, LBIQ performs well in two difficult scenarios we are specifically interested in: 1) when label are sparse (Fig. 8) 2) when there are many distortion types (Fig. 9).

Overall, our model is a multi-layer network that learns a regression function from images to a single scalar quality score for each image, as in [13, 10]. Fig. 2 shows the configuration of our model. There are two specific components of the model: the first component is a Gaussian Process that regresses the final quality score given activations from a trained neural network. The second component is a neural network whose goal is to provide a feature representation that is informative for image quality assessment.

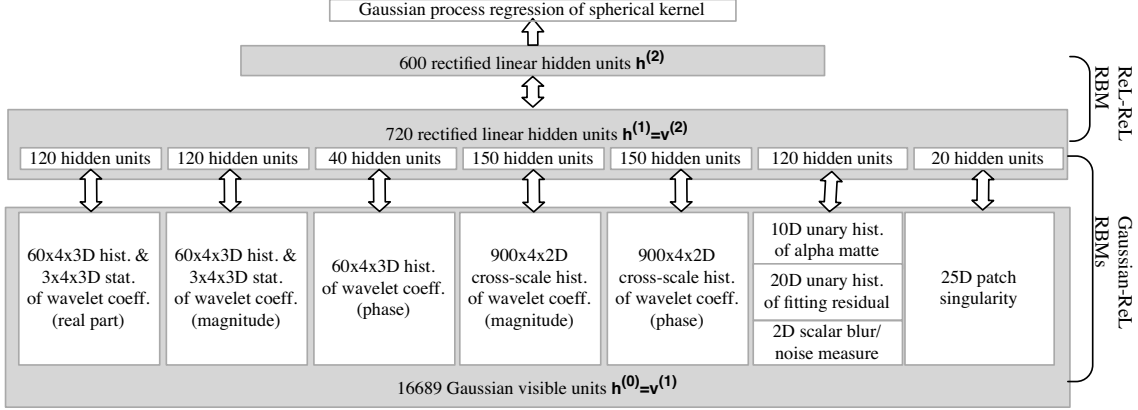


Figure 2: Configuration of our model. The number of hidden units are determined by cross-validation.

The neural network is a deep belief network (DBN) [4] of three layers. The bottom most layer (Fig. 2) is the input layer with Gaussian visible units (i.e. activation is linear to input $\mathcal{A}(z) = z$) and connects to the second layer of rectified linear hidden units (ReLU, i.e. activation rectifies negative input to zero $\mathcal{A}(z) = \max(z, 0)$). Note that these two layers are comprised of restricted Boltzmann machines (RBMs) that encodes the input data with the hidden layer activities. We propose to use seven different groupings of features and train a separate RBM for each of the features. This enables us to reduce not only the number of free parameters but also the computational cost of training. The output of these RBMs are concatenated as the input to the next layer, which is just a single RBM with ReL units in both layers (720 inputs and 600 outputs for the case of Fig. 2). The detailed input and hidden dimensionality is indicated in Fig. 2.

Intuitively, the rectifier network that we propose performs a non-negative factorization of the input features. The ReL units we use introduce non-linearity to the model, and the hard threshold we use encourages zeros in the unit responses. We have also tried to train an alternative model of standard binary hidden units, but it needs several magnitudes more units to express the same data because the resulting model includes many units of shared weights and different bias offset. Since ReL units are an approximation of such an ensemble of binary units [9], we find using ReL units a natural choice.

Our parameter learning for the neural network is inline with that of Hinton et.al [3], where they propose training similar models with an greedy layer-wise unsupervised pre-training stage followed by a supervised fine-tuning stage. In the pre-training stage, we learn a generative model of the unlabeled input features. Consequently, the objective of pre-training is simply to adjust the parameters to maximize the likelihood of the unlabeled data. Eventually when the RBMs reach their equilibrium state, the hidden units that

are activated by the visible units can reproduce the visible data in turn. On the other hand, the fine-tuning stage further strengthens the discriminative power of the model by aligning the network output with the available labels of the training data. Specifically, in the fine-tuning stage, we adjust weights to maximize the Gaussian process likelihood on the labeled data set. We provide the details of these two stages of the training algorithm in Sec. 4 and Sec. 5.

One of the most critical parameters in the proposed model is the number of hidden units in each layer. We determined the number of hidden units in the RBMs by cross-validation. Specifically, we divided the data into a training and validation set whose reference images are disjoint. Then, for each RBM in the network, we first trained a series of models of different number of hidden units with the training set and then evaluated the reconstruction on the validation set. We choose the number of hidden units of Gaussian-ReL RBMs in the range of $\{5, 10, 20, 30, 40, 50, 60, 80, 100, 120, 150, 200\}$, and the number of hidden units of the top RBM with ReL visible units in the range of $\{100, 200, 300, \dots, 1000\}$. The model with the smallest reconstruction error was chosen as the optimal model.

The validation set is also used for monitoring the training of the neural network model. As the neural network may over-fit to the training set in both the pre-training and fine-tuning stage, we stop the learning process when we observe that the reconstruction error (in the pre-training stage) or the regression performance (in the fine-tuning stage) starts to decrease.

3.1. Using unlabeled data

As noted previously, image quality assessment methods generally suffer from lack of sufficient training data. However, it has been shown that learning in multi-layer models can benefit from using unlabeled data for structuring their mid-level representations [18].



Figure 3: Examples of reference images in the unlabeled dataset we simulate.

The key requirement of semi-supervised learning is that the simulated and labeled data set should have an identical and independent distributions. That is to say the distortion types and levels should be the same as the labeled data set, and the reference images should be about the same size as those in the labeled data set. It is straightforward to expand the data coverage from reference images under this assumption.

We crawl 80 high quality images from internet (as shown in Fig. 3) and simulate distortions of various levels and types on these images. This considerably expands the training data along the dimension of reference images. In comparison, the LIVE and TID2008 data sets generate distorted images from only 20 – 30 images. Note that a considerable portion of that data is used for performance evaluation and not available for training. So the actual labeled training set is at most about 20 reference images \times 5(LIVE) or 17(TID2008) distortion types \times 4 – 5 distortion levels.

Prior to simulating distortions, we first reduce the size of images to 512×384 pixels to match the resolution of the labeled data. In experiments for the LIVE dataset, our unlabeled data is limited to the five included distortion types. In experiments for the TID2008 dataset, we simulate 13 of 17 distortion types in the TID dataset. The 4 distortion types that we did not simulate are non eccentricity pattern noise, local block-wise intensity change, global intensity shift, and contrast change. We exclude these types from the model because the LBIQ features we use are invariant to these distortions and the inclusion of them hinders pre-training.¹ In total, this results in an unlabeled dataset of 5200 distorted images generated from 80 distortion-free images.

4. Pre-training without labels

In the pre-training stage, we learn a generative model of the features using the neural network representation. Such a generative perspective is feasible due to the fact that we can simply view the deep belief network as a stack of RBMs. The seven RBMs connecting the hidden input layer and the first hidden layer learn a generative model of each of the input image features. Similarly, the next hidden layer learns the joint distribution of all features. Note that we can continue to add more hidden layers and model even higher-level

¹For fairness of evaluation, we only ignore these types during pretraining and include all distortion types in the fine-tuning and evaluation stage.

Input: B random batches of training samples v_1, v_2, \dots, v_B ;

Output: model parameters $\Theta = \{\theta_1, \theta_2, \dots, \theta_K\}$;

Parameters: learning rate η , momentum τ ;

initialize Θ (see text for details);

for $t = 1 \dots T$ **do**

for $v_+ = v_1, \dots, v_B$ **do**

compute hidden unit mean $\mathbf{h}_+ = \max(\mathbf{W} \frac{\mathbf{v}_+}{\sigma} + \lambda, 0)$;

compute visible unit mean $\mathbf{v}_- = \sigma \mathbf{W}^\top \mathbf{h}_+ + \mathbf{b}$;

compute hidden unit mean $\mathbf{h}_- = \max(\mathbf{W} \frac{\mathbf{v}_-}{\sigma} + \lambda, 0)$;

compute CD-1 gradients $\Delta \theta_k^{(t)}$:

$$\Delta \mathbf{W} = \mathbf{CD} \left(\frac{\mathbf{v}_+^\top \mathbf{h}_-}{\sigma} \right), \quad (1)$$

$$\Delta \lambda = \mathbf{CD}(\mathbf{h}_-), \quad \Delta \mathbf{b} = \mathbf{CD} \left(\frac{(\mathbf{v}_- - \mathbf{b})}{\sigma^2} \right) \quad (2)$$

$$\Delta \sigma = \mathbf{CD} \left(\frac{(\mathbf{v}_- - \mathbf{b})^2}{\sigma^3} - \mathbf{h}_-^\top \mathbf{W} \frac{\mathbf{v}_-}{\sigma^2} \right) \quad (3)$$

compute ADAGRAD learning rate $\gamma_k = \eta / \sqrt{\sum_t \Delta \theta_k^{(t)}}$;

add momentum $\hat{\Delta \theta}_k^{(t)} = \Delta \theta_k^{(t)} + \tau \Delta \theta_k^{(t-1)}$;

adjust model $\theta_k = \theta_k + \hat{\Delta \theta}_k^{(t)} \gamma_k, k = 1 \dots K$;

Figure 4: Pretrain RBMs of linear visible and ReL hidden units.

representation of features. However, we limited ourselves to two hidden layers as our experiments (see Sec. 6.1) did not show an advantage for a deeper model architecture with the specific features we used.

We pre-train the model greedily in a layer-by-layer manner. Specifically, we first train the seven Gaussian-ReL RBMs in the bottom layer and then the ReL-ReL RBM in the top-layer by maximizing the likelihood of the data by stochastic gradient descendant.² The top level RBM of rectified input is parameterized by network weights \mathbf{W} , hidden layer bias λ , and visible layer bias \mathbf{b} . The Gaussian-rectified RBMs are parameterized by network weights and biases \mathbf{W} , λ , and b as well as visible unit variance σ .

Fig. 4 and 5 outlines the steps to learn these parameters given a number of samples for the input layer v_+ . The high-level idea is to optimize the model parameters. The core part is to compute the gradient of likelihood by 1-step contrastive divergence (CD-1), which uses Gibbs sampling to approximate the intractable true gradient. Specifically, we simulate the model driven by the input data v for 1.5 cycles and collect mean activation of visible and hidden units h_+ , h_- , and v_- for the first and last upward half-cycle, and the gradients are computed from the difference in statistics between two cycles

$$\mathbf{CD}(z) = \mathbf{E}_+(z) - \mathbf{E}_-(z). \quad (6)$$

²In fact, the likelihood of RBMs of ReL unit is not well defined, but the contrastive divergence gradients of such RBMs are well defined by interpreting the ReL units as sum of binary units of shared weights – see [9] for details.

Input: B random batches of training samples v_1, v_2, \dots, v_B ;

Output: model parameters $\Theta = \{\theta_1, \theta_2, \dots, \theta_K\}$;

Parameters: learning rate η , momentum τ ;

initialize Θ (see text for details);

for $t = 1 \dots T$ **do**

for $\mathbf{v}_+ = v_1, \dots, v_B$ **do**

 compute hidden unit mean $h_+ = \max(\mathbf{W}\mathbf{v}_+ + \lambda, 0)$;

 compute visible unit mean $v_- = \max(\mathbf{W}^\top \mathbf{h}_+ + \mathbf{b}, 0)$;

 compute hidden unit mean $h_- = \max(\mathbf{W}\mathbf{v}_- + \lambda, 0)$;

 compute CD-1 gradients $\Delta\theta^{(t)}$;

$$\Delta\mathbf{W} = \mathbf{CD}(\mathbf{v}^\top \mathbf{h}), \quad (4)$$

$$\Delta\lambda = \mathbf{CD}(h), \quad \Delta\mathbf{b} = \mathbf{CD}(\mathbf{v}) \quad (5)$$

 compute ADAGRAD learning rate $\gamma_k = \eta / \sqrt{\sum_t \Delta\theta_k^{(t)}}$;

 add momentum $\Delta\hat{\theta}_k^{(t)} = \Delta\theta_k^{(t)} + \tau\Delta\theta_k^{(t-1)}$;

 adjust model $\theta_k = \theta_k + \Delta\hat{\theta}_k^{(t)}\gamma_k, k = 1 \dots K$;

Figure 5: Pretrain RBMs of ReL visible and hidden units.

We refer readers interested in deeper technical details, such as the probabilistic model of RBMs and derivation of CD-1 gradients, to the work by Hinton and et. al [3].

In our implementation, we set the bias on the hidden units to zero. The bias and conditional variance on the visible units and units are initialized as the mean and variance of the training data. The bipartite weights of the RBM are initialized from random samples in the uniform distribution $\mathcal{U}(-0.005, 0.005)$ for all RBMs.

We randomly divide the training data into $B = 5$ batches. The step-size for adjusting each parameter is determined by ADAGRAD[1] which depends on a global learning rate $\eta = 0.005$ for Gaussian-RBMs and a much smaller rate of $\eta = 0.0005$ for the top-level RBM. To accelerate learning, we add momentum to the gradient with $\tau = 0.9$. The training generally converges in $T = 1000$ epochs.

5. Learning image quality with labeled data

Given the pre-trained deep belief network of L layers, we can generate mid-level representations corresponding to the labeled distorted image data x for which we know their ground truth quality scores y . We then model the joint distribution of y and x as via a Gaussian process regression formulation using a simple squared-exponential kernel

$$\mathcal{K}_{ij} = \exp\left(-\frac{1}{2D}|\mathbf{h}_i^{(L)} - \mathbf{h}_j^{(L)}|^2\right). \quad (7)$$

where D is the dimension of the neural network output.

The vector $\mathbf{h}_i^{(L)}$ corresponds to the activation of the neural network network due to input x_i

$$\mathbf{h}_i^{(l)} = \max(\mathbf{W}^{(l)}\mathbf{v}_i^{(l)} + \lambda^{(l)}, 0) \quad (8)$$

$$\mathbf{v}_i^{(l)} = \mathbf{h}_i^{(l-1)}, \mathbf{v}_i^{(0)} = x_i. \quad (9)$$

We adjust the weights of the neural network and hyper-parameters of the Gaussian process δ to maximize the log likelihood function

$$\mathcal{L} = \log p(x, y) \propto -\frac{1}{2} \log |\mathcal{K}_\delta| - \frac{1}{2} y^\top (\mathcal{K}_\delta)^{-1} y, \quad (10)$$

$$\text{where } \mathcal{K}_\delta = \mathcal{K} + \delta \mathbf{I}. \quad (11)$$

The parameter δ denotes the regression noise model variance for the GP likelihood.

The partial derivatives of the log likelihood can be analytically computed by the chain rule.

$$\frac{\partial \mathcal{L}}{\partial \delta} = \text{tr}\left(\frac{\partial \mathcal{L}}{\partial \mathcal{K}_\delta}\right) \quad (12)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(k)}} = \sum_i \frac{\partial \mathcal{L}}{\partial \mathbf{h}_i^{(L)}} \prod_{l=k}^{L-1} \frac{\partial \mathbf{h}_i^{(l+1)}}{\partial \mathbf{h}_i^{(l)}} \mathcal{S}(\mathbf{W}^{(l)} \mathbf{h}^{(l)} + \lambda^{(l)}) \mathbf{h}^{(l)} \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda^{(k)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_i^{(l)}} \prod_{l=k}^{L-1} \frac{\partial \mathbf{h}_l(x_i)}{\partial \mathbf{h}_{l+1}(x_i)} \mathcal{S}(\mathbf{W}_l \mathbf{h}_l + \lambda_l) \quad (14)$$

where $\mathcal{S}(z)$ is the step function whose value is 1 when $z > 0$ and 0 otherwise, and

$$\frac{\partial \mathcal{L}}{\partial \mathcal{K}_\delta} = -\frac{1}{2} ((\mathcal{K}_\delta^{-1} y)^\top (\mathcal{K}_\delta^{-1} y) - \mathcal{K}_\delta^{-1}), \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_i^{(L)}} = \frac{1}{D} \sum_j \left[\frac{\partial \mathcal{L}}{\partial \mathcal{K}_\delta} \cdot \mathcal{K} \right]_{ij} (\mathbf{h}_j^{(L)} - \mathbf{h}_i^{(L)}), \quad (16)$$

$$\frac{\partial \mathbf{h}_i^{(l+1)}}{\partial \mathbf{h}_i^{(l)}} = \mathbf{W}^{(l)\top} \mathcal{S}(\mathbf{W}^{(l)} \mathbf{h}^{(l)} + \lambda^{(l)}) \quad (17)$$

Again, we use a momentum of 0.9 to accelerate training. We perform gradient descent to fine tune the model by adjusting the parameters with a constant rate $\tau = 0.01$.

Prediction: During the testing phase, we again use the deep belief network to compute the mid-level representation and then make a prediction about the image quality via Gaussian Process regression. The regression module then predicts the image quality scores y_n for features x_n of unseen images as a Gaussian distribution of mean

$$\bar{y}_n = \mathcal{K}(\mathbf{h}^{(N)}(x_n), \mathbf{h}^{(N)})^\top (\mathcal{K} + \delta \mathbf{I})^{-1} y. \quad (18)$$

6. Results

In this section, we empirically evaluate our model in supervised and semi-supervised settings and also present an evaluation of generalizing across distortion types. We perform evaluation on the LIVE and TID2008 datasets. In the semi-supervised setting, we use simulated data (as discussed in Sec. 3.1) to pre-train the model. The criterion we use for evaluation is Spearman correlation that reflects how

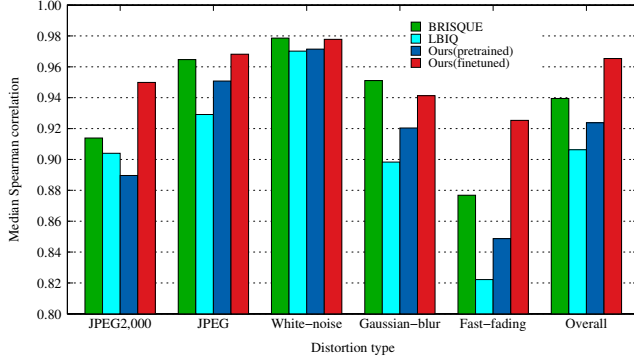


Figure 6: Performance of our fine-tuned and pretrained model and two state-of-art methods on the LIVE dataset.

well rankings are preserved in the predicted measure. To eliminate bias due to division of the data, we perform a repeated random subsampling experiment of 1000 trials in all experiments on both the LIVE and the TID2008 datasets. We use mean Spearman correlation to evaluate the quality of estimate. Larger correlation indicates higher relevance.

6.1. Performance under supervised setting

We first evaluate our model assuming sufficient labeled data are used for training on the LIVE dataset. Distorted images of 23 out of 29 references are used for training and the remaining images are used for testing.

We evaluate the performance of two variations of our model as well as two state-of-art blind measures:

- **A pretrain-only model** is learned by first pretraining the neural network layer by layer and then maximize the Gaussian process likelihood by adjusting the hyper-parameter δ only.
- **Our final fine-tuned model** as described in Sec. 3.
- **The baseline LBIQ measure** performs SVM regression after PCA of the same features as our measure.
- **The BRISQUE measure**[6] performs SVM regression on a different set of features, and is the one of best-performing blind measures to our knowledge.

Fig. 6 shows the Spearman correlation between the predicted quality measure and the groundtruth human quality scores on the 5 distortion types as well as the overall correlation. We observe that even without fine-tuning, overall our model performs significantly better than the LBIQ model. This indicates the non-linearity of the RBM model greatly contributes to the success of our model and verifies our assumption that the unsupervised pre-training only can lead to good features for regression. With refinement, our model performs even better and achieves a performance slightly better than the state-of-art BRISQUE measure (0.965 v.s.

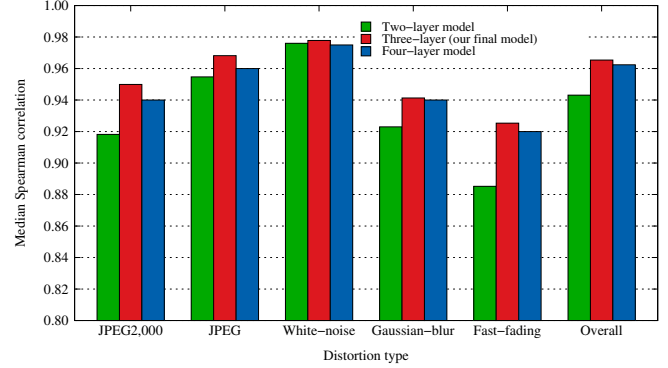


Figure 7: Performance of pre-trained model on LIVE dataset with different neural network depths.

0.939). This supports our argument that a good learning model plays a critical role in image quality measures.

Second, we investigate how the performance of the model relates to the depth of the neural network. We compare our three-layer model with a four-layer model with an additional layer of 600 hidden units (the number of hidden units is decided by cross-validation) and a two-layer model that simply feeds the output of the seven RBMs to the Gaussian process model. We observe that the three-layer model performs slightly better than the other two. This indicates that it is necessary to have a mid-layer representation to combine the seven RBMs together, but the non-linearity captured by the three-layer model is sufficient to represent the structure of the specific data and features we use.

6.2. Performance with semi-supervised setting

We first conduct experiments to investigate the performance variation of our model and state-of-art models with different amount of labeled and unlabeled data. In each trial, we pretrain our model with unlabeled distorted images of $\{5, 20, 80\}$ reference images and fine-tuned by labeled distorted images of $\{3, 7, 11, 15, 19, 23\}$ reference images. For the comparison algorithms (LBIQ and BRISQUE), we perform regression on the same labeled images alone.

Fig. 8 plots the performance degradation with the reduction of labeled data for our method (under three conditions of amount of unlabeled data used) and those from previous work. We make two observations about the results. First, our model appears more robust than the LBIQ and BRISQUE measure with the decrease of labeled data. To achieve a performance comparable to state-of-art (Spearman correlation $\geq .9$) BIQA measures, both LBIQ and BRISQUE needs at least labeled distortion image of 15 reference images, yet with sufficient unlabeled data, our model performs well with just labeled images of 7 references, and the performance of our semi-supervised model with just 3 labeled images is still in the reasonable range (0.85 –

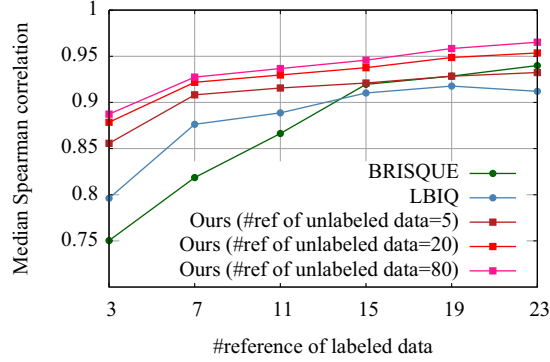


Figure 8: Performance variation on LIVE dataset with different amount of labeled data for our method and two previous methods: BRISQUE and LBIQ.

0.9). Second, with increasing unlabeled data to pretrain the model, the regression performance steadily improves. When pretraining with unlabeled data of 20 reference images, the performance is nearly as good as when pretraining with the entire unlabeled dataset of 80 references.

Finally, we compare the performance of our model and the state-of-art methods on the TID2008 dataset³. A repeated random subsampling experiment of 1000 trials is conducted for fair comparison. Due to the difficulty of this dataset, we use the entire unlabeled dataset for pretraining and distortion images of 20 reference images for testing. As shown in Fig. 9, the overall performance of our method is 0.841, and it is much better than state-of-the-art methods: LBIQ is 0.74 and BRISQUE is 0.61. Though it is not as good as the original LBIQ metric for a few distortion types. This is because the Gaussian process in trying to reconcile among the 17 distortion types, sacrifices the performance on individual distortion types.⁴

6.3. Generalization across distortion types

Finally, we explore the ability of our model to generalize across distortion types.

We first visualize the low-dimensional embedding of the LIVE and TID2008 dataset in Fig. 10 to gain some intuition on why generalization across distortion types are possible. We scatter the eigenspace projection of the neural network output of LIVE and TID2008 dataset and color code the corresponding distortion types or subjective image quality. Fig. 10 shows not only continuity in subjective image quality in the eigenspace but also a clustering across similar

³We exclude the 25th image of the dataset from testing because it is a synthetic image and its feature significantly differs from natural images.

⁴We have also tried to pick images of distortion types that are well modeled by the LBIQ features for training and all images for testing. This results in better performance than LBIQ in the specific distortion types we use for training but worse overall performance.

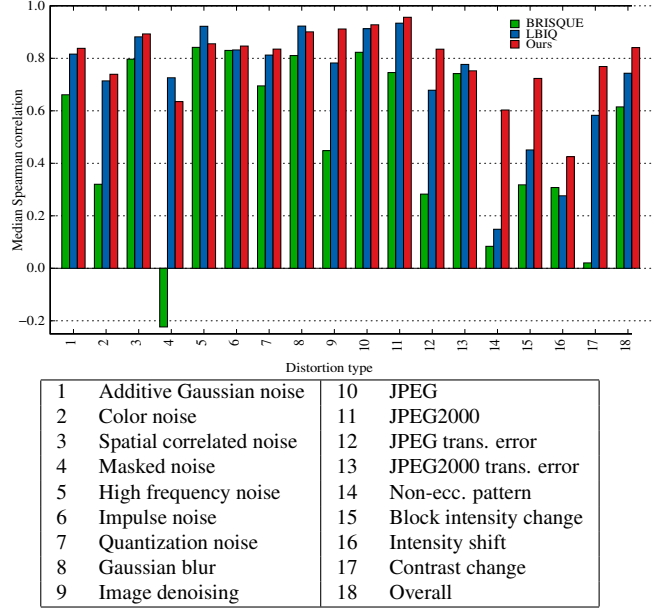


Figure 9: Performance of LBIQ, BRISQUE and our method on the TID2008 dataset. Our model is trained using $80 \times 13 \times 5$ unlabeled and $20 \times 17 \times 5$ labeled images. LBIQ and BRISQUE are trained using the same labeled data.

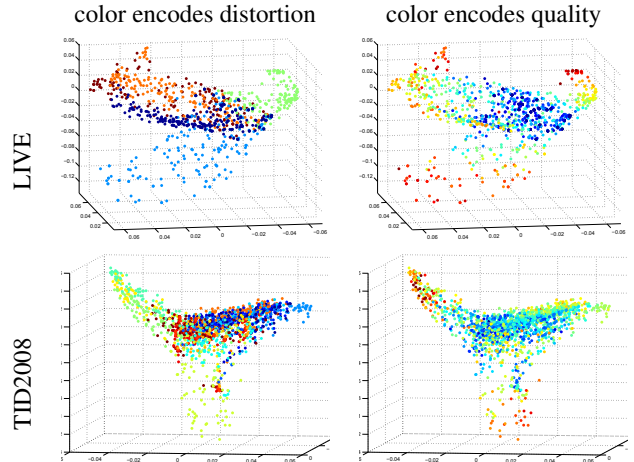


Figure 10: Kernel eigen-space embedding of the LIVE and TID2008 dataset. Similar distortion types are clustered together and quality changes smoothly in this eigenspace.

lar distortion types (different types of noise, distortion, and compression). This indicates that the quality measure of a distortion type can be generalized from labels for similar distortion types.

To validate this idea, we perform a “leave-one-distortion-type-out” experiment on the LIVE dataset. For each distortion type, we use all data for other distortion types to train

	distortion type	Spearman correlation
1	JPEG2000	0.9580
2	JPEG	0.9512
3	White noise	-0.4879
4	Gaussian blur	0.9692
5	Fast fading	0.9436

Figure 11: Spearman correlation of predicted and subjective quality in leave-one-distortion-type-out setting.

the model and estimate the quality measure of the left-out distortion type. As shown by Fig. 11, the blur and compression types are predicted well. Yet the white-noise type is poorly measured because it is very different from others.

7. Conclusions and Future Work

State-of-art image quality assessment techniques use kernel regression methods to measure image quality using training sets of distorted images. The success of such methods rely on the kernel function to adequately embed the training data into a quality-relevant sub-space. It also requires a reasonable amount of training data to exemplify the image quality measure.

We propose to represent the kernel function for image quality assessment with a rectifier neural network. The many degrees of freedom and non-linearity of such a model allows it to represent the structure of image distortions with flexibility. The ability to perform unsupervised pre-training of the model allows us to use a large volume of unlabeled image data to train the model without being restricted by the limited access to human scores. Experiments show that our method leads to significant improvement over previous methods (both blind and non-blind) for two challenging datasets, and robustness to reduction of labels.

Our work shows the potential to exploit advanced learning models to overcome challenges in blind image quality assessment. We believe an adequate learning model is as important as other aspects of IQA measures, such as feature crafting and pooling, which have been deeply explored.

Although the model we propose may generally benefit a wide range regression-based IQA metrics, our exploration is currently limited in two aspects. First, we have only applied our model to a specific set of features [14]. We believe that by building on an existing, well known feature set makes the learning contribution clear, and making the learning algorithm the only variable was best scientifically. We believe though that our model can benefit a wide range of handcrafted features as long as they sufficiently express the structure of the data. Therefore as future work, we think using other existing features, combined, or new features is a good direction; as is investigating the benefits of specific features by looking at the weights of the top level RBM.

Second, we only increase unlabeled samples in the dimension of reference images. Therefore, it is not clear whether or how our method can be extended to handle the

increase of distortion types. However, we do show that our current model can handle some unseen distortion types to some extent as it exploits statistical co-dependencies across various kinds of distortions. Exploring and addressing the above limitations is a promising direction for future work.

References

- [1] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. of Machine Learning Research*, 12:2121–2159, July 2011.
- [2] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *NIPS’2010 Workshop on Deep Learning and Unsupervised Feature Learning*, Apr. 2010.
- [3] G. E. Hinton, S. Osindero, and Y. Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- [4] G. E. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 2006.
- [5] C. Li and A. C. Bovik. Content-partitioned structural similarity index for image quality assessment. *Image Communication*, 25(7):517–526, 2010.
- [6] A. Mittal, A. Moorthy, and A. Bovik. No-reference image quality assessment in the spatial domain. *Image Processing, IEEE Transactions on*, 21(12):4695–4708, 2012.
- [7] A. Mittal, R. Soundararajan, and A. Bovik. Making a completely blind image quality analyzer. *Signal Processing Letters, IEEE*, 20(3):209–212, 2013.
- [8] A. K. Moorthy and A. C. Bovik. A two-step framework for constructing blind image quality indices. *IEEE Signal Processing Letters*, 17(5):513–516, 2010.
- [9] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- [10] N. Ponomarenko, V. Lukin, A. Zelensky, K. Egiazarian, M. Carli, and F. Battisti. Tid2008 - a database for evaluation of full-reference visual quality assessment metrics. *Advances of Modern Radioelectronics*, 10:30–45, 2009.
- [11] R. Salakhutdinov and G. Hinton. Using deep belief nets to learn covariance kernels for gaussian processes. In *NIPS*, 2007.
- [12] R. Salakhutdinov and G. Hinton. Semantic hashing. *Int. J. Approx. Reasoning*, 50:969–978, July 2009.
- [13] H. R. Sheikh, Z. Wang, L. Cormack, and A. C. Bovik. *LIVE image quality assessment database release 2*. <http://live.ece.utexas.edu/research/quality>.
- [14] H. Tang, N. Joshi, and A. Kapoor. Learning a blind measure of perceptual image quality. In *Computer Vision and Pattern Recognition, IEEE Conference on*, pages 305–312, 2011.
- [15] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Trans. on Image Processing*, 3:600–612, 2004.
- [16] W. Xue, L. Zhang, and X. Mou. Learning without human scores for blind image quality assessment. In *CVPR*, 2013.
- [17] P. Ye, J. Kumar, L. Kang, and D. S. Doermann. Unsupervised feature learning framework for no-reference image quality assessment. In *CVPR*, pages 1098–1105, 2012.
- [18] J. Zhu. Semi-supervised learning literature survey. 2008.