

Image Super-Resolution via Sparse Representation

Jianchao Yang, *Student Member, IEEE*, John Wright, *Student Member, IEEE* Thomas Huang, *Life Fellow, IEEE*
and Yi Ma, *Senior Member, IEEE*

Abstract—This paper presents a new approach to single-image superresolution, based on sparse signal representation. Research on image statistics suggests that image patches can be well-represented as a sparse linear combination of elements from an appropriately chosen over-complete dictionary. Inspired by this observation, we seek a sparse representation for each patch of the low-resolution input, and then use the coefficients of this representation to generate the high-resolution output. Theoretical results from compressed sensing suggest that under mild conditions, the sparse representation can be correctly recovered from the downsampled signals. By jointly training two dictionaries for the low resolution and high resolution image patches, we can enforce the similarity of sparse representations between the low resolution and high resolution image patch pair with respect to their own dictionaries. Therefore, the sparse representation of a low resolution image patch can be applied with the high resolution image patch dictionary to generate a high resolution image patch. The learned dictionary pair is a more compact representation of the patch pairs, compared to previous approaches which simply sample a large amount of image patch pairs, reducing the computational cost substantially. The effectiveness of such a sparsity prior is demonstrated for general image super-resolution and also for the special case of face hallucination. In both cases, our algorithm generates high-resolution images that are competitive or even superior in quality to images produced by other similar SR methods, but with faster processing speed.

I. INTRODUCTION

Super-resolution image reconstruction is currently a very active area of research, as it offers the promise of overcoming some of the inherent resolution limitations of low-cost imaging sensors (e.g. cell phone or surveillance cameras) allowing better utilization of the growing capability of high-resolution displays (e.g. high-definition LCDs). Such resolution-enhancing technology may also prove to be essential in medical imaging and satellite imaging where diagnosis or analysis from low-quality images can be extremely difficult. Conventional approaches to generating a super-resolution (SR) image normally require as input *multiple* low-resolution images of the same scene, which are aligned with sub-pixel accuracy. The SR task is cast as the inverse problem of recovering the original high-resolution image by fusing the low-resolution images, based on reasonable assumptions or prior knowledge about the observation model that maps the high-resolution image to the low-resolution images. The fundamental reconstruction constraint for SR is that the recovered image, after applying the same generation model, should reproduce the observed low

resolution images. However, SR image reconstruction is generally a severely ill-posed problem because of the insufficient number of low resolution images, ill-conditioned registration and unknown blurring operators, and the solution from the reconstruction constraint is not unique. Various regularization methods have been proposed to further stabilize the inversion this ill-posed problem, such as [14], [11], [25].

However, the performance of these reconstruction-based super-resolution algorithms degrades rapidly when the desired magnification factor is large or the number of available input images is small. In these cases, the result may be overly smooth, lacking important high-frequency details [2]. Another class of SR approach is based on interpolation [29], [6], [27]. While simple interpolation methods such as bilinear or bicubic interpolation tend to generate overly smooth images with ringing and jagged artifacts, interpolation by exploiting the natural image priors will generally produce more favorable results. Dai *et al.* [6] represented the local image patches using the background/foreground descriptors and reconstructed the sharp discontinuity between the two. Sun *et al.* [27] explored the gradient profile prior for local image structures and applied it to super-resolution. Such approaches are effective in preserving the edges in the zoomed image. However, they are limited in modeling the visual complexity of the real images. For natural images with fine textures or smooth shading, these approaches tend to produce watercolor-like artifacts.

A third category of SR approach is based on machine learning techniques, which attempt to capture the co-occurrence prior between low-resolution and high-resolution image patches. [12] proposed an example-based learning strategy that applies to generic images where the low-resolution to high-resolution prediction is learned via a Markov Random Field (MRF) solved by belief propagation. [23] extends this approach by using the Primal Sketch priors to enhance blurred edges, ridges and corners. Nevertheless, the above methods typically require enormous databases of millions of high-resolution and low-resolution patch pairs, and are therefore computationally intensive. [5] adopts the philosophy of LLE [22] from manifold learning, assuming similarity between the two manifolds in the high-resolution patch space and the low-resolution patch space. Their algorithm maps the local geometry of the low-resolution patch space to the high-resolution patch space, generating high-resolution patch as a linear combination of neighbors. Using this strategy, more patch patterns can be represented using a smaller training database. However, using a fixed number K neighbors for reconstruction often results in blurring effects, due to over- or under-fitting.

While the mentioned approaches above were proposed for generic image super-resolution, specific image priors can be

Jianchao Yang and Thomas Huang are with Beckman Institute, University of Illinois Urbana-Champaign, Urbana, IL 61801 USA (email: jyang29@ifp.uiuc.edu; huang@ifp.uiuc.edu).

John Wright and Yi Ma are with CSL, University of Illinois Urbana-Champaign, Urbana, IL 61801 USA (email: jnwright@uiuc.edu; yima@uiuc.edu).

incorporated when tailored to SR applications for specific domains such as human faces. This *face hallucination* problem was addressed in the pioneering work of Baker and Kanade [30]. However, the gradient pyramid-based prediction introduced in [30] does not directly model the face prior, and the pixels are predicted individually, causing discontinuities and artifacts. Liu *et al.* [16] proposed a two-step statistical approach integrating the global PCA model and a local patch model. Although the algorithm yields good results, the holistic PCA model tends to yield results like the mean face and the probabilistic local patch model is complicated and computationally demanding. Wei Liu *et al.* [32] proposed a new approach based on TensorPatches and residue compensation. While this algorithm adds more details to the face, it also introduces more artifacts.

This paper focuses on the problem of recovering the super-resolution version of a given low-resolution image. Similar to the aforementioned learning-based methods, we will rely on patches from the input image. However, instead of working directly with the image patch pairs sampled from high resolution and low resolution images [33], we learn a compact representation for these patch pairs to capture the co-occurrence prior, significantly improving the speed of the algorithm. Our approach is motivated by recent results in sparse signal representation, which suggest that the linear relationships among high-resolution signals can be accurately recovered from their low-dimensional projections [3], [9]. Although the super-resolution problem is very ill-posed, making precise recovery impossible, the image patch sparse representation demonstrates both effectiveness and robustness in regularizing the inverse problem.

a) *Basic Ideas:* To be more precise, let $\mathbf{D} \in \mathbb{R}^{n \times K}$ be an overcomplete dictionary of K bases, and suppose a signal $\mathbf{x} \in \mathbb{R}^n$ can be represented as a sparse linear combination with respect to \mathbf{D} . That is, the signal \mathbf{x} can be written as $\mathbf{x} = \mathbf{D}\boldsymbol{\alpha}_0$ where $\boldsymbol{\alpha}_0 \in \mathbb{R}^K$ is a vector with very few ($\ll K$) nonzero entries. In practice, we might only observe a small set of measurements \mathbf{y} of \mathbf{x} :

$$\mathbf{y} \doteq \mathbf{L}\mathbf{x} = \mathbf{L}\mathbf{D}\boldsymbol{\alpha}_0, \quad (1)$$

where $\mathbf{L} \in \mathbb{R}^{k \times n}$ with $k < n$ is a projection matrix. In our super-resolution context, \mathbf{x} is a high-resolution image (patch), while \mathbf{y} is its low-resolution counterpart (or features extracted from it). If the dictionary \mathbf{D} is overcomplete, the equation $\mathbf{x} = \mathbf{D}\boldsymbol{\alpha}$ is underdetermined for the unknown coefficients $\boldsymbol{\alpha}$. The equation $\mathbf{y} = \mathbf{L}\mathbf{D}\boldsymbol{\alpha}$ is even more dramatically underdetermined. Nevertheless, under mild conditions, the sparsest solution $\boldsymbol{\alpha}_0$ to this equation will be unique. Furthermore, if \mathbf{D} satisfies an appropriate near-isometry condition, then for a wide variety of matrices \mathbf{L} , any sufficiently sparse linear representation of a high-resolution image patch \mathbf{x} in terms of the \mathbf{D} can be recovered (almost) perfectly from the low-resolution image patch [9], [21]. Figure 1 shows an example that demonstrates the capabilities of our method derived from this principle. The image of the raccoon face is blurred and downsampled to half of the original size. And then we zoom the image to the original size using our method. Even for such

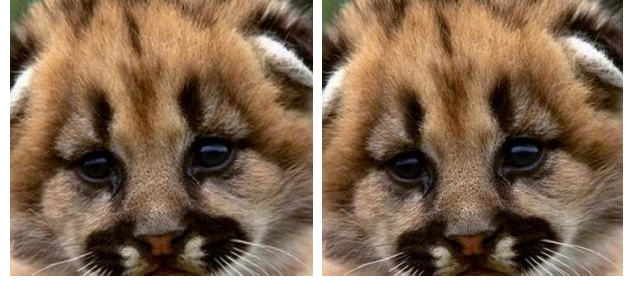


Fig. 1. Reconstruction of a raccoon face with magnification factor 2. Left: result by our method. Right: the original image. There is little noticeable difference.

a complicated texture, sparse representation recovers a visually appealing reconstruction of the original signal.

Recently sparse representation has been successfully applied to many other related inverse problems in image processing, such as denoising [10] and restoration [17], often improving on the state-of-the-art. For example in [10], the authors use the K-SVD algorithm [1] to learn an overcomplete dictionary from natural image patches and successfully apply it to the image denoising problem. In our setting, we do not directly compute the sparse representation of the high-resolution patch. Instead, we will work with two coupled dictionaries, \mathbf{D}_h for high-resolution patches, and \mathbf{D}_l for low-resolution patches. The sparse representation of a low-resolution patch in terms of \mathbf{D}_l will be directly used to recover the corresponding high-resolution patch from \mathbf{D}_h . We obtain a locally consistent solution by allowing patches to overlap and demanding that the reconstructed high-resolution patches agree on the overlapped areas. Unlike the K-SVD algorithm, we try to learn the two overcomplete dictionaries in a probabilistic model similar to [28]. To enforce that the image patch pairs have the same sparse representations with respect to \mathbf{D}_h and \mathbf{D}_l , we learn the two dictionaries simultaneously by concatenating them with normalization. The learned compact dictionaries will be applied to both generic image super-resolution and face hallucination to demonstrate their effectiveness.

Compared to the aforementioned learning-based methods, our algorithm requires only two compact learned dictionaries, instead of a large training patch database. The computation, mainly based on linear programming or convex optimization, is much more efficient and scalable, compared with [12], [23], [5]. The online recovery of the sparse representation uses the low-resolution dictionary only – the high-resolution dictionary is used only to calculate the final high-resolution image. The computed sparse representation adaptively selects the most relevant patch bases in the dictionary to best represent each patch of the given low-resolution image. This leads to superior performance, both qualitatively and quantitatively, compared to methods [5] that use a fixed number of nearest neighbors, generating sharper edges and clearer textures. In addition, the sparse representation is robust to noise as suggested in [10], and thus our algorithm is more robust to noise in the test image, while most other methods cannot perform denoising and super-resolution simultaneously.

b) Organization of the Paper: The remainder of this paper is organized as follows. Section II details our formulation and solution to the image super-resolution problem based on sparse representation. Specifically, we study how to apply sparse representation for both generic image super-resolution and face hallucination. In Section III, we discuss how to learn the two dictionaries for the high-resolution and low-resolution image patches respectively. Various experimental results in Section IV demonstrate the efficacy of sparsity as a prior for regularizing image super-resolution.

c) Notations: Specifically \mathbf{X} and \mathbf{Y} denote the high resolution and low resolution image respectively, and \mathbf{x} and \mathbf{y} denote the high resolution and low resolution image patch respectively. We use bold uppercase \mathbf{D} to denote the dictionary for sparse coding, especially we use \mathbf{D}_h and \mathbf{D}_ℓ to denote the dictionaries for high resolution and low resolution image patches respectively. Bold lowercase letters denote vectors. Unbold uppercase letters denote regular matrices, i.e., D is used as a downsampling operation in matrix form. And unbold lowercase letters are used as scalars.

II. IMAGE SUPER-RESOLUTION FROM SPARSITY

The single-image super-resolution problem asks: given a low-resolution image \mathbf{Y} , recover a higher-resolution image \mathbf{X} of the same scene. Two constraints are modeled in this work to regularize this ill-posed problem: 1) reconstruction constraint, which requires that the recovered \mathbf{X} should be consistent with the input \mathbf{Y} with respect to the image generation model; and 2) sparsity prior, which assumes that the high resolution patches can be sparsely represented in an appropriately chosen overcomplete dictionary, and that their sparse representations can be recovered from the low resolution observation.

1) *Reconstruction constraint:* The observed low-resolution image \mathbf{Y} is a blurred and downsampled version of the high resolution image \mathbf{X} :

$$\mathbf{Y} = \mathbf{D}\mathbf{H}\mathbf{X} \quad (2)$$

Here, \mathbf{H} represents a blurring filter, and \mathbf{D} the downsampling operator.

Super-resolution remains extremely ill-posed, since for a given low-resolution input \mathbf{Y} , infinitely many high-resolution images \mathbf{X} satisfy the above reconstruction constraint. We further regularize the problem via the following prior on small patches \mathbf{x} of \mathbf{X} :

2) *Sparsity prior:* The patches \mathbf{x} of the high-resolution image \mathbf{X} can be represented as a sparse linear combination in a dictionary \mathbf{D}_h trained from high-resolution patches sampled from training images:¹

$$\mathbf{x} \approx \mathbf{D}_h \boldsymbol{\alpha} \quad \text{for some } \boldsymbol{\alpha} \in \mathbb{R}^K \text{ with } \|\boldsymbol{\alpha}\|_0 \ll K. \quad (3)$$

The sparse representation $\boldsymbol{\alpha}$ will be recovered by representing patches \mathbf{y} of the input image \mathbf{Y} , with respect to a low resolution dictionary \mathbf{D}_ℓ co-trained with \mathbf{D}_h . The dictionary training process will be discussed in Sec. III.

We apply our approach to both generic images and face images. For generic image super-resolution, we divide the

problem into two steps. First, as suggested by the sparsity prior (3), we find the sparse representation for each local patch, respecting spatial compatibility between neighbors. Next, using the result from this local sparse representation, we further regularize and refine the entire image using the reconstruction constraint (2). In this strategy, a local model from the sparsity prior is used to recover lost high-frequency for local details. The global model from the reconstruction constraint is then applied to remove possible artifacts from the first step and make the image more consistent and natural. The face images differ from the generic images in that the face images have more regular structure and thus reconstruction constraints in the face subspace can be more effective. For face image super-resolution, we reverse the above two steps to make better use of the global face structure as a regularizer. We first find a suitable subspace for human faces, and apply the reconstruction constraints to recover a medium resolution image. We then recover the local details using the sparsity prior for image patches.

The remainder of this section is organized as follows: in Sec. II-A, we discuss super-resolution for generic images. We will introduce the local model based on sparse representation and global model based on reconstruction constraints. In Sec. II-B we discuss how to introduce the global face structure into this framework to achieve more accurate and visually appealing super-resolution for face images.

A. Generic Image Super-Resolution from Sparsity

1) *Local model from sparse representation:* Similar to the patch-based methods mentioned previously, our algorithm tries to infer the high-resolution image patch for each low-resolution image patch from the input. For this local model, we have two dictionaries \mathbf{D}_h and \mathbf{D}_ℓ , which are trained to have the same sparse representations for each high-resolution and low-resolution image patch pair. We subtract the mean pixel value for each patch, so that the dictionary represents image textures rather than absolute intensities.

For each input low-resolution patch \mathbf{y} , we find a sparse representation with respect to \mathbf{D}_ℓ . The corresponding high-resolution patch bases \mathbf{D}_h will be combined according to these coefficients to generate the output high-resolution patch \mathbf{x} . The problem of finding the sparsest representation of \mathbf{y} can be formulated as:

$$\min \|\boldsymbol{\alpha}\|_0 \quad \text{s.t.} \quad \|\mathbf{F}\mathbf{D}_\ell \boldsymbol{\alpha} - \mathbf{F}\mathbf{y}\|_2^2 \leq \epsilon, \quad (4)$$

where \mathbf{F} is a (linear) feature extraction operator. The main role of \mathbf{F} in (4) is to provide a perceptually meaningful constraint² on how closely the coefficients $\boldsymbol{\alpha}$ must approximate \mathbf{y} . We will discuss the choice of \mathbf{F} in Section III.

Although the optimization problem (4) is NP-hard in general, recent results [7], [8] suggest that as long as the desired coefficients $\boldsymbol{\alpha}$ are sufficiently sparse, they can be efficiently recovered by instead minimizing the ℓ^1 -norm, as follows:

$$\min \|\boldsymbol{\alpha}\|_1 \quad \text{s.t.} \quad \|\mathbf{F}\mathbf{D}_\ell \boldsymbol{\alpha} - \mathbf{F}\mathbf{y}\|_2^2 \leq \epsilon. \quad (5)$$

¹Similar mechanisms – sparse coding with an overcomplete dictionary – are also believed to be employed by the human visual system [19].

²Traditionally, one would seek the sparsest $\boldsymbol{\alpha}$ s.t. $\|\mathbf{D}_\ell \boldsymbol{\alpha} - \mathbf{y}\|_2 \leq \epsilon$. For super-resolution, it is more appropriate to replace this 2-norm with a quadratic norm $\|\cdot\|_{\mathbf{F}\mathbf{T}\mathbf{F}}$ that penalizes visually salient high-frequency errors.

Lagrange multipliers offer an equivalent formulation

$$\min \lambda \|\alpha\|_1 + \frac{1}{2} \|FD_\ell \alpha - Fy\|_2^2, \quad (6)$$

where the parameter λ balances sparsity of the solution and fidelity of the approximation to y . Notice that this is essentially a linear regression regularized with ℓ^1 -norm on the coefficients, known in statistical literature as the Lasso [24].

Solving (6) individually for each local patch does not guarantee the compatibility between adjacent patches. We enforce compatibility between adjacent patches using a one-pass algorithm similar to that of [13].³ The patches are processed in raster-scan order in the image, from left to right and top to bottom. We modify (5) so that the super-resolution reconstruction $D_h \alpha$ of patch y is constrained to closely agree with the previously computed adjacent high-resolution patches. The resulting optimization problem is

$$\min \|\alpha\|_1 \quad \text{s.t.} \quad \begin{aligned} \|FD_\ell \alpha - Fy\|_2^2 &\leq \epsilon_1, \\ \|PD_h \alpha - w\|_2^2 &\leq \epsilon_2, \end{aligned} \quad (7)$$

where the matrix P extracts the region of overlap between current target patch and previously reconstructed high-resolution image, and w contains the values of the previously reconstructed high-resolution image on the overlap. The constrained optimization (7) can be similarly reformulated as:

$$\min \lambda \|\alpha\|_1 + \frac{1}{2} \|\tilde{D}\alpha - \tilde{y}\|_2^2, \quad (8)$$

where $\tilde{D} = \begin{bmatrix} FD_\ell \\ \beta PD_h \end{bmatrix}$ and $\tilde{y} = \begin{bmatrix} Fy \\ \beta w \end{bmatrix}$. The parameter β controls the tradeoff between matching the low-resolution input and finding a high-resolution patch that is compatible with its neighbors. In all our experiments, we simply set $\beta = 1$. Given the optimal solution α^* to (8), the high-resolution patch can be reconstructed as $x = D_h \alpha^*$.

2) *Enforcing global reconstruction constraint*: Notice that (5) and (7) do not demand exact equality between the low-resolution patch y and its reconstruction $D_\ell \alpha$. Because of this, and also because of noise, the high-resolution image X_0 produced by the sparse representation approach of the previous section may not satisfy the reconstruction constraint (2) exactly. We eliminate this discrepancy by projecting X_0 onto the solution space of $DHX = Y$, computing

$$X^* = \arg \min_X \|X - X_0\| \quad \text{s.t.} \quad DHX = Y. \quad (9)$$

The solution to this optimization problem can be efficiently computed using the back-projection method, originally developed in computer tomography and applied to super-resolution in [15], [4]. The update equation for this iterative method is

$$X_{t+1} = X_t + ((Y - DHX_t) \uparrow s) * p, \quad (10)$$

where X_t is the estimate of the high-resolution image after the t -th iteration, p is a “backprojection” filter, and $\uparrow s$ denotes upsampling by a factor of s .

We take result X^* from backprojection as our final estimate of the high-resolution image. This image is as close as

³There are different ways to enforce compatibility. In [5], the values in the overlapped regions are simply averaged, which will result in blurring effects. The one-pass algorithm [13] is shown to work almost as well as the use of a full MRF model [12].

Algorithm 1 (Super-Resolution via Sparse Representation).

- 1: **Input**: training dictionaries D_h and D_ℓ , a low-resolution image Y .
- 2: **For** each 3×3 patch y of Y , taken starting from the upper-left corner with 1 pixel overlap in each direction,
 - Solve the optimization problem with \tilde{D} and \tilde{y} defined in (8): $\min \lambda \|\alpha\|_1 + \frac{1}{2} \|\tilde{D}\alpha - \tilde{y}\|_2^2$.
 - Generate the high-resolution patch $x = D_h \alpha^*$. Put the patch x into a high-resolution image X_0 .
- 3: **End**
- 4: Using back-projection, find the closest image to X_0 which satisfies the reconstruction constraint:

$$X^* = \arg \min_X \|X - X_0\| \quad \text{s.t.} \quad DHX = Y.$$

- 5: **Output**: super-resolution image X^* .
-

possible to the initial super-resolution X_0 given by sparsity, while satisfying the reconstruction constraint. The entire super-resolution process is summarized as Algorithm 1.

3) *Global optimization interpretation*: The simple SR algorithm outlined in the previous two subsections can be viewed as a special case of a more general sparse representation framework for inverse problems in image processing. Related ideas have been profitably applied in image compression, denoising [10], and restoration [17]. In addition to placing our work in a larger context, these connections suggest means of further improving the performance, at the cost of increased computational complexity.

Given sufficient computational resources, one could in principle solve for the coefficients associated with all patches *simultaneously*. Moreover, the entire high-resolution image X itself can be treated as a variable. Rather than demanding that X be perfectly reproduced by the sparse coefficients α , we can penalize the difference between X and the high-resolution image given by these coefficients, allowing solutions that are not perfectly sparse, but better satisfy the reconstruction constraints. This leads to a large optimization problem:

$$\begin{aligned} X^* = \arg \min_{X, \{\alpha_{ij}\}} & \left\{ \|DHX - Y\|_2^2 + \eta \sum_{i,j} \|\alpha_{ij}\|_0 \right. \\ & \left. + \gamma \sum_{i,j} \|D_h \alpha_{ij} - P_{ij} X\|_2^2 + \tau \rho(X) \right\}. \end{aligned} \quad (11)$$

Here, α_{ij} denotes the representation coefficients for the $(i, j)_{th}$ patch of X , and P_{ij} is a projection matrix that selects the $(i, j)_{th}$ patch from X . $\rho(X)$ is a penalty function that encodes prior knowledge about the high-resolution image. This function may depend on the image category, or may take the form of a generic regularization term (e.g., Huber MRF, Total Variation, Bilateral Total Variation).

Algorithm 1 can be interpreted as a computationally efficient approximation to (11). The sparse representation step recovers the coefficients α by approximately minimizing the sum of the second and third terms of (11). The sparsity term $\|\alpha_{ij}\|_0$ is relaxed to $\|\alpha_{ij}\|_1$, while the high-resolution fidelity term $\|D_h \alpha_{ij} - P_{ij} X\|_2$ is approximated by its low-resolution

version $\|F\mathbf{D}_\ell\boldsymbol{\alpha}_{ij} - F\mathbf{y}_{ij}\|_2$.

Notice, that if the sparse coefficients $\boldsymbol{\alpha}$ are fixed, the third term of (11) essentially penalizes the difference between the super-resolution image \mathbf{X} and the reconstruction given by the coefficients: $\sum_{i,j} \|\mathbf{D}_h\boldsymbol{\alpha}_{ij} - P_{ij}\mathbf{X}\|_2^2 \approx \|\mathbf{X}_0 - \mathbf{X}\|_2^2$. Hence, for small γ , the back-projection step of Algorithm 1 approximately minimizes the sum of the first and third terms of (11).

Algorithm 1 does not, however, incorporate any prior besides sparsity of the representation coefficients – the term $\rho(\mathbf{X})$ is absent in our approximation. In Section IV we will see that sparsity in a relevant dictionary is a strong enough prior that we can already achieve good super-resolution performance. Nevertheless, in settings where further assumptions on the high-resolution signal are available, these priors can be incorporated into the global reconstruction step of our algorithm.

B. Face super-resolution from Sparsity

Face resolution enhancement is usually desirable in many surveillance scenarios, where there is always a large distance between the camera and the objects (people) of interest. Unlike the generic image super-resolution discussed earlier, face images are more regular and thus should be easier to handle. Indeed, for face super-resolution, we can deal with lower resolution input images. The basic idea is first to use the face prior to zoom the image to a reasonable medium resolution, and then to employ the local sparsity prior model to recover details. To be more precise, the solution is also approached in two steps: 1) global model: use reconstruction constraint to recover a medium high-resolution face image, but the solution is searched only in the face subspace; and 2) local model: use the local sparse model to recover the image details.

a) Non-negative matrix factorization: In face super-resolution, the most frequently used subspace method for modeling the human face is Principal Component Analysis (PCA), which chooses a low-dimensional subspace that captures as much of the variance as possible. However, the PCA bases are holistic, and tend to generate smooth faces similar to the mean. Moreover, because principal component representations allow negative coefficients, the PCA reconstruction is often hard to interpret.

Even though faces are objects with lots of variance, they are made up of several relatively independent parts such as eyes, eyebrows, noses, mouths, checks and chins. Nonnegative Matrix Factorization (NMF) [31] seeks a representation of the given signals as an additive combination of local features. To find such a part-based subspace, NMF is formulated as the following optimization problem:

$$\begin{aligned} \arg \min_{U,V} \|Z - UV\|_2^2 \\ \text{s.t. } U \geq 0, V \geq 0, \end{aligned} \quad (12)$$

where $Z \in \mathbb{R}^{n \times m}$ is the data matrix, $U \in \mathbb{R}^{n \times r}$ is the basis matrix and $V \in \mathbb{R}^{r \times m}$ is the coefficient matrix. In our context here, A simply consists of a set of high-resolution training face images as its column vectors. The number of the bases r can

be chosen as $n * m / (n + m)$ which is smaller than n and m , meaning a more compact representation. It can be shown that a locally optimum of (12) can be obtained via the following update rules:

$$\begin{aligned} V_{ij} &\leftarrow V_{ij} \frac{(U^T Z)_{ij}}{(U^T U V)_{ij}} \\ U_{ki} &\leftarrow U_{ki} \frac{(Z V^T)_{ki}}{(U V V^T)_{ki}}, \end{aligned} \quad (13)$$

where $1 \leq i \leq r$, $1 \leq j \leq m$ and $1 \leq k \leq n$. The obtained basis matrix U is often sparse and localized.

b) Two step face super-resolution: Let \mathbf{X} and \mathbf{Y} denote the high resolution and low resolution faces respectively. \mathbf{Y} is obtained from \mathbf{X} by smoothing and downsampling as in Eq. 2, Given \mathbf{Y} , we can achieve the optimal solution for \mathbf{X} based on the Maximum *a Posteriori* (MAP) criteria,

$$\mathbf{X}^* = \arg \max_{\mathbf{X}} p(\mathbf{Y}|\mathbf{X})p(\mathbf{X}). \quad (14)$$

Using the rules in (13), we can obtain the basis matrix U , which is composed of sparse bases. Let Ω denote the face subspace spanned by U . Then in the subspace Ω , the super-resolution problem in (14) can be formulated using the reconstruction constraints as:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{D}H\mathbf{U}\mathbf{c} - \mathbf{Y}\|_2^2 + \lambda \rho(\mathbf{U}\mathbf{c}) \quad \text{s.t. } \mathbf{c} \geq 0, \quad (15)$$

where $\rho(\mathbf{U}\mathbf{c})$ is a prior term regularizing the high resolution solution, $\mathbf{c} \in \mathbb{R}^{r \times 1}$ is the coefficient vector in the subspace Ω for estimated the high resolution face, and λ is a parameter used to balance the reconstruction fidelity and the penalty of the prior term. In this paper, we simply use a generic image prior requiring that the solution be smooth. Let Γ denote a matrix performing high-pass filtering. The final formulation for (15) is:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{D}H\mathbf{U}\mathbf{c} - \mathbf{Y}\|_2^2 + \lambda \|\Gamma \mathbf{U}\mathbf{c}\|_2 \quad \text{s.t. } \mathbf{c} \geq 0. \quad (16)$$

The medium high-resolution image $\hat{\mathbf{X}}$ is approximated by $\mathbf{U}\mathbf{c}^*$. The prior term in (16) suppresses the high frequency components, resulting in over-smoothness in the solution image. We rectify this using the local patch model based on sparse representation mentioned earlier in Sec. II-A1. The complete framework of our algorithm is summarized as Algorithm 2.

III. LEARNING THE DICTIONARY PAIR

In the previous section, we discussed regularizing the super-resolution problem using sparse prior that each pair of high- and low-resolution image patches have the same sparse representations with respect to the two dictionaries \mathbf{D}_h and \mathbf{D}_ℓ . A straightforward way to obtain two such dictionaries is to sample image patch pairs directly, which preserves the correspondence between the high resolution and low resolution patch items [33]. However, such a strategy will result in large dictionaries and hence expensive computation. This section will focus on learning a more compact dictionary pair for speeding up the computation.

Algorithm 2 (Face Hallucination via Sparse Representation).

- 1: Input: sparse basis matrix U , training dictionaries \mathbf{D}_h and \mathbf{D}_l , a low-resolution image \mathbf{Y} .
 - 2: Find a smooth high-resolution face $\hat{\mathbf{X}}$ from the subspace spanned by U through:
 - Solve the optimization problem in (16):
 $\arg \min_{\mathbf{c}} \|\mathbf{M}\mathbf{U}\mathbf{c} - \mathbf{Y}\|_2 + \lambda \|\mathbf{U}\mathbf{c}\|_2 \quad \text{s.t.} \quad \mathbf{c} \geq 0.$
 - $\hat{\mathbf{X}} = \mathbf{U}\mathbf{c}^*$.
 - 3: For each patch \mathbf{y} of $\hat{\mathbf{X}}$, taken starting from the upper-left corner with 1 pixel overlap in each direction,
 - Solve the optimization problem with $\tilde{\mathbf{D}}$ and $\tilde{\mathbf{y}}$ defined in (8): $\min_{\alpha} \|\alpha\|_1 + \frac{1}{2} \|\tilde{\mathbf{D}}\alpha - \tilde{\mathbf{y}}\|_2^2.$
 - Generate the high-resolution patch $\mathbf{x} = \mathbf{D}_h\alpha^*$. Put the patch \mathbf{x} into a high-resolution image \mathbf{X}^* .
 - 4: Output: super-resolution face \mathbf{X}^* .
-

A. Single Dictionary Training

Sparse coding is the problem of learning a dictionary \mathbf{D} that can sparsely represent a given set of training examples $\mathbf{X} = \{x_1, x_2, \dots, x_t\}$. Generally, it is hard to learn a compact dictionary which guarantees that sparse representation of (4) can be recovered from ℓ_1 minimization in (5). Fortunately, many sparse coding algorithms proposed previously suffice for practical applications. In this work, we focus on the following formulation:

$$\begin{aligned} \mathbf{D} = \arg \min_{\mathbf{D}, \mathbf{Z}} \|\mathbf{X} - \mathbf{D}\mathbf{Z}\|_2^2 + \lambda \|\mathbf{Z}\|_1 \\ \text{s.t.} \quad \|\mathbf{D}_i\|_2^2 \leq 1, i = 1, 2, \dots, K. \end{aligned} \quad (17)$$

where the ℓ_1 norm $\|\mathbf{Z}\|_1$ is to enforce sparsity, and the ℓ_2 norm constraints on each columns remove the scaling ambiguity⁴. This particular formulation has been studied extensively [19], [28], [34]. (17) is not convex in both \mathbf{D} and \mathbf{Z} , but is convex in one of them with the other fixed. The basic idea is to minimize (17) alternatively over the codes \mathbf{Z} for a given dictionary \mathbf{D} , and then over \mathbf{D} for a given \mathbf{Z} , leading to a local minimum of the overall objective function.

B. Joint Dictionary Training

Given the sampled training image patch pairs $P = \{\mathbf{X}^h, \mathbf{Y}^l\}$, where $\mathbf{X}^h = \{x_1, x_2, \dots, x_n\}$ are the set of sampled high resolution image patches and $\mathbf{Y}^l = \{y_1, y_2, \dots, y_n\}$ are the corresponding low resolution image patches (or features), our goal is to learn dictionaries for high resolution and low resolution image patches, so that the sparse representation of the high resolution patch is the same as the sparse representation of the corresponding low resolution patch. This is a difficult problem, due to the ill-posed nature of super-resolution. The individual sparse coding problems in the high-resolution and low-resolution patch spaces are

$$\mathbf{D}_h = \arg \min_{\{\mathbf{D}^h, \mathbf{Z}\}} \|\mathbf{X}^h - \mathbf{D}^h\mathbf{Z}\|_2^2 + \lambda \|\mathbf{Z}\|_1, \quad (18)$$

⁴Note that without the norm constraints the cost can always be reduced by dividing \mathbf{Z} by $c > 1$ and multiplying \mathbf{D} by $c > 1$.

and

$$\mathbf{D}_l = \arg \min_{\{\mathbf{D}^l, \mathbf{Z}\}} \|\mathbf{Y}^l - \mathbf{D}^l\mathbf{Z}\|_2^2 + \lambda \|\mathbf{Z}\|_1, \quad (19)$$

respectively. We combine these objectives, forcing the high-resolution and low-resolution representations to share the same codes, instead writing

$$\begin{aligned} \min_{\{\mathbf{D}^h, \mathbf{D}^l, \mathbf{Z}\}} \frac{1}{N} \|\mathbf{X}^h - \mathbf{D}^h\mathbf{Z}\|_2^2 + \frac{1}{M} \|\mathbf{Y}^l - \mathbf{D}^l\mathbf{Z}\|_2^2 \\ + \lambda \left(\frac{1}{N} + \frac{1}{M} \right) \|\mathbf{Z}\|_1, \end{aligned} \quad (20)$$

where N and M are the dimensions of the high resolution and low resolution image patches in vector form. Here, $1/N$ and $1/M$ balance the two cost terms of (18) and (19). (20) can be rewritten as

$$\min_{\{\mathbf{D}^h, \mathbf{D}^l, \mathbf{Z}\}} \|\mathbf{X}_c - \mathbf{D}_c\mathbf{Z}\|_2^2 + \lambda \left(\frac{1}{N} + \frac{1}{M} \right) \|\mathbf{Z}\|_1, \quad (21)$$

or equivalently

$$\min_{\{\mathbf{D}^h, \mathbf{D}^l, \mathbf{Z}\}} \|\mathbf{X}_c - \mathbf{D}_c\mathbf{Z}\|_2^2 + \hat{\lambda} \|\mathbf{Z}\|_1, \quad (22)$$

where

$$\mathbf{X}_c = \begin{bmatrix} \frac{1}{\sqrt{N}} \mathbf{X}^h \\ \frac{1}{\sqrt{M}} \mathbf{Y}^l \end{bmatrix}, \quad \mathbf{D}_c = \begin{bmatrix} \frac{1}{\sqrt{N}} \mathbf{D}_h \\ \frac{1}{\sqrt{M}} \mathbf{D}_l \end{bmatrix}. \quad (23)$$

Thus we can use the same learning strategy in the single dictionary case for training the two dictionaries for our super-resolution purpose. Note that since we are using features from the low resolution image patches, \mathbf{D}^h and \mathbf{D}^l are not simply connected by a linear transform, otherwise the training process of (22) will depend on the high resolution image patches only (for detail, refer to Part III-C). Fig. 2 shows the dictionary learned by (22) for generic images.⁵ The learned dictionary demonstrates basic patterns of the image patches, such as orientated edges, instead of raw patch prototypes, due to its compactness.

C. Feature Representation for Low Resolution Image Patches

In (4), we use a feature transformation F to ensure that the computed coefficients fit the most relevant part of the low-resolution signal, and hence have a more accurate prediction for the high resolution image patch reconstruction. Typically, F is chosen as some kind of high-pass filter. This is reasonable from a perceptual viewpoint, since people are more sensitive to the high-frequency content of the image. The high-frequency components of the low-resolution image are also arguably the most important for predicting the lost high-frequency content in the target high-resolution image.

In the literature, people have suggested extracting different features for the low resolution image patch in order to boost the prediction accuracy. Freeman et al. [12] used a high-pass filter to extract the edge information from the low-resolution input patches as the feature. Sun et. al. [23] used a set of Gaussian derivative filters to extract the contours in the low-resolution patches. Chang et. al. [5] used the first-order and

⁵We omit the dictionary for the low resolution image patches because we are training on features instead the patches themselves.

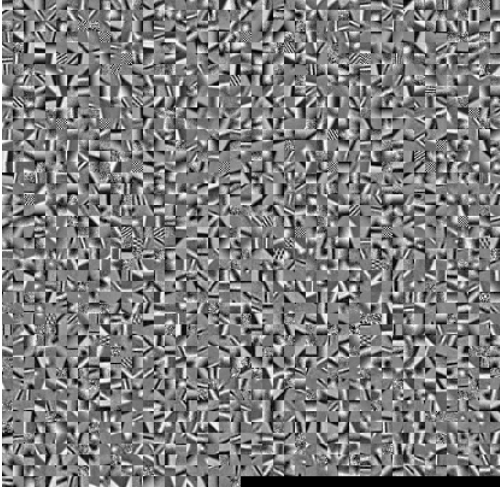


Fig. 2. The high resolution image patch dictionary trained by (22) using 100,000 high resolution and low resolution image patch pairs sampled from the generic training images.

second-order gradients of the patches as the representation. In this paper, we also use the first-order and second-order derivatives as the feature for the low-resolution patch due to their simplicity and effectiveness. The four 1-D filters used to extract the derivatives are:

$$\begin{aligned} \mathbf{f}_1 &= [-1, 0, 1], & \mathbf{f}_2 &= \mathbf{f}_1^T, \\ \mathbf{f}_3 &= [1, 0, -2, 0, 1], & \mathbf{f}_4 &= \mathbf{f}_3^T, \end{aligned} \quad (24)$$

where the superscript “ T ” means transpose. Applying these four filters yields four feature vectors for each patch, which are concatenated into one vector as the final representation of the low-resolution patch. In our implementation, the four filters are not applied directly to the sampled low resolution image patches; instead, we apply the four filters to the training images. And thus for each low resolution training image we get four gradient maps, and we sample four patches from these gradient maps at each location. Therefore, the feature representation for each low resolution image patch also encodes its neighboring information, which is beneficial for promoting compatibility among adjacent patches in the final super-resolution image.

In practice, we find that it works better to extract the features from the upsampled version of the low-resolution image instead of the original one. That is, we first upsample the low resolution image by factor of two⁶ using bicubic interpolation, and then extract gradient features from it. Since we know all the zoom ratios, it is easy to track the correspondence between high resolution image patches and the upsampled low resolution image patches both for training and testing. Because of the way of extracting features from the low resolution image patches, the two dictionaries \mathbf{D}_h and \mathbf{D}_l are not simply linearly connected, making the joint learning process in Eq. 22 more reasonable.

⁶We choose 2 mainly for dimension considerations. For example, if we work on 3-by-3 patches in the low resolution image, by upsampling the image by ratio of 2, the final feature for the 9 dimensional low resolution patch will be $6 \times 6 \times 4 = 144$.

IV. EXPERIMENTAL RESULTS

In this section, we first demonstrate the super-resolution results obtained by applying the above methods on both generic images and face images. Then, we move on to talk about the effects of dictionary size on the reconstructed image and finally we argue that our method is more robust to noise that typically exists in the input image. In our experiments, we magnify the input low resolution image by a factor of 3 for generic images and 4 for face images, which is commonplace in the literature of super-resolution. In generic image super-resolution, for the low-resolution images, we always use 3×3 low-resolution patches (upsampled to 6×6), with overlap of 1 pixel between adjacent patches, corresponding to 9×9 patches with overlap of 3 pixels for the high-resolution patches. In face super-resolution, we choose the patch size as 5×5 pixels for both low resolution and high resolution face image. For color images, we apply our algorithm to the illuminance channel only, since humans are more sensitive to illuminance changes.

A. Single Image Super-Resolution

1) *Generic image super-resolution*: We apply our methods to generic images such as flowers, human faces and architectures. The two dictionaries for high resolution and low resolution image patches are trained from 100,000 patch pairs sampled from natural images collected from the internet. We fix the dictionary size as 1024 in all our experiments, which is a balance between computation and image quality (In Sec. IV-B we will examine the effects of different dictionary sizes). In the super-resolution algorithm Eq. 8, the choice of λ depends on the level of noise in the input image, which we will discuss further in Section IV-C. For generic low-noise images, we always set $\lambda = 0.01$ in all our experiments, which generally yields satisfactory results.

Fig. 3 compares the outputs of our method with those of the neighborhood embedding method [5]. The neighborhood embedding method is similar to ours in the sense that both methods use the linear combination weights derived from the low resolution image patch to generate the underlying high resolution image patch. Unlike our method, the neighborhood embedding method uses fixed k nearest neighbors to find the reconstruction supports and does not including a dictionary training phase. To make a fair comparison, we use the same 100,000 patch pairs for the neighborhood embedding and try different k 's to get the most visually appealing results. Using a compact dictionary pair, our method is much faster and yet generates shaper results. As the reconstructed images show in Fig. 3, there are noticeable differences in the texture of the leaves, the fuzz on the leaf stalk, and also the freckles on the face of the girl.

In Figure 4, we compare our method with several more state-of-the-art methods on an image of the Parthenon used in [6], including back projection [15], neighbor embedding [5], and the recently proposed method based on a learned soft edge prior [6]. The result from back projection has many jagged effects along the edges. Neighbor embedding generates sharp edges in places, but blurs the texture on the temple's facade. The soft edge prior method gives a decent reconstruction,



Fig. 3. The flower and girl image magnified by a factor of 3. Left to right: input, bicubic interpolation, neighbor embedding [5], our method, and the original.



Fig. 4. Results on an image of the Parthenon with magnification factor 3. Top row: low-resolution input, bicubic interpolation, back projection. Bottom row: neighbor embedding [5], soft edge prior [6], and our method.

but introduces undesired smoothing that is not present in our result.

2) *Face super-resolution*: In this part, we evaluate our proposed super-resolution algorithm on frontal views of human faces. The experiments are conducted on the face database FRGC Ver 1.0 [35]. All these face images were aligned by an automatic alignment algorithm using the eye positions, and then cropped to the size of 100×100 pixels. To obtain the face subspace Ω spanned by W , we select 540 face images as training, covering both genders, different races, varying ages and different facial expressions (Figure 5). To prepare the coupled dictionaries needed for our sparse representation algorithm, we also sample 100,000 patch pairs from the training images and train the dictionary pair of size 1024. 30 new face images (from people not in the training set) are chosen as our test cases, and are blurred and downsampled to the size of 25-by-25 pixels.

As earlier mentioned, face image super-resolution can handle more challenging tasks than generic image super-resolution due to the regular face structure. Indeed, it is not an easy job



Fig. 5. Example training faces for the face super-resolution algorithm. The training images cover faces of both genders, different ages, different races and various facial expressions.

to zoom the 25×25 low resolution face image by 4 times using the method for generic image super-resolution. First, the downsampling process loses so much information that it is difficult to predict well a 12×12 high resolution patch given only a 3×3 image patch. Second, the resolution of the face image is so low that the structures of the face that are useful for super-resolution inference (such as corners and edges) collapses into only a couple of pixels. The two-step approach for face super-resolution, on the other hand, can compensate for the lost information in the first step using the redundancy of the face structures by searching the solution in the face subspace respecting the reconstruction constraints.

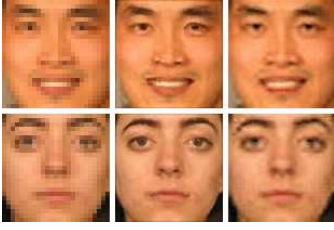


Fig. 6. The comparison between the two-step face hallucination algorithm with the generic image super-resolution algorithm applied to low resolution face images. From left to right: input image, super-resolution result using the two step approach, and super-resolution result using the generic approach.

The local model from sparse representation then can be further employed to enhance the edges and textures to achieve shaper results. We also apply the method for generic image directly to the face images, and compare the results with the proposed two-step approach, as shown in Fig. 6. Since the resolution of the input face image is so low that direct applying the generic approach does not seem to generate satisfying image.

In our experiments with face images, we also set $\lambda = 0.01$ for sparsity regularization. We compare our algorithm with bicubic interpolation [29] and back-projection [15]. The results are shown in Fig. 7, which indicate that our method can generate much higher resolution faces. From columns 4 and 5, we can also see that the local patch method based on sparse representation further enhances the edges and textures.



Fig. 7. Results of our algorithm compared to other methods. From left to right columns: low resolution input; bicubic interpolation; back projection; sparse coding via NMF followed by bilateral filtering; sparse coding via NMF and Sparse Representation; Original.

B. Effects of Dictionary Size

The above experimental results show that the sparsity prior for image patches is very effective in regularizing the otherwise ill-posed super-resolution problem. In those results, we have fixed the dictionary size to be 1024. Intuitively, larger dictionaries should have more expressive power (in the extreme, we can use the sampled patches as the dictionary directly) and thus can give more accurate approximation, while increasing the computation cost. In this section, we evaluate the effect of dictionary size on generic image super-resolution. From the sampled 100,000 image patch pairs, we train four dictionaries of size 256, 512, 1024 and 2048, and apply

Images	Bicubic	D256	D512	D1024	D2048
Girl	6.858	6.739	6.579	6.431	6.374
Flower	4.097	4.104	3.925	3.760	3.720
Lena	8.578	8.043	7.842	7.645	7.423
Statue	11.490	11.125	10.601	10.017	9.869

TABLE I
THE RMS ERRORS OF THE RECONSTRUCTED IMAGES USING
DICTIONARIES OF DIFFERENT SIZES.

them to the same input image. The results are evaluated both visually and quantitatively (in RMS errors).

Fig. 8 shows the reconstructed results for the Lena image using dictionaries of different sizes. While there is not much visual difference for the results using different dictionary sizes from 256 to 2048, we indeed observe that with larger dictionary the reconstruction artifacts visible in those with small dictionary size will gradually diminish. The visual observation is also supported by the RMS errors of the recovered images. In Table I, we list the RMS errors of the reconstructed images for dictionaries of different sizes. As shown in the table, using larger dictionaries will always yield smaller RMS error, and all of them have smaller RMS errors than those by bicubic interpolation. However, the computation is approximately linear to the size of dictionary; larger dictionaries will result in heavier computation. In practice, one chooses the dictionary size as a trade-off between reconstruction quality and computation. We find that dictionary size 1024 can yield decent outputs, while allowing fast computation.

C. Robustness to Noise

Most super-resolution algorithms assume that the input images are clean and free of noise, an assumption which is likely to be violated in real applications. To deal with noisy data, previous algorithms usually divide the recovery process into two disjoint steps: first denoising and then super-resolution. However, the results of such a strategy depend on the specific denoising techniques, and any artifacts during denoising on the low-resolution image will be kept or even magnified in the latter super-resolution process. Here we demonstrate that by formulating the problem into our sparse representation model, our method is much more robust to noise with input and thus can handle super-resolution and denoising simultaneously. Note that in (6) the parameter λ depends on the noise level of the input data; the more noisy the data, the larger the value of λ should be.

Fig. 10 shows how λ influence the reconstructed results given the same noiseless input image. The larger λ , the smoother the result image texture gets. This can be seen from the Eq. 8, which can be formulated into a probabilistic model:

$$\begin{aligned}
 \alpha^* &= \arg \min \lambda \|\alpha\|_1 + \frac{1}{2} \|\tilde{D}\alpha - \tilde{y}\|_2^2 \\
 &= \arg \min \frac{\lambda}{\sigma^2} \|\alpha\|_1 + \frac{1}{2\sigma^2} \|\tilde{D}\alpha - \tilde{y}\|_2^2 \quad (25) \\
 &= \arg \max P(\alpha) \cdot P(\tilde{y}|\alpha, \tilde{D})
 \end{aligned}$$

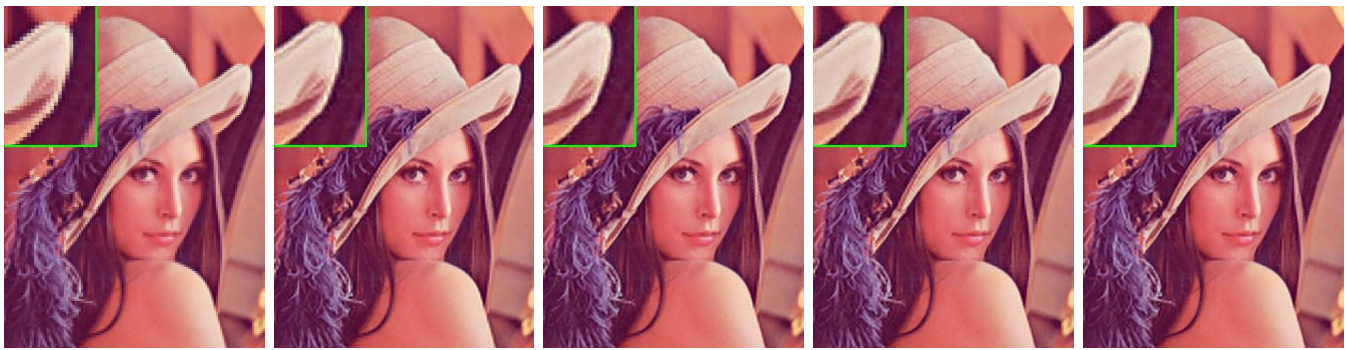


Fig. 8. The effects of dictionary size on the super-resolution reconstruction of Lena. From left to right: low resolution input, our method using dictionary size 256, 512, 1024 and 2048 respectively.

with

$$P(\alpha) = \frac{1}{2b} \exp\left(-\frac{\|\alpha\|_1}{b}\right) \quad (26)$$

$$P(\tilde{y}|\alpha, \tilde{D}) = \frac{1}{2\sigma^2} \exp\left(-\frac{1}{2\sigma^2} \|\tilde{D}\alpha - \tilde{y}\|_2^2\right)$$

where $b = \sigma^2/\lambda$ is the variance of the Laplacian prior on α , and σ^2 is the variance of the noise assumed on the data \tilde{y} . Suppose the Laplacian variance b is fixed, the more noisy of the data (σ^2 is larger), the larger of the value λ should be. On the other hand, given the input image, the larger value of λ we set, the more noisy the model will assume of the data, and thus tend to generate smoother results

To test the robustness to noise of our algorithm, we added different levels of Gaussian noise to the low resolution input image. The standard deviation of the Gaussian noise ranges from 2 to 6. The regularization parameter λ is empirically set to be one tenth of the standard deviation. In Fig. 11, we show the results of our algorithm applying to the Liberty statue image with different levels of Gaussian noise. For comparison, we also show the results of using neighborhood embedding. The number of neighbors for Neighbor Embedding [5] is chosen as decreasing as the noise becomes heavier. As shown in the figure, the neighborhood embedding method is good at preserving edges, but fails to distinguish the signal from noise, and therefore the reconstructed images are very noisy. Table 9 shows the RMS errors of the reconstructed images from different levels of noisy data. Comparison is made with both Bicubic interpolation and Neighbor Embedding. In terms of RMS error, our method outperforms Neighbor Embedding in all cases, and outperforms Bicubic interpolation if the data is not too noisy. However, when the input is too noisy, the problem is so ill-posed that super-resolution effort actually will increase the RMS error.

V. CONCLUSION

The experimental results of the previous section demonstrate the effectiveness of sparsity as a prior for patch-based super-resolution. However, one of the most important questions for future investigation is to determine, in terms of the within-category variation, the number of raw sample patches required to generate a dictionary satisfying the sparse representation prior. Tighter connections to the theory of compressed sensing

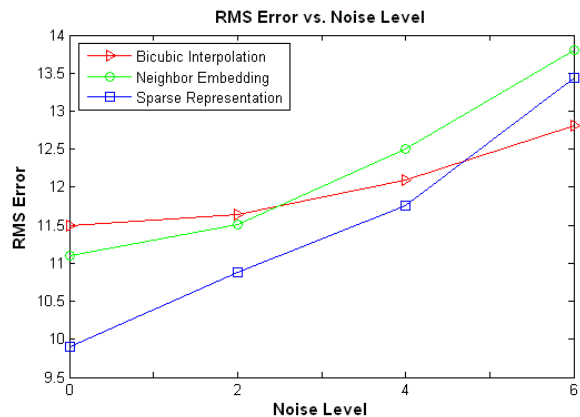


Fig. 9. The RMS errors of the reconstructed images from different levels of noisy inputs.

may also yield conditions on the appropriate patch size or feature dimension.

REFERENCES

- [1] M. Aharon, M. Elad, and A. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, Vol. 54, No. 11, November 2006.
- [2] S. Baker and T. Kanade. Limits on super-resolution and how to break them. *IEEE TPAMI*, 24(9):1167-1183, 2002.
- [3] E. Candes. Compressive sensing. *Proc. International Congress of Mathematicians*, 2006.
- [4] D. Capel. Image mosaicing and super-resolution. Ph.D. Thesis, Department of Eng. Science, University of Oxford, 2001.
- [5] H. Chang, D.-Y. Yeung, and Y. Xiong. Super-resolution through neighbor embedding. *CVPR*, 2004.
- [6] S. Dai, M. Han, W. Xu, Y. Wu, and Y. Gong. Soft edge smoothness prior for alpha channel super resolution. *Proc. ICCV*, 2007.
- [7] D. L. Donoho. For most large underdetermined systems of linear equations, the minimal ℓ^1 -norm solution is also the sparsest solution. *Comm. on Pure and Applied Math*, Vol. 59, No. 6, 2006.
- [8] D. L. Donoho. For most large underdetermined systems of linear equations, the minimal ℓ^1 -norm near-solution approximates the sparsest near-solution. Preprint, accessed at <http://www-stat.stanford.edu/~donoho/>, 2004.
- [9] D. L. Donoho. Compressed sensing. Preprint, accessed at <http://www-stat.stanford.edu/~donoho/>, 2005.
- [10] M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE TIP*, Vol. 15, No. 12, 2006.
- [11] S. Farsiu, M. D. Robinson, M. Elad, and P. Milanfar. Fast and robust multiframe super-resolution. *IEEE TIP*, 2004.
- [12] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael. Learning low-level vision. *IJCV*, 2000.



Fig. 10. The effects of λ on the recovered image given the input. From left to right, $\lambda = 0.01, 0.05, 0.1, 0.3$. The larger λ is, the smoother the result image gets. Note that the results are generated from the local model only.

- [13] W. T. Freeman, T. R. Jones, and E. C. Pasztor. Example-based super-resolution. *IEEE Computer Graphics and Applications*, Vol. 22, Issue 2, 2002.
- [14] R.C. Hardie, K.J. Barnard, and E.A. Armstrong. Joint MAP registration and high-resolution image estimation using a sequence of undersampled images. *IEEE TIP*, 1997.
- [15] M. Irani and S. Peleg. Motion analysis for image enhancement: resolution, occlusion and transparency. *JVCI*, 1993.
- [16] C. Liu, H. Y. Shum, and W. T. Freeman. Face hallucination: theory and practice. *IJCV*, Vol. 75, No. 1, pp. 115-134, October, 2007.
- [17] J. Mairal, G. Sapiro, and M. Elad. Learning multiscale sparse representations for image and video restoration. *submitted to SIAM Multiscale Modeling and Simulation*, 2007.
- [18] E. Nowak, F. Jurie, and B. Triggs. Sampling strategies for bag-of-features image classification. *Proc. ECCV*, 2006.
- [19] B. Olshausen and D. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37:3311-3325, 1997.
- [20] L. C. Pickup, S. J. Roberts, and A. Zisserman. A sampled texture prior for image super-resolution. *Proc. NIPS*, 2003.
- [21] H. Rauhut, K. Schnass, and P. Vandergheynst. Compressed sensing and redundant dictionaries. Preprint, accessed at <http://homepage.univie.ac.at/holger.rauhut/>. 2007.
- [22] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500): 2323-2326, 2000.
- [23] J. Sun, N.-N. Zheng, H. Tao, and H. Shum. Image hallucination with primal sketch priors. *Proc. CVPR*, 2003.
- [24] R. Tibshirani. Regression shrinkage and selection via the Lasso. *J. Royal Statist. Soc B.*, Vol. 58, No. 1, pages 267-288, 1996.
- [25] M. E. Tipping and C. M. Bishop. Bayesian image super-resolution. *Proc. NIPS*, 2003.
- [26] Q. Wang, X. Tang, and H. Shum. Patch based blind image super resolution. *Proc. ICCV*, 2005.
- [27] J. Sun, Z. Xu and H. Shum. Image super-resolution using gradient profile prior. *Proc. CVPR*, 2008.
- [28] Honglak Lee, Alexis Battle, Rajat Raina and Andrew Y. Ng. Efficient sparse coding algorithms. In *Proceedings of the Neural Information Processing Systems (NIPS)*, 2007.
- [29] H. S. Hou and H. C. Andrews. Cubic spline for image interpolation and digital filtering. *IEEE Trans. on SP*, , 1978.
- [30] S. Baker and T. Kanade. Hallucinating Faces. *IEEE International Conference on Automatic Face and Gesture Recognition*, March 2000.
- [31] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature* 401, 6755 (October 1999), 788-791.
- [32] Wei Liu, Dahua Lin, and Xiaoou Tang. Face Hallucinating Faces: TensorPatch Super-Resolution and Coupled Residue Compensation. *Proc. CVPR*, 2006.
- [33] Jianchao Yang, John Wright, Thomas Huang and Yi Ma. Image Super-Resolution as Sparse Representation of Raw Image Patches. *Proc. CVPR*, 2008.
- [34] J. Mutch and K. Kreutz-Delgado. Learning sparse overcomplete codes for images. *The Journal of VLSI Signal Processing*, 45:97-110, 2008.
- [35] P. Phillips, P. Flynn, T. Scruggs, K. Bowyer, J. Chang, K. Hoffman, J. Marques, J. Min and W. Worek. Overview of Face Recognition Grand Challenge. *Proc. CVPR*, 2005.

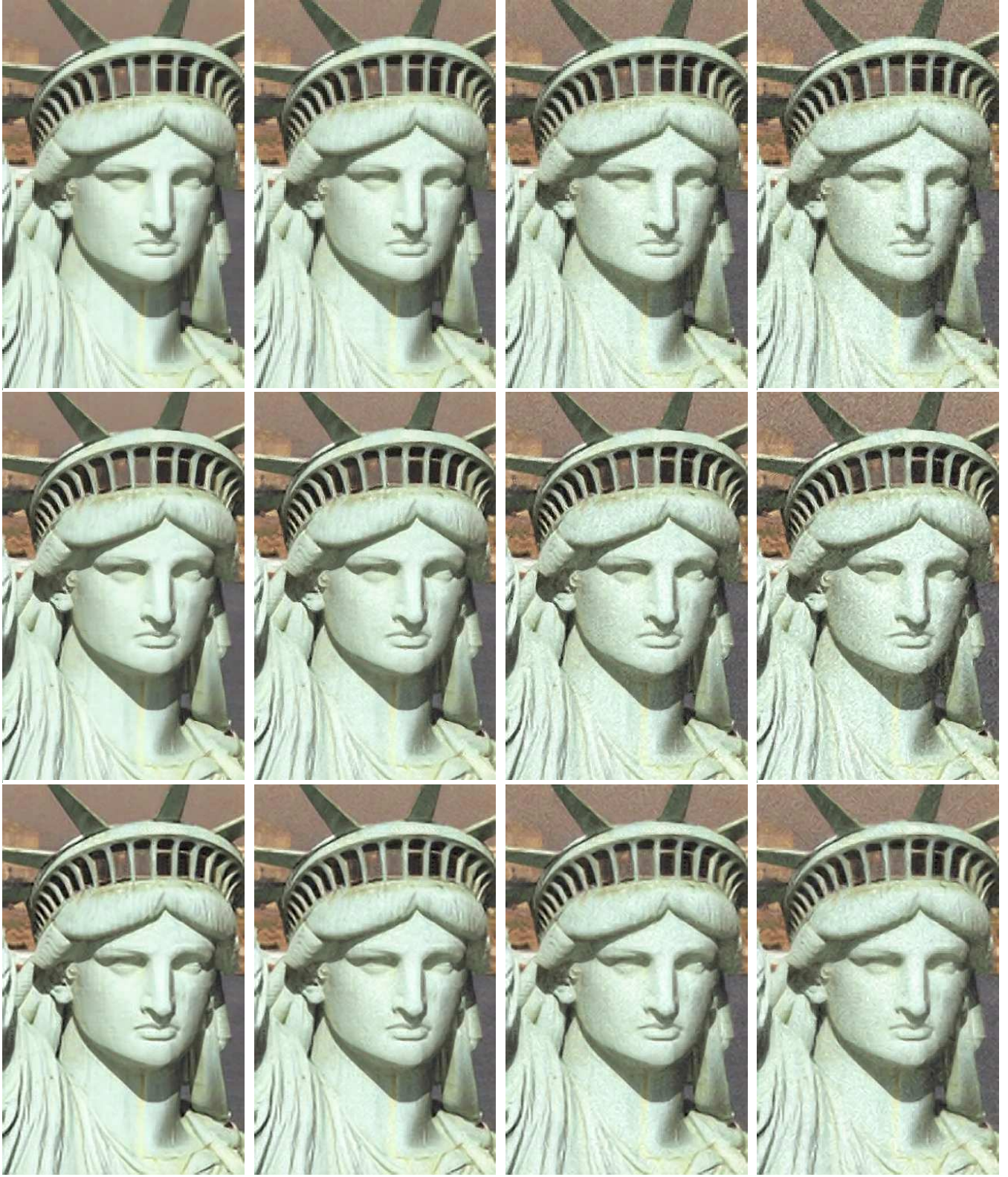


Fig. 11. Performance evaluation of our proposed algorithm on noisy data. Noise level (standard deviation of Gaussian noise) from left to right: 0, 2, 4 and 6. Top row: input images. Middle row: recovered images using neighbor embedding [5] ($k = 15, 12, 9, 7$). Bottom row: recovered images using our method ($\lambda = 0.05, 0.2, 0.4, 0.6$).