# Joint Regression and Ranking for Image Enhancement

Parag Shridhar Chandakkar
Arizona State University
pchandak@asu.edu

Baoxin Li
Arizona State University
baoxin.li@asu.edu

## Abstract

*Research on automated image enhancement has gained momentum in recent years, partially due to the need for easy-to-use tools for enhancing pictures captured by ubiquitous cameras on mobile devices. Many of the existing leading methods employ machine-learning-based techniques, by which some enhancement parameters for a given image are found by relating the image to the training images with known enhancement parameters. While knowing the structure of the parameter space can facilitate search for the optimal solution, none of the existing methods has explicitly modeled and learned that structure. This paper presents an end-to-end, novel joint regression and ranking approach to model the interaction between desired enhancement parameters and images to be processed, employing a Gaussian process (GP). GP allows searching for ideal parameters using only the image features. The model naturally leads to a ranking technique for comparing images in the induced feature space. Comparative evaluation using the ground-truth based on the MIT-Adobe FiveK dataset plus subjective tests on an additional data-set were used to demonstrate the effectiveness of the proposed approach.*

## 1. Introduction

The corpus of images on the Web is exponentially increasing in size with close to two billion photos being added or circulated each day[1]. Image sharing has become an integral part of daily life for many people, and they want their photos to look good without doing too much manual editing. Some tools for easy image enhancement have already been deployed on popular social networking platforms, such as Instagram, or mobile devices such as iPhone. However, most such tools are essentially based on some pre-defined image filters for obtaining certain visual effects. Recent research efforts on automated image enhancement employing machine learning techniques for improved functionalities such as content adaptivity and personalization.

Such solutions range from learning a tone mapping between the spaces of low-quality and high-quality images[2] to building a ranking relation between these two spaces [2, 5, 13, 9, 12, 14, 6], although we are yet to see such techniques being deployed on a popular platform.

Many recent approaches follow the pipeline shown in Fig. 1. During training, these approaches learn a model to assign a score for a given image quantifying its *visual appeal*. To enhance a new image, a nearest training image is found. Then a dense sampling around the parameters[3] of the high-quality counterpart of that training image gives the candidate set of enhancement parameters for the new image. A set of candidate images is then generated by applying these enhancement parameters to the new image. The next step is extracting features of all candidate images and use the learned model to select the highest-quality image.

There are two major drawbacks in the above processing flow. First, it is computationally expensive at the testing phase since a search through the entire training data is needed. The training data for such applications could be of huge size and is usually hosted on a server. Thus hundreds of thousands people querying the server per second is undesirable and uncalled-for. Second, the set of candidate parameters which would enhance the original image is found in a sub-optimal manner by doing a $k$NN search. It does not provide any structured way to search for the optimal parameters and thus it becomes necessary to search the entire training set and create a lot of candidate images, resulting a computational bottleneck for the testing phase.

In this paper, we develop a joint regression and ranking approach to address the above drawbacks. Our approach employs GP regression to predict the mean and variance of the candidate parameters *only* from the feature vector of a low-quality image. We also simultaneously train a ranking model on the GP-covariance-kernel-induced feature space. To achieve this, we derive and use the dual-form of ranking SVM [10] with the GP kernel integrated into it. Thus

---

[1] http://www.kpcb.com/internet-trends

[2] We call the images before enhancement as low-quality and those after the enhancement as high-quality in the rest of this article. We also refer to the process of enhancing a new picture as "the testing stage".

[3] The brightness, saturation and contrast are referred to as "parameters" of an image in this article

(a) A typical machine-learned image enhancement approach



(b) Proposed approach. Note that our approach has no interaction with the training data during the testing stage
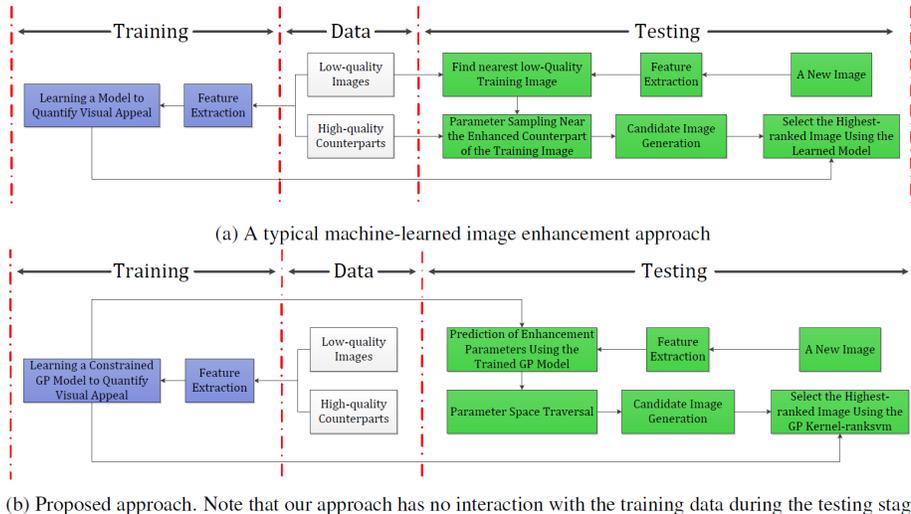
Figure 1. Pipelines of image enhancement approaches.

the kernel builds a relation between the image feature space and its corresponding enhancement parameter space. Along with that, the kernel learns to give more weight to the image features which are highly responsible for making an image to be of higher quality. Finally, we learn the GP kernel in such a way that all the high-quality counterparts of a low-quality image form a cluster. This allows exploration of the image parameter space in a structured manner for obtaining the optimal solution.

In the testing stage (i.e., while enhancing a new image), the model provides the expected value and variances of the enhancement parameters, drastically reducing required computation since there is no need to sift through the training set. We can generate/show some enhanced images by applying parameters that are $k$ standard deviations away from the expected values of the parameters, where $k$ is a user-defined and can be changed on-the-fly, so that user can choose from some good candidate images. The same kernel can be used to rank images if the user wants to see a single image. Through extensive experiments, we show that our approach is computationally efficient at the testing phase, and that it predicts parameters correctly for new images and that it also predicts the ranking relations between the new images and its enhanced counterparts.

## 2. Related Work

Automated image enhancement has recently been an active research area. Various solutions have been proposed for this task. We review those works which aim to improve the visual appeal of an image using automated techniques. A novel tone-operator was proposed to solve the tone reproduction problem [18]. A database named MIT-Adobe FiveK of corresponding low and high-quality images was published in [5]. They also proposed algorithm to solve the

problem of global tonal adjustment. The tone adjustment problem only manipulates the luminance channel. In [11], an approach was presented, focusing on correcting images containing faces. They built a system to align faces between a "good" and a "bad" photo and then use the good faces to correct the bad ones.

Content-aware enhancement approaches have been developed which aim to improve a specific image region. Some examples of such approaches are [2, 14]. A drawback of these is the reliance on obtaining segmented regions that are to be enhanced, which itself may prove difficult. Pixel-level enhancement was performed by using local scene descriptors. First, images similar to the input are retrieved from the training set. Then for each pixel in the input, a set of pixels was retrieved from the training set and they were used to improve the input pixel. Finally, Gaussian random fields are used to maintain the spatial smoothness in the enhanced image. This approach does not take the global information of an image into account and hence the local adjustments may not look right when viewed globally. A deep-learning based approach was presented in [26]. In [12], users were required to enhance a small amount of images to augment the current training data.

Two closely related and recent works involve training a ranking model from low and high-quality image pairs [25, 6]. In a recent state-of-art method [25], a dataset of 1300 corresponding image pairs was reported, where even the intermediate enhancement steps are recorded. A ranking model trained with this information can quantify the (enhancement) quality of an image. In [6], non-corresponding low and high-quality image pairs were used to train a ranking model. Both the approaches use $k$NN search at the test time to create a pool of candidate images first. After extracting features and ranking all of them, the best image is

presented to the user.

Now we briefly review Gaussian process based methods which are relevant in this context. GP has been effectively used to obtain good performance for applications where complex relationships have to be learned using a small amount of data (in the order of several hundreds) [23]. In [7], it was used for view-invariant facial recognition. A GP latent variable model was used to learn a discriminative feature space using LDA prior where examples from similar classes are project nearby. In [20], GP regression was used to map the non-frontal facial points to the frontal view. Then facial expression methods can be used using these projected frontal view points. Coupled GP have been used to capture dependencies between the mappings learned between non-frontal and frontal poses, which improves the facial expression recognition performance [19].

Our effort in this paper deals with enhancement considering contrast, saturation and brightness of an image. We attempt to explicitly model interactions between parameters controlling these factors and features extracted from the underlying image, employing GP. Our approach of joint regression and ranking allows us to learn the complex mapping from the image features to the regions corresponding to desired enhancement in the parameters space, without actually generating several hundreds of enhanced candidate images. The expected value of the parameters and their standard deviations provide us with a way to systematically explore the parameters space. In the next section, we detail our proposed approach. To the best of our knowledge, this is the first attempt on incorporating GP regression and ranking for image enhancement.

## 3. Proposed Approach

Our problem is to predict the set of image parameters which would enhance a given image. Our proposed approach consists of two objectives: 1. Given a low-quality image feature, probabilistically estimate the parameters that could generate the enhanced counterpart. 2. Produce a ranking in the GP-kernel-induced feature space and thereby discover the features responsible for making an image of higher-quality.

We have pairs of low and high-quality images along with their parameters for training. The feature representation of an image will be discussed in detail in section 3.6. Features of $N$ low-quality images are represented by $\boldsymbol{F} = \{\boldsymbol{f}_1, \boldsymbol{f}_2, \ldots, \boldsymbol{f}_N\}^4$. We have $p$ high-quality versions for each low-quality image. Its features are represented by $\boldsymbol{F}^+ = \{\boldsymbol{F}_1^+, \ldots, \boldsymbol{F}_N^+\}$, where $\boldsymbol{F}_i^+ = \{\boldsymbol{f}_{i1}^+, \ldots, \boldsymbol{f}_{ip}^+\}$, and $\boldsymbol{f}_i, \boldsymbol{f}_{ij}^+ \in \mathbb{R}^{D \times 1} \, \forall \, i, j$. We also have $p$ sets of high-quality parameters for a low-quality image. However, for illus-

---

$^4$We represent vectors by lower-case bold letters. Matrices are represented by upper-case bold letters. Scalars are denoted by non-bold letters.

tration, we predict parameters only for the first set. Note that we still use all the $p$ sets of high-quality images to train a ranking model. The parameters of low and high-quality images are represented by $\boldsymbol{Y} = \{\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N\}$ and $\boldsymbol{Y}^+ = \{\boldsymbol{y}_1^+, \ldots, \boldsymbol{y}_N^+\}$ respectively. We use three image parameters, namely, brightness, contrast and saturation and hence $\boldsymbol{y}_i, \boldsymbol{y}_i^+ \in \mathbb{R}^{3 \times 1} \, \forall \, i$. Our task is to obtain $\boldsymbol{y}_i^+$ using only $\boldsymbol{f}_i$ and $\boldsymbol{y}_i$. We predict each parameter using a separate GP. To that end, we collect the $m^{th}$ parameter of all low and high-quality images into, $\bar{\boldsymbol{y}}_m = (y_{1m}, \ldots, y_{Nm})^T$ and $\bar{\boldsymbol{y}}_m^+ = (y_{1m}^+, \ldots, y_{Nm}^+)^T$, respectively and train a separate GP model to predict each parameter. We now outline the proposed joint GP regression and ranking.

### 3.1. GP Regression

GPs define a prior distribution over functions which becomes a posterior over functions after observing the data. GPs assume that this distribution over functions is jointly Gaussian with a mean and a positive definite covariance kernel function. GPs provide well-calibrated, probabilistic outputs which are particularly useful and necessary in our application [15]. If we let the prior on regression function be a GP, then it can be denoted as $GP(m(\boldsymbol{f}), \kappa(\boldsymbol{f}, \boldsymbol{f}'))$ where $\boldsymbol{f}$ and $\boldsymbol{f}'$ are image features $\in \mathbb{R}^{D \times 1}$ as defined previously, $m(\boldsymbol{f})$ is a mean function and $\kappa(\boldsymbol{f}, \boldsymbol{f}')$ is a covariance function. It is well-known that GPs are flexible enough to model an arbitrary mean. It can be shown that the posterior predictive density for a single test input is:

$$p(\bar{y}_{*m}^+ | \boldsymbol{f}_*, \boldsymbol{F}, \boldsymbol{Y}) = \mathcal{N}(\bar{y}_{*m}^+ | \boldsymbol{k}_*^T \boldsymbol{K}_y^{-1} \bar{\boldsymbol{y}}_m^+, k_{**} - \boldsymbol{k}_*^T \boldsymbol{K}_y^{-1} \boldsymbol{k}_*)$$

(1)

where $\boldsymbol{k}_* = [\kappa(\boldsymbol{f}_*, \boldsymbol{f}_1), \ldots, \kappa(\boldsymbol{f}_*, \boldsymbol{f}_N)]$, $N$ is the number of samples, $k_{**} = \kappa(\boldsymbol{f}_*, \boldsymbol{f}_*)$ and $\boldsymbol{K}_y = \boldsymbol{K} + \sigma_y^2 \boldsymbol{I}_N$. $\boldsymbol{K}$ is a kernel function between all training inputs $\boldsymbol{f}$, and $(\cdot)_*$ denotes a new data sample. The noise or uncertainty in the output is modeled by the noise variance, $\sigma_y^2$.

It can be further shown that the log-likelihood function for a GP regression model is easily obtained by using a standard multivariate Gaussian distribution as follows,

$$\log p(\bar{\boldsymbol{y}}_m^+ | \boldsymbol{F}) = -0.5 \left(\bar{\boldsymbol{y}}_m^+\right)^T \boldsymbol{K}_y^{-1} \bar{\boldsymbol{y}}_m^+ - 0.5 log |\boldsymbol{K}_y| - $$
$$0.5 N log(2\pi)$$

(2)

We choose a standard squared exponential kernel for our application. It is as follows,

$$\kappa(\boldsymbol{f}_i, \boldsymbol{f}_j) = \sigma_f^2 \exp(-\frac{1}{2}(\boldsymbol{f}_i - \boldsymbol{f}_j)^T \cdot \boldsymbol{\Lambda} \cdot (\boldsymbol{f}_i - \boldsymbol{f}_j)) + \sigma_y^2 \delta_{ij}$$

(3)

Here, $\sigma_f^2$ controls the vertical scale of the regression function, $\sigma_y^2$ models uncertainty, $\boldsymbol{\Lambda}$ is a diagonal matrix with entries $\{\theta_1, \ldots, \theta_D\}$ and $\delta_{pq}$ is a Kronecker delta function, which takes the value 1 if $p = q$, it is zero everywhere else. We call $\boldsymbol{h} = \{\sigma_f^2, \boldsymbol{\Lambda}, \sigma_y^2\}$ as hyper-parameters. It is easy

to see that the prediction in Equation 1 is dependent on the kernel and in turn on the hyper-parameters: $\sigma_f$, $\boldsymbol{\Lambda}$ and $\sigma_y^2$. We will see later how to obtain optimal hyper-parameters. Now, we explain inclusion of ranking into our formulation.

## 3.2. GP Ranking

We build a ranking relation in the GP-kernel-induced feature space. Thus the GP kernel discovers the features responsible for making an image of higher-quality and assigns higher weight to them by adjusting the hyper-parameters. The primal form of rank SVM [10] is given by:

$$\min_{w,\xi_{ij}} \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C\sum_{i,j}\xi_{ij}, \text{ subject to: } \boldsymbol{u}_i \succ \boldsymbol{u}_j \ \forall \ (i,j)$$
(4)

where $\boldsymbol{u}_i \succ \boldsymbol{u}_j$ indicates that $\boldsymbol{u}_i$ is ranked higher than $\boldsymbol{u}_j$.

For ranking, we observe that only building a relation between low and high-quality images does not provide good ranking accuracy on new images. The reason is that the enhanced images often possess high saturation, brightness and/or contrast. Thus the ranking model sometimes assigns a higher score to over-saturated and over-exposed images. This would not be a problem if one had intermediate information about the enhancement steps being performed [25]. Thus we deteriorate our original low-quality images by shifting the image parameters to both extremes. The amount of shifting for an image is decided by first deteriorating 20 images manually and then heuristically defining a relation between existing image parameters and the amount of shifting needed to significantly deteriorate the image. Let's call these images as *poor-quality* images. We also generate $p$ poor-quality images for every low-quality image. We now have features for poor, low and high quality images, denoted by $\boldsymbol{F}^-$, $\boldsymbol{F}$ and $\boldsymbol{F}^+$ respectively. Primal form for our ranking model can be written as follows,

$$\min_{w,\xi_{ij}} \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C_1\sum_{i,j}\xi_{ij} + C_2\sum_{i,k}\xi'_{ik},$$
$$\text{subject to: } \boldsymbol{w}^T\boldsymbol{f}_{ij}^+ \geq \boldsymbol{w}^T\boldsymbol{f}_i + 1 - \xi_{ij},$$
$$\text{subject to: } \boldsymbol{w}^T\boldsymbol{f}_i \geq \boldsymbol{w}^T\boldsymbol{f}_{ik}^- + 1 - \xi'_{ik},$$
$$\text{subject to: } \boldsymbol{w}^T\boldsymbol{f}_{ij}^+ \geq \boldsymbol{w}^T\boldsymbol{f}_{ik}^- + 1 - \xi''_{ik}, \xi_{ij}, \ \xi'_{ik}, \ \xi''_{ik} \geq 0$$
$$\forall \ i = \{1,\dots,N\}, \forall j = \{1,\cdots,p\}, \forall k = \{1,\cdots,p\}.$$
(5)

Now we derive the dual form to incorporate the GP-kernel, $\kappa$. It would be cumbersome to derive the dual of Equation 5 as it stands. Instead, we can define a new set of data $\boldsymbol{D}$ consisting of $\boldsymbol{f}_i - \boldsymbol{f}_{ij}^+$, $\boldsymbol{f}_{ik}^- - \boldsymbol{f}_i$ and $\boldsymbol{f}_{ik}^- - \boldsymbol{f}_{ij}^+ \ \forall \ i,j,k$. Now, $\boldsymbol{D}$ has $N' = N(2p + p^2)$ elements, we can write the primal

form as follows:

$$\min_{w,\xi_i} \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C\sum_i \xi_i,$$
$$\text{subject to: } \boldsymbol{w}^T\boldsymbol{D}_i + 1 - \xi_i \leq 0, \xi_i \geq 0, \forall i = \{1,\dots,N'\}.$$
(6)

We use Lagrangian multipliers to convert the above equation into an unconstrained optimization problem.

$$L(\boldsymbol{w},\boldsymbol{\alpha},\boldsymbol{\beta}) = \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C\sum_i \xi_i +$$
$$\sum_i \alpha_i(\boldsymbol{w}^T\boldsymbol{D}_i + 1 - \xi_i) - \sum_i \beta_i\xi_i$$
(7)

Differentiating with respect to $\boldsymbol{w}$ and $\xi$ and equating them to zero, we get,

$$\nabla_w L(\boldsymbol{w},\boldsymbol{\alpha},\boldsymbol{\beta}) = 0 \Rightarrow \boldsymbol{w} = -\sum_i \alpha_i\boldsymbol{D}_i$$
$$\nabla_\xi L(\boldsymbol{w},\boldsymbol{\alpha},\boldsymbol{\beta}) = C - \alpha_i - \beta_i = 0 \Rightarrow \alpha_i \leq C.$$
(8)

Substituting $\boldsymbol{w}$ back into Equation 6 and doing some algebraic manipulation, we get a dual maximization problem as,

$$\max_{\boldsymbol{\alpha}} \sum_i \alpha_i - \frac{1}{2}\sum_i\sum_j \alpha_i\alpha_j\boldsymbol{D}_i^T\boldsymbol{D}_j, \text{ subj. to: } 0 \leq \alpha_i \leq C.$$
(9)

The inner product in the above equation can be replaced with GP kernel by employing the kernel trick. Thus the final optimization problem to get $\boldsymbol{\alpha}$ becomes,

$$\max_{\boldsymbol{\alpha}} \boldsymbol{1}^T\boldsymbol{\alpha} - \frac{1}{2}\boldsymbol{\alpha}^T\boldsymbol{K}_y\boldsymbol{\alpha}.$$
(10)

Here, $\boldsymbol{1}$ is a column vector of ones. The length of both $\boldsymbol{\alpha}$ and $\boldsymbol{1}$ is $N(2p + p^2)$. The dimensions of $\boldsymbol{K}_y$ are $N(2p + p^2) \times N(2p + p^2)$. The $(i,j)^{th}$ element of $\boldsymbol{K}_y$ is $\kappa(\boldsymbol{D}_i,\boldsymbol{D}_j)$.

## 3.3. Clustering high-quality images together

We turn our attention to the third constraint. Given a low-quality image: 1. it forces all its high-quality counterparts to form a cluster and 2. it tries to maximize the distance between poor-quality and high-quality images in the GP-kernel-induced feature space. The intuition behind this is as follows. Ultimately, given a new image, we would not only like to predict the parameters for its enhanced counterpart, but we also wish to traverse the parameter space and explore more of such enhancement parameters. The traversing of the parameter space is, in our opinion, essential since the choices of people vary by a great amount and no model would do justice with just one set of predicted parameters. Note that this constraint tries to minimize distance between $\boldsymbol{f}_i$ and $\boldsymbol{f}_{ij}^+ \ \forall j$, so by definition of GP, the distance between the corresponding output parameters, $\boldsymbol{y}_{ij}^+ \ \forall j$, will

be reduced, which in turn achieves the aforementioned effective traversal. The second part of the constraint tries to push the predicted parameters away from the parameters of the poor-quality images. The details of traversing the parameter space after getting the GP predictions are discussed later. The constraint can be formulated as follow.

$$\min_{h} \left( \sum_i ||\boldsymbol{K}_y^{\boldsymbol{F}_i^+}||_F^2 - ||\boldsymbol{K}_y^{\boldsymbol{F}_i^+, \boldsymbol{F}_i^-}||_F^2 \right), \qquad (11)$$

where $|| \cdot ||_F^2$ indicates squared Frobenius norm. The term $\boldsymbol{K}_y^{\boldsymbol{F}_i^+, \boldsymbol{F}_i^-}$ is a $p \times p$ matrix defined as follows,

$$\boldsymbol{K}_y^{\boldsymbol{F}_i^+, \boldsymbol{F}_i^-} = \begin{bmatrix} \kappa(\boldsymbol{f}_{i1}^+, \boldsymbol{f}_{i1}^-) & \cdots & \kappa(\boldsymbol{f}_{i1}^+, \boldsymbol{f}_{ip}^-) \\ \kappa(\boldsymbol{f}_{i2}^+, \boldsymbol{f}_{i1}^-) & \cdots & \kappa(\boldsymbol{f}_{i2}^+, \boldsymbol{f}_{ip}^-) \\ \vdots & \ddots & \vdots \\ \kappa(\boldsymbol{f}_{ip}^+, \boldsymbol{f}_{i1}^-) & \cdots & \kappa(\boldsymbol{f}_{ip}^+, \boldsymbol{f}_{ip}^-) \end{bmatrix} \qquad (12)$$

The term $\boldsymbol{K}_y^{\boldsymbol{F}_i^+}$ is equal to $\boldsymbol{K}_y^{\boldsymbol{F}_i^+, \boldsymbol{F}_i^+}$.

We combine Equations 2, 10 and 11 to form our objective function. It is as follows,

$$\min_{h} \ Z = \frac{1}{2} \left( \bar{\boldsymbol{y}}_m^+ \right)^T \boldsymbol{K}_y^{-1} \bar{\boldsymbol{y}}_m^+ + \frac{1}{2} \log |\boldsymbol{K}_y| - \boldsymbol{1}^T \boldsymbol{\alpha} + \\ \frac{1}{2} \boldsymbol{\alpha}^T \boldsymbol{K}_y \boldsymbol{\alpha} + \sum_i \left( ||\boldsymbol{K}_y^{\boldsymbol{F}_i^+}||_F^2 - ||\boldsymbol{K}_y^{\boldsymbol{F}_i^+, \boldsymbol{F}_i^-}||_F^2 \right) \qquad (13)$$

Note that we have removed the constant term. We now focus on how to solve Equation 10 and 13 to get $\boldsymbol{\alpha}$ and $\boldsymbol{h}$.

## 3.4. Optimization

Our optimization problem is separable in $\boldsymbol{\alpha}$ and $\boldsymbol{h}$. First we optimize $\boldsymbol{\alpha}$, which can be done by using a standard rank-SVM solver. It could also be solved by using quadratic programming, however, that would be memory inefficient. In particular, we use a rank-SVM implementation which uses the LASVM algorithm proposed in [4]. LASVM employs active example selection to significantly reduce the accuracy after just one pass over the training examples.

After optimizing $\boldsymbol{\alpha}$, we turn our attention to Equation 13. We find its local minimizer, $\boldsymbol{h}^*$, by using scaled conjugate gradient descent (SCG) algorithm. SCG is chosen due to its ability to handle tens of thousands of variables. SCG has also been widely used in previous approaches involving GPs [17, 7, 19]. We use chain rule to compute $\frac{\partial Z}{\partial \boldsymbol{h}}$ by evaluating first $\frac{\partial Z}{\partial \boldsymbol{K}_y}$ and then $\frac{\partial \boldsymbol{K}_y}{\partial \boldsymbol{h}}$. The matrix calculus identities from [16] are used while computing the following expressions.

$$\frac{\partial Z}{\partial \boldsymbol{K}_y} = -\frac{1}{2} \boldsymbol{K}_y^{-1} \boldsymbol{y}_m^+ \left( \boldsymbol{y}_m^+ \right)^T \boldsymbol{K}_y^{-1} + \frac{1}{2} \boldsymbol{K}_y^{-1} + \frac{1}{2} \boldsymbol{\alpha} \boldsymbol{\alpha}^T + \\ 2 \sum_i \left( \boldsymbol{K}_y^{\boldsymbol{F}_i^+} - \boldsymbol{K}_y^{\boldsymbol{F}_i^+, \boldsymbol{F}_i^-} \right),$$

$$\left[ \frac{\partial \boldsymbol{K}_y}{\partial \theta_q} \right]_{ij} = -\frac{1}{2} \sigma_f^2 \exp \left( -\frac{1}{2} (\boldsymbol{f}_i - \boldsymbol{f}_j)^T \Lambda (\boldsymbol{f}_i - \boldsymbol{f}_j) \right) \cdot \\ (\boldsymbol{f}_i^{(q)} - \boldsymbol{f}_j^{(q)})^2,$$

$$\frac{\partial \boldsymbol{K}_y}{\partial \sigma_f^2} = \sigma_f^2 \exp(-\frac{1}{2} (\boldsymbol{f}_i - \boldsymbol{f}_j)^T \cdot \Lambda \cdot (\boldsymbol{f}_i - \boldsymbol{f}_j)),$$

$$\left[ \frac{\partial \boldsymbol{K}_y}{\partial \sigma_y^2} \right]_{ij} = \delta_{ij},$$

$$\frac{\partial Z}{\partial \theta_q} = \text{tr} \left[ \left( \frac{\partial Z}{\partial \boldsymbol{K}_y} \right)^T \left( \frac{\partial \boldsymbol{K}_y}{\partial \theta_q} \right) \right] \quad \forall q \in \{1, \ldots, D\},$$

(14)

where tr denotes matrix trace. Similarly, $\frac{\partial Z}{\partial \sigma_f^2}$ and $\frac{\partial Z}{\partial \sigma_y^2}$ are computed to construct $\frac{\partial Z}{\partial \boldsymbol{h}} \in \mathbb{R}^{D+2}$. This derivative can now be used to obtain the optimal set of hyper-parameters, $\boldsymbol{h}$. In practice, all the matrix inverses are implemented using Cholesky decomposition. We alternately optimize for $\boldsymbol{\alpha}$ and $\boldsymbol{h}$ till Equation 13 converges or the maximum cycles are reached. We set the convergence criterion to be $10^{-3}$ and the maximum cycles to 20.

## 3.5. Testing

Once we get the optimal $\boldsymbol{\alpha}$ and $\boldsymbol{h}$, we can predict the parameters, $\{\bar{y}_{*1}^+, \bar{y}_{*2}^+, \bar{y}_{*3}^+\}$, for the enhanced counterpart by using three trained GP models in Equation 1. Let us call the mean and variances of the predicted parameters as $\boldsymbol{m} = \{m_1, m_2, m_3\}$ and $\boldsymbol{s} = \{s_1, s_2, s_3\}$ respectively. With their availability, we now explain our parameter space traversal.

As mentioned before, people's choices vary a lot in such applications. Thus, it is essential to explore the parameter space to generate additional enhancement parameters. The first advantage of our approach is that we can generate such parameters without referring to the training set. Since we explore the parameter space in a structured manner (with a certain mean and variance), we can afford to generate only 32 parameters per image instead of hundreds as done in conventional $k$NN-based heuristic methods.

The First step in parameter space traversal is to determine lower and upper bounds. Those can be decided heuristically. For example, we decrease the saturation, brightness and contrast at most by an amount of $\{15\%, 15\%, 5\%\}$ and increase it at most by $\{35\%, 35\%, 20\%\}$ of the original image parameter values. We observed that these limits are not absolutely critical to the quality since the generated images will be ranked later using the learned $\boldsymbol{\alpha}$ and the images with extreme parameter settings will usually be filtered out.

Now, we change (increase and decrease) the mean value of the parameters by $\mu s$ till it reaches the pre-specified thresholds. Intuitively, we think that $s$ gives us the direction of our stride in the parameter space and $\mu$ gives us the length of that stride. The value of $\mu$ is determined by the number of enhanced counterparts the user wants to generate for each low-quality image. We set that value to be 30. This value could be decreased if the user is on a mobile device with a smaller screen and similarly increased when operating on a desktop. These settings can be changed on-the-fly.

### 3.6. Image feature representation

We extract 432-D color histogram with 12 bins for hue, 6 bins each for saturation and value, which acts as a global feature. We then divide the image into a $12 \times 12$ grid. For each grid, we calculate its saturation, value by taking the mean values of those image blocks in the HSV color space. We also calculate RMS contrast on that grid. These act as localized features of 144-D each. We finally append the image parameters, which are average saturation, value and RMS contrast. Appending the image parameters allows GP to express the parameters of the enhanced counterparts as a function of both, the low-quality parameters and its feature vector. We finally get a 867-D ($= 432 + 3 \times 144 + 3$) representation for every image.

### 3.7. Implementation Details and Efficiency

GPs are known to be computationally-intensive. They take about $O(N^3)$ time for training, where $N$ are the number of training examples. The matrix inversion of an $N \times N$ matrix and the computation of the derivative of the kernel are the bottlenecks in the GP training procedure. We train a GP model using about 1200 low-quality images and six counterparts per image in about 18 hours on an Intel Xeon @2.4 GHz $\times$ 16. The computational efficiency can be improved by using GP regression techniques proposed for large data [8, 1] or using efficient data-structures such as kd-trees [22]. During testing, our approach is extremely fast. We tested it on two systems, Intel Xeon and a modern desktop system with Intel i7 @3.7GHz. It can predict all the three parameters for 3150 and 1287 images per second using Intel Xeon and i7 systems respectively. A built-in $k$NN-search function processes only 224 images per second when asked to find one nearest neighbor in 5000 image data-set on the Intel Xeon system. All the implementations are done in MATLAB. Since our approach need not query the training database, it could be portable and potentially allow for enhancements being performed on mobile devices.

## 4. Data-sets and Experimental Setup

In this section, we present describe the data and experimental setup. Results of these experiments are presented in the Section 5. We perform four kinds of experiments.

The first experiment provides a weak quantitative measure of the accuracy of our approach. We use the MIT-Adobe FiveK [5] data-set for this experiment. This data-set has 5000 low-quality images with 5 expert-enhanced counterparts for each image. This is the largest such data-set available. We use 1200 images and 6 counterparts (three each for poor and high-quality) per low-quality image to train our GP models. We use 1500 and 800 images for validation and testing respectively. We predict the parameters (i.e. brightness, contrast and saturation) for the first enhanced counterpart of all the images in the test set. Then we calculate the root mean square error (RMSE) and a more stringent criterion - Pearson's correlation - between the ground-truth parameters computed from expert-enhanced image and our predicted parameters. We compare our quantitative results against twin Gaussian processes (TGP) [3]. TGP is a strutured prediction method which considers correlation between both input and output to produce predictions. Though a low RMSE between ground truth and predicted parameters does not guarantee that the enhancement will be visually appealing (unless the RMSE tends to zero), it gives us a confirmation that the prediction is lying near the ground-truth in the parameter space. Also, this experiment validates the effectiveness of the GP regressor.

Second experiment is a qualitative measure of the image quality produced by the proposed and the competing algorithms - $k$NN, Picasa and [25]. The metric of $L2$ error in the L*ab space was adopted in [25]. We believe that it is a poor indicator of the enhancement quality and instead opt for Visual Information Fidelity (VIF) metric [21]. This metric has the ability to predict whether the visual quality of the other image has been enhanced in comparison with the reference image by producing a value greater than one. This is unlike other quality metrics such as SSIM [24], FSIM [28], VSI [27] etc. We use the publicly available implementation of VIF[5]. We calculate the VIF between the proposed enhancement (reference image) and the enhancement by 1. $k$NN 2. Picasa and 3. the approach of [25] (other image). Thus VIF $< 1$ implies that the proposed enhancement is better than the one produced by the competing algorithm and vice-versa. This comparison is done for 60 pairs where 15 images each are enhanced using Picasa and [25], whereas the remaining 30 images are enhanced using $k$NN approach.

Third experiment is aimed towards evaluating the effectiveness of GP ranking. For each image we generate only 32 enhanced versions. Our GP ranker selects the highest ranked image out of those 32 and presents it to the user. The highest ranked image is supposed to have the best quality. We compute the VIF metric between the best image selected by the ranker (reference image) and the other 31 images (other images). Ideally, for all these 31 images, we should get values less than one indicating that GP ranker

---

[5]available at `live.ece.utexas.edu/research/quality/`

has indeed selected the best image.

We also carry out a subjective evaluation test to assess if people prefer the enhanced counterparts generated by our approach. We compare our approach against three other methods. First one is the $k$NN-based approach. Given a low-quality image, we search for the nearest non-duplicate image from the 5000 images of MIT-Adobe dataset. The parameters of the expert-enhanced counterparts of the nearest image are applied to the given low-quality image. In this manner, we generate 5 enhanced counterparts per low-quality image. Note that, $k$NN utilizes all other 4999 images whereas we only use the model trained on 1200 images for prediction. Then we compare against Picasa's one-touch-enhance tool. The third approach is from [25], which also is a learning-to-rank based image enhancement approach that uses the pipeline shown at the top in Fig. 1.

We use 60 images for the subjective test which was performed by 15 people. Thirty images are selected from our testing set of the MIT-Adobe data-set. The rest of the images are from the data-set used in the paper [25]. Since we only have access to their testing set, we use that data-set solely for subjective test purposes. It contains 124 images out of which we randomly select 30 images. We enhance all the 60 images using our approach. The comparison against other methods is done as follows.

The first 30 images from the MIT-Adobe data is split into two halves. The first half is enhanced using the $k$NN approach and the second half is enhanced using Picasa. The remaining 30 images from [25] are split into two halves. The first half is enhanced using the $k$NN approach and for the second half, we directly use the high-resolution results of the test data-set of [25]. Thus each person compares 60 image pairs. One of the image in that pair has been enhanced using our approach and the other image has been enhanced using either $k$NN approach, Picasa or the approach of [25]. The subject has to choose the image which he/she finds "visually-appealing". If the subject feels that both images have almost the same visual appeal, a third option of preferring neither image is provided. The order in which the images appear in front of a subject is always randomized. The pairing order is also randomized. The subjects do the evaluation test in standard lighting conditions and at a comfortable and constant distance from the screen.

## 5. Results

We present results of our quantitative analysis first. We train three GP models to predict saturation, brightness and contrast for 800 images from the test set of MIT-Adobe data-set. When compared with the parameters of expert-enhanced counterparts, we achieve RMSEs of $0.0057, 0.0022, 0.0037$ and correlations of $0.5359, 0.5553, 0.8023$ respectively, for the above three parameters. TGP gets an average RMSE of 0.0022 but it suf-
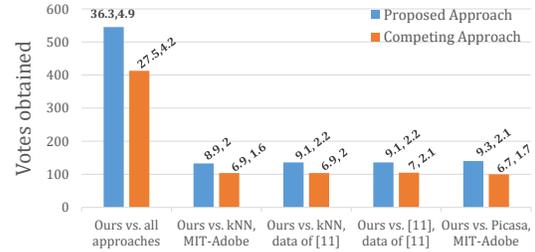


Figure 2. Subjective evaluation test metrics.

fers while producing an average correlation of only 0.3326. We can see that it is relatively easier for a GP to relate the contrast to the image quality, which is intuitive since contrast variation changes the image drastically and it also makes the image look vibrant or dull. This in turn contributes most to the visual appeal of an image.

The left bar chart in Fig. 3 shows the results of second experiment. VIF between our enhancement and competing enhancements produces values which are, in most cases, less than one. Thus according to VIF metric, our approach produces better enhancements than Picasa, $k$NN-based heuristics and [25]. For third experiment, we get 32 VIF values for each image, which correspond to 32 enhanced versions generated by our approach. The GP ranker selects one, as mentioned earlier. We compute the average VIF value and its standard deviation over 31 other images. This process is repeated for all the 60 images and the VIF values are shown in the right bar chart of Fig. 3.

We now analyze the results of our subjective tests. We provide the following five metrics about our subjective test in Fig. 2. 1. we count votes gathered by our approach and by all other competing approaches bundled into one. This is a coarse measure of how much preference people have towards enhancements generated by our approach. 2. We count votes gathered by our approach and by $k$NN approach on the MIT-Adobe data-set. 3. comparison of votes gathered by our approach and by the $k$NN approach on the data-set of [25]. 4. comparing our approach against the results of [25] on their data. 5. Lastly, we compare our approach versus Picasa on the MIT-Adobe data. Fig. 2 shows all these metrics. On top of each bar, we indicate the mean and standard deviation for that particular approach and metric. For example, the second set of bars denote that for the MIT-Adobe data, our approach gathered 133 votes against 104 votes gathered by the $k$NN approach. The average number of votes obtained per user for our and the $k$NN approach were 8.9 and 6.9 with the standard deviations of 2 and 1.6, respectively. Fig. 2 shows that people consistently prefer our approach over other state-of-art approaches.

Fig. 4 shows some of the results obtained by ours and the approach of [25], $k$NN and Picasa's auto-enhance tool. The first and the third row illustrate that the $k$NN approach is
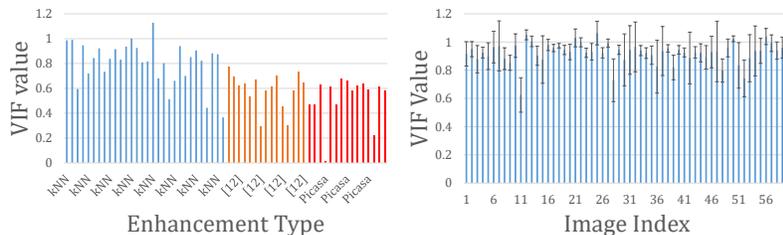
Figure 3. Left plot shows VIF values comparing proposed enhancement and enhancements produced by competing algorithms. The right plot shows the mean and standard deviation of the VIF values between the best enhancements and 31 other enhancements "rejected" by GP ranker. VIF values $< 1$ are desirable in both the cases.
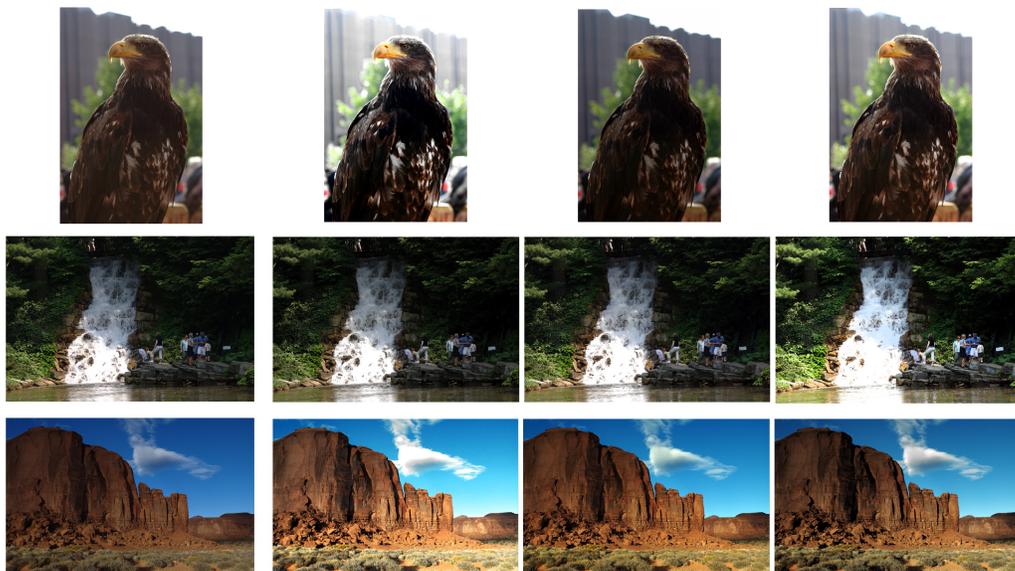


Figure 4. The left column always contains an original low-quality image. **Row 1 and 3**: Columns 2-4 contain images enhanced by $k$NN, Picasa and GP respectively. **Row 2**: The right three columns contain enhanced versions generated by GP. Please read text for details.

not always effective and sometimes may give over(under)-exposed results due to its dependence on the nearest training image parameters. The second row shows three representative versions generated by GP. We can see that the image in the fourth column is over-exposed. However, our ranking model successfully filters out that image and selects the one in the third column. In general, we observed that $k$NN can only get comparable results to Picasa and our approach if it finds a good match in the training set. Thus $k$NN is unlikely to scale to large-scale enhancement tasks.

## 6. Conclusion

We presented a novel approach to image enhancement using joint regression and ranking by employing GPs. We train our GP models on the pairs formed from poor, low and high-quality images. The learned GP models predict the desired parameters for a low-quality image from its features, which may produce its enhanced counterparts. We also described a strategy to traverse the parameter space without referring to the training images, which makes our approach efficient during testing. The GP prediction is defined by the covariance kernel, on which we impose two constraints. The first one enables the kernel to learn the feature dimensions responsible for making an image of higher-quality. The other constraint clusters all the enhancement parameters corresponding to a low-quality image, thereby allowing for effective parameter traversal. We perform quantitative and subjective evaluation experiments on two-data sets to assess the effectiveness of our approach, first one being the MIT-Adobe data [5] and the another one proposed in [25]. Quantitative experiments show that our predictions produce a low RMSE when compared with the ground-truth parameters of the MIT-Adobe data. The results show that people consistently prefer the enhancements produced by the proposed approach over the other state-of-art approaches.

# References

[1] S. Ambikasaran, D. Foreman-Mackey, L. Greengard, D. W. Hogg, and M. O'Neil. Fast direct methods for gaussian processes and the analysis of nasa kepler mission data. *arXiv preprint arXiv:1403.6015*, 2014.

[2] F. Berthouzoz, W. Li, M. Dontcheva, and M. Agrawala. A framework for content-adaptive photo manipulation macros: Application to face, landscape, and global manipulations. *ACM Trans. Graph.*, 30(5):120, 2011.

[3] L. Bo and C. Sminchisescu. Twin gaussian processes for structured prediction. *International Journal of Computer Vision*, 87(1-2):28–52, 2010.

[4] A. Bordes, S. Ertekin, J. Weston, and L. Bottou. Fast kernel classifiers with online and active learning. *The Journal of Machine Learning Research*, 6:1579–1619, 2005.

[5] V. Bychkovsky, S. Paris, E. Chan, and F. Durand. Learning photographic global tonal adjustment with a database of input/output image pairs. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 97–104. IEEE, 2011.

[6] P. S. Chandakkar, Q. Tian, and B. Li. Relative learning from web images for content-adaptive enhancement. In *Multimedia and Expo (ICME), 2015 IEEE International Conference on*, pages 1–6. IEEE, 2015.

[7] S. Eleftheriadis, O. Rudovic, and M. Pantic. Discriminative shared gaussian processes for multiview and view-invariant facial expression recognition. *Image Processing, IEEE Transactions on*, 24(1):189–204, 2015.

[8] J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. *arXiv preprint arXiv:1309.6835*, 2013.

[9] S. J. Hwang, A. Kapoor, and S. B. Kang. Context-based automatic local image enhancement. In *Computer Vision–ECCV 2012*, pages 569–582. Springer, 2012.

[10] T. Joachims. Optimizing search engines using clickthrough data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 133–142. ACM, 2002.

[11] N. Joshi, W. Matusik, E. H. Adelson, and D. J. Kriegman. Personal photo enhancement using example images. *ACM Trans. Graph*, 29(2):12, 2010.

[12] S. B. Kang, A. Kapoor, and D. Lischinski. Personalization of image enhancement. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1799–1806. IEEE, 2010.

[13] A. Kapoor, J. C. Caicedo, D. Lischinski, and S. B. Kang. Collaborative personalization of image enhancement. *International Journal of Computer Vision*, 108(1-2):148–164, 2014.

[14] L. Kaufman, D. Lischinski, and M. Werman. Content-aware automatic photo enhancement. In *Computer Graphics Forum*, volume 31, pages 2528–2540. Wiley Online Library, 2012.

[15] K. P. Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.

[16] K. B. Petersen, M. S. Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7:15, 2008.

[17] C. E. Rasmussen. Gaussian processes for machine learning. 2006.

[18] E. Reinhard, M. Stark, P. Shirley, and J. Ferwerda. Photographic tone reproduction for digital images. In *ACM Transactions on Graphics (TOG)*, volume 21, pages 267–276. ACM, 2002.

[19] O. Rudovic, M. Pantic, and I. Patras. Coupled gaussian processes for pose-invariant facial expression recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(6):1357–1369, 2013.

[20] O. Rudovic, I. Patras, and M. Pantic. Regression-based multi-view facial expression recognition. In *Pattern Recognition (ICPR), 2010 20th International Conference on*, pages 4121–4124. IEEE, 2010.

[21] H. R. Sheikh and A. C. Bovik. Image information and visual quality. *Image Processing, IEEE Transactions on*, 15(2):430–444, 2006.

[22] Y. Shen, A. Ng, and M. Seeger. Fast gaussian process regression using kd-trees. In *Proceedings of the 19th Annual Conference on Neural Information Processing Systems*, number EPFL-CONF-161316, 2006.

[23] R. Urtasun and T. Darrell. Discriminative gaussian process latent variable model for classification. In *Proceedings of the 24th international conference on Machine learning*, pages 927–934. ACM, 2007.

[24] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *Image Processing, IEEE Transactions on*, 13(4):600–612, 2004.

[25] J. Yan, S. Lin, S. B. Kang, and X. Tang. A learning-to-rank approach for image color enhancement. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 2987–2994. IEEE, 2014.

[26] Z. Yan, H. Zhang, B. Wang, S. Paris, and Y. Yu. Automatic photo adjustment using deep neural networks. *arXiv preprint arXiv:1412.7725*, 2014.

[27] L. Zhang, Y. Shen, and H. Li. Vsi: A visual saliency-induced index for perceptual image quality assessment. *Image Processing, IEEE Transactions on*, 23(10):4270–4281, 2014.

[28] L. Zhang, L. Zhang, X. Mou, and D. Zhang. Fsim: a feature similarity index for image quality assessment. *Image Processing, IEEE Transactions on*, 20(8):2378–2386, 2011.