

Characterization of human faces under illumination variations using rank, integrability, and symmetry constraints*

S. Kevin Zhou, Rama Chellappa, and David W. Jacobs

Center for Automation Research
University of Maryland, College Park MD 20742, USA

Abstract. Photometric stereo algorithms use a Lambertian reflectance model with a varying albedo field and involve the appearances of only one object. This paper extends photometric stereo algorithms to handle all the appearances of all the objects in a class, in particular the class of human faces. Similarity among all facial appearances motivates a rank constraint on the albedos and surface normals in the class. This leads to a factorization of an observation matrix that consists of exemplar images of different objects under different illuminations, which is beyond what can be analyzed using bilinear analysis. Bilinear analysis requires exemplar images of different objects under same illuminations. To fully recover the *class-specific* albedos and surface normals, integrability and face symmetry constraints are employed. The proposed linear algorithm takes into account the effects of the varying albedo field by approximating the integrability terms using only the surface normals. As an application, face recognition under illumination variation is presented. The rank constraint enables an algorithm to separate the illumination source from the observed appearance and keep the illuminant-invariant information that is appropriate for recognition. Good recognition results have been obtained using the PIE dataset.

1 Introduction

Recovery of albedos and surface normals has been studied in the computer vision community for a long time. Usually a Lambertian reflectance model, ignoring both attached and cast shadows, is employed. Early works from the shape from shading (SFS) literature assume a constant albedo field: this assumption is not valid for many real objects and thus limits the practical applicability of the SFS algorithms. Early photometric stereo approaches require knowledge of lighting conditions, but full control of the lighting sources is also constraining. Recent research efforts [1, 2, 4–8] attempt to go beyond these restrictions by (i) using a varying albedo field, a more accurate model of the real world, and (ii) assuming no prior knowledge or requiring no control of the lighting sources. As a consequence, the complexity of the problem has also increased significantly.

* Partially supported by the DARPA/ONR Grant N00014-03-1-0520.

If we fix the imaging geometry and only move the lighting source to illuminate one object, the observed images (ignoring the cast and attached shadows) lie in a subspace completely determined by three images illuminated by three independent lighting sources [4]. If an ambient component is added [6], this subspace becomes 4-D. If attached shadows are considered as in [1, 2], the subspace dimension grows to infinity but most of its energy is packed in a limited number of harmonic components, thereby leading to a low-dimensional subspace approximation. However, all the photometric-stereo-type approaches (except [5]) commonly restrict themselves to using *object-specific* samples and cannot handle the appearances not belonging to the object. In this paper, we extend the photometric stereo algorithms to handle all the appearances of all the objects in a class, in particular the human face class.

To this end, we impose a rank constraint (i.e. a linear generalization) on the albedos and surface normals of all human faces. This rank constraint enables us to accomplish a factorization of the observation matrix that decomposes a *class-specific* ensemble into a product of two matrices: one encoding the albedos and surfaces normals for a class of objects and the other encoding blending linear coefficients and lighting conditions. A class-specific ensemble consists of exemplar images of different objects under different illuminations, which can not be analyzed using bilinear analysis [14]. Bilinear analysis requires exemplar images of different objects under the same illuminations. Because a factorization is always up to an invertible matrix, a full recovery of the albedos and surface normals is not a trivial task and requires additional constraints. The surface integrability constraint [9, 10] has been used in several approaches [6, 8] to successfully perform the recovery task. The symmetry constraint has also been employed in [7, 11] for face images. In Section 3, we present an approach which fuses these constraints to recover the albedos and surface normals for the face class, even in the presence of shadows. More importantly, this approach takes into account the effects of the varying albedo field by approximating the integrability terms using only the surface normals instead of the product of the albedos and the surface normals. Due to the nonlinearity embedded in the integrability terms, regular algorithms such as steepest descent are inefficient. We derive a linearized algorithm to find the solution.

In addition, the blending linear coefficients offer an illuminant-invariant identity signature, which is appropriate for face recognition under illumination variation. In Section 4, we first present a method for computing such coefficients and then report face recognition results using the PIE database [12].

1.1 Notations

In general, we denote a scalar by a , a vector by \mathbf{a} , and a matrix with r rows and c columns by $\mathbf{A}_{r \times c}$. The matrix transpose is denoted by \mathbf{A}^T , the pseudo-inverse by \mathbf{A}^\dagger . The matrix L_2 -norm is denoted by $\|\cdot\|_2$.

The following notations are introduced for the sake of notational conciseness and emphasis of special structure.

- Concatenation notations: \Rightarrow and \Downarrow .
 \Rightarrow and \Downarrow mean horizontal and vertical concatenations, respectively. For example, we can represent a $n \times 1$ vector $\mathbf{a}_{n \times 1}$ by $\mathbf{a} = [a_1, a_2, \dots, a_n]^T = [\Downarrow_{i=1}^n a_i]$ and its transpose by $\mathbf{a}^T = [a_1, a_2, \dots, a_n] = [\Rightarrow_{i=1}^n a_i]$. We can use \Rightarrow and \Downarrow to concatenate matrices to form a new matrix. For instance, given a collection of matrices $\{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n\}$ of size $r \times c$, we construct a $r \times cn$ matrix¹ $[\Rightarrow_{i=1}^n \mathbf{A}_i] = [\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n]$ and a $rn \times c$ matrix $[\Downarrow_{i=1}^n \mathbf{A}_i] = [\mathbf{A}_1^T, \mathbf{A}_2^T, \dots, \mathbf{A}_n^T]^T$. In addition, we can combine \Rightarrow and \Downarrow to achieve a concise notation. Rather than representing a matrix $\mathbf{A}_{r \times c}$ as $[a_{ij}]$, we represent it as $\mathbf{A}_{r \times c} = [\Downarrow_{i=1}^r [\Rightarrow_{j=1}^c a_{ij}]] = [\Rightarrow_{j=1}^c [\Downarrow_{i=1}^r a_{ij}]]$. Also we can easily construct ‘big’ matrices using ‘small’ matrices $\{\mathbf{A}_{11}, \mathbf{A}_{12}, \dots, \mathbf{A}_{1n}, \dots, \mathbf{A}_{mn}\}$ of size $r \times c$. The matrix $[\Downarrow_{i=1}^m [\Rightarrow_{j=1}^n \mathbf{A}_{ij}]]$ is of size $rm \times cn$, the matrix $[\Rightarrow_{i=1}^m [\Rightarrow_{j=1}^n \mathbf{A}_{ij}]]$ of size $r \times cmn$.
- Kronecker (tensor) product: \otimes .
 It is defined as $\mathbf{A}_{m \times n} \otimes \mathbf{B}_{r \times c} = [\Downarrow_{i=1}^m [\Rightarrow_{j=1}^n a_{ij} \mathbf{B}]]_{mr \times nc}$.

2 Setting and constraints

We assume a Lambertian imaging model with a varying albedo field and no shadows. A pixel h is represented as

$$h = p \mathbf{n}^T \mathbf{s} = \mathbf{t}^T \mathbf{s}, \quad (1)$$

where p is the albedo at the pixel, $\mathbf{n} \doteq [\hat{a}, \hat{b}, \hat{c}]^T$ is the unit surface normal vector at the pixel, $\mathbf{t}_{3 \times 1} \doteq [a \doteq p\hat{a}, b \doteq p\hat{b}, c \doteq p\hat{c}]^T$ is the product of albedo and surface normal, and \mathbf{s} specifies a distant illuminant (a 3×1 unit vector multiplied by the light’s intensity).

For an image \mathbf{h} , a collection of d pixels $\{h_i, i = 1, \dots, d\}$ ², by stacking all the pixels into a column vector \mathbf{h} , we have

$$\mathbf{h}_{d \times 1} \doteq [\Downarrow_i h_i] = [\Downarrow_i \mathbf{t}_i^T] \mathbf{s} = \mathbf{T}_{d \times 3} \mathbf{s}_{3 \times 1}, \quad (2)$$

where $\mathbf{T} \doteq [\Downarrow_i \mathbf{t}_i^T]$ contains all albedo and surface normal information about the object. We call the \mathbf{T} matrix the *object-specific albedo-shape* matrix.

In the case of photometric stereo, we have n images of the *same* object, say $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n\}$, observed at a fixed pose illuminated by n different lighting sources, forming an *object-specific* ensemble. Simple algebraic manipulation gives:

$$\mathbf{H}_{d \times n} \doteq [\Rightarrow_i \mathbf{h}_i] = \mathbf{T} [\Rightarrow_i \mathbf{s}_i] = \mathbf{T}_{d \times 3} \mathbf{S}_{3 \times n}, \quad (3)$$

¹ We do not need the size of $\{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n\}$ to be exactly same. We use matrices of the same size in this example for simplicity. For example, for $[\Rightarrow_{i=1}^n \mathbf{A}_i]$, we only need the number of rows of these matrices to be same.

² The index i corresponds to a spatial position $\mathbf{x} \doteq (x, y)$. If no confusion, we will interchange both notations. For instance, we might also use $\mathbf{x} = 1, \dots, d$.

where \mathbf{H} is the *observation matrix* and $\mathbf{S} \doteq [\Rightarrow_i \mathbf{s}_i]$ encodes the information about the illuminants. Hence photometric stereo is rank-3 constrained. Therefore, given at least three exemplar images for one object under three different independent illuminations, we can determine the identity of a new probe image by checking if it lies in the linear span of the three exemplar images [4]. This requires obtaining at least three images for each object in the gallery set, which may not be possible in many applications. Note that in this recognition setting, there is no need for the training set that is defined below; in other words, the training set is equivalent to the gallery set.

We follow [13] in defining a typical recognition protocol for face recognition algorithms. Three sets are needed: Gallery, probe, and training sets. The gallery set consists of images with known identities. The probe set consists of images whose identities are to be determined by matching with the gallery set. In addition, the training set is provided for the recognition algorithm to learn characteristic features of the face images. In general, we assume no identity overlap between the gallery set and the training set and often store only one exemplar image for each object in the gallery set. However, the training set can have more than one image for each object. In order to generalize from the training set to the gallery and probe sets, we note that all images in the training, gallery, and probe sets belong to the same face class, which naturally leads to a rank constraint.

2.1 The rank constraint

We impose the rank constraint on the \mathbf{T} matrix by assuming that any \mathbf{T} matrix is a linear combination of some basis matrices $\{\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_m\}$, i.e., there exist coefficients $\{f_j; j = 1, \dots, m\}$ such that

$$\mathbf{T}_{d \times 3} = \sum_{j=1}^m f_j \mathbf{T}_j = [\Rightarrow_j \mathbf{T}_j](\mathbf{f} \otimes \mathbf{I}_3) = \mathbf{W}_{d \times 3m}(\mathbf{f}_{m \times 1} \otimes \mathbf{I}_3), \quad (4)$$

where $\mathbf{f} \doteq [\Downarrow_j f_j]$, $\mathbf{W} \doteq [\Rightarrow_j \mathbf{T}_j]$, and \mathbf{I}_n denotes an identity matrix of dimension $n \times n$. Since the \mathbf{W} matrix encodes all albedos and surface normals for a class of objects, we call it a *class-specific albedo-shape* matrix. Similar rank constraints are widely found in the literature; see for example [20, 21].

Substitution of (4) into (2) yields

$$\mathbf{h}_{d \times 1} = \mathbf{T}\mathbf{s} = \mathbf{W}(\mathbf{f} \otimes \mathbf{I}_3)\mathbf{s} = \mathbf{W}(\mathbf{f} \otimes \mathbf{s}) = \mathbf{W}_{d \times 3m} \mathbf{k}_{3m \times 1}, \quad (5)$$

where $\mathbf{k} \doteq \mathbf{f} \otimes \mathbf{s}$. This leads to a two-factor bilinear analysis [14]. Recently, a multilinear analysis has been proposed in [20, 22].

With the availability of n images $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n\}$ for *different* objects, observed at a fixed pose illuminated by n different lighting sources, forming a *class-specific* ensemble, we have

$$\mathbf{H}_{d \times n} = [\Rightarrow_i \mathbf{h}_i] = \mathbf{W}[\Rightarrow_i (\mathbf{f}_i \otimes \mathbf{s}_i)] = \mathbf{W}[\Rightarrow_i \mathbf{k}_i] = \mathbf{W}_{d \times 3m} \mathbf{K}_{3m \times n}, \quad (6)$$

where $\mathbf{K} \doteq [\Rightarrow_i (\mathbf{f}_i \otimes \mathbf{s}_i)] = [\Rightarrow_i \mathbf{k}_i]$. It is a rank- $3m$ problem. Notice that \mathbf{K} takes a special form.

The rank constraint generalizes many approaches in the literature and is quite easily satisfied. If $m = 1$, this reduces to the case of photometric stereo; if the surface normal is fixed and the albedo field lies in a rank- m linear subspace, we have (4) satisfied too. Interestingly, the ‘Eigenface’ approach [15] is just a special case of our approach, but this is only for a fixed illumination source. Suppose that the illuminant vector is $\tilde{\mathbf{s}}$. (5) and (6) reduce to equations:

$$\mathbf{h}_{d \times 1} = \mathbf{W}(\mathbf{f} \otimes \tilde{\mathbf{s}}) = \tilde{\mathbf{W}}_{d \times m} \mathbf{f}_{m \times 1}; \quad \mathbf{H}_{d \times n} = [\Rightarrow_i \mathbf{h}_i] = \tilde{\mathbf{W}}[\Rightarrow_i \mathbf{f}_i] = \tilde{\mathbf{W}}_{d \times m} \mathbf{F}_{m \times n}, \quad (7)$$

where $\tilde{\mathbf{W}} \doteq [\Rightarrow_i \mathbf{T}_i \tilde{\mathbf{s}}]$. Therefore, our approach can be regarded as a generalized ‘Eigenface’ analysis able to handle illumination variation.

Our immediate goal is to estimate \mathbf{W} and \mathbf{K} from the observation matrix \mathbf{H} . The first step is to invoke an SVD factorization, say $\mathbf{H} = \mathbf{U}\mathbf{A}\mathbf{V}^T$, and retain the top $3m$ components as $\mathbf{H} = \mathbf{U}_{3m}\mathbf{A}_{3m}\mathbf{V}_{3m}^T = \hat{\mathbf{W}}\hat{\mathbf{K}}$, where $\hat{\mathbf{W}} = \mathbf{U}_{3m}$ and $\hat{\mathbf{K}} = \mathbf{A}_{3m}\mathbf{V}_{3m}^T$. Thus, we can recover \mathbf{W} and \mathbf{K} up to a $3m \times 3m$ invertible matrix \mathbf{R} with $\mathbf{W} = \hat{\mathbf{W}}\mathbf{R}$, $\mathbf{K} = \mathbf{R}^{-1}\hat{\mathbf{K}}$. Additional constraints are required to determine the \mathbf{R} matrix. We use the integrability and symmetry constraints, both related to \mathbf{W} . Moreover, $\mathbf{K} = [\Rightarrow_i (\mathbf{f}_i \otimes \mathbf{s}_i)]$.

2.2 The integrability constraint

One common constraint used in SFS research is the integrability of the surface [9, 10, 6, 8]. Suppose that the surface function is $z = z(\mathbf{x})$ with $\mathbf{x} \doteq (x, y)$, we must have $\frac{\partial}{\partial x} \frac{\partial z}{\partial y} = \frac{\partial}{\partial y} \frac{\partial z}{\partial x}$. If given the unit surface normal vector $\mathbf{n}_{(\mathbf{x})} \doteq [\hat{a}_{(\mathbf{x})}, \hat{b}_{(\mathbf{x})}, \hat{c}_{(\mathbf{x})}]^T$ at pixel \mathbf{x} , we have $\frac{\partial}{\partial x} \frac{\hat{b}_{(\mathbf{x})}}{\hat{c}_{(\mathbf{x})}} = \frac{\partial}{\partial y} \frac{\hat{a}_{(\mathbf{x})}}{\hat{c}_{(\mathbf{x})}}$. In other words, with $\alpha_{(\mathbf{x})}$ defined as an integrability constraint term,

$$\alpha_{(\mathbf{x})} \doteq \hat{c}_{(\mathbf{x})} \frac{\partial \hat{b}_{(\mathbf{x})}}{\partial x} - \hat{b}_{(\mathbf{x})} \frac{\partial \hat{c}_{(\mathbf{x})}}{\partial x} + \hat{a}_{(\mathbf{x})} \frac{\partial \hat{c}_{(\mathbf{x})}}{\partial y} - \hat{c}_{(\mathbf{x})} \frac{\partial \hat{a}_{(\mathbf{x})}}{\partial y} = 0. \quad (8)$$

If instead are given the product of albedo and surface normal $\mathbf{t}_{(\mathbf{x})} \doteq [a_{(\mathbf{x})}, b_{(\mathbf{x})}, c_{(\mathbf{x})}]^T$ with $a_{(\mathbf{x})} \doteq p_{(\mathbf{x})}\hat{a}_{(\mathbf{x})}$, $b_{(\mathbf{x})} \doteq p_{(\mathbf{x})}\hat{b}_{(\mathbf{x})}$, and $c_{(\mathbf{x})} \doteq p_{(\mathbf{x})}\hat{c}_{(\mathbf{x})}$, (8) still holds with \hat{a} , \hat{b} , and \hat{c} replaced by a , b , and c , respectively. Practical algorithms approximate the partial derivatives by forward or backward differences or other differences that use an inherent smoothness assumption. Hence, the approximations based on $\mathbf{t}_{(\mathbf{x})}$ are very rough especially at places where abrupt albedo variation exists (e.g. the boundaries of eyes, iris, eyebrow, etc) since the smoothness assumption is seriously violated. We should by all means use $\mathbf{n}_{(\mathbf{x})}$ in order to remove this effect.

2.3 The symmetry constraint

For a face image in a frontal view, one natural constraint is its symmetry about the central y -axis as proposed in [7, 11]:

$$p(x,y) = p(-x,y); \hat{a}(x,y) = -\hat{a}(-x,y); \hat{b}(x,y) = \hat{b}(-x,y); \hat{c}(x,y) = \hat{c}(-x,y), \quad (9)$$

which is equivalent to, using $\mathbf{x} \doteq (x, y)$ and its symmetric point $\bar{\mathbf{x}} \doteq (-x, y)$,

$$a(\mathbf{x}) = -a(\bar{\mathbf{x}}); b(\mathbf{x}) = b(\bar{\mathbf{x}}); c(\mathbf{x}) = c(\bar{\mathbf{x}}). \quad (10)$$

If a face image is in a non-frontal view, such a symmetry still exists but the coordinate system should be modified to take into account the view change.

3 The recovery of albedos and surface normals

The recovery task is to find from the observation matrix \mathbf{H} the *class-specific albedo-shape* matrix \mathbf{W} (or equivalently \mathbf{R}), which satisfies both the integrability and symmetry constraints, as well as the matrices \mathbf{F} and \mathbf{S} . Denote $\mathbf{R} \doteq [\Rightarrow_{j=1}^m [r_{aj}, r_{bj}, r_{cj}]]$ and $\hat{\mathbf{W}} \doteq [\Downarrow_{\mathbf{x}=1}^d \hat{\mathbf{w}}_{(\mathbf{x})}^T]$. As $\mathbf{W} \doteq [\Downarrow_{\mathbf{x}=1}^d [\Rightarrow_{j=1}^m [a_j(\mathbf{x}), b_j(\mathbf{x}), c_j(\mathbf{x})]]] = \hat{\mathbf{W}}\mathbf{R}$, we have

$$a_j(\mathbf{x}) = \hat{\mathbf{w}}_{(\mathbf{x})}^T \mathbf{r}_{aj}, \quad b_j(\mathbf{x}) = \hat{\mathbf{w}}_{(\mathbf{x})}^T \mathbf{r}_{bj}, \quad c_j(\mathbf{x}) = \hat{\mathbf{w}}_{(\mathbf{x})}^T \mathbf{r}_{cj}; \quad j = 1, \dots, m. \quad (11)$$

Practical systems must take into account attached and cast shadows as well as sensor noise. The existence of shadows in principle increases the rank to infinity. But, if we exclude shadowed pixels or set them as missing values, we still have rank 3. Performing a SVD with missing values is discussed in [3, 16].

In view of the above circumstances, we formulate the following optimization problem: minimizing over \mathbf{R} , \mathbf{F} , and \mathbf{S} the cost function \mathcal{E} defined as

$$\begin{aligned} \mathcal{E}(\mathbf{R}, \mathbf{F}, \mathbf{S}) &= \frac{1}{2} \sum_{i=1}^n \sum_{\mathbf{x}=1}^d \mathbf{i}_i(\mathbf{x}) \{h_i(\mathbf{x}) - \hat{\mathbf{w}}(\mathbf{x})^T \mathbf{R}(\mathbf{f}_i \otimes \mathbf{s}_i)\}^2 \\ &\quad + \frac{\lambda_1}{2} \sum_{j=1}^m \sum_{\mathbf{x}=1}^d \{\alpha_j(\mathbf{x})\}^2 + \frac{\lambda_2}{2} \sum_{j=1}^m \sum_{\mathbf{x}=1}^d \{\beta_j(\mathbf{x})\}^2, \\ &= \mathcal{E}_0(\mathbf{R}, \mathbf{F}, \mathbf{S}) + \lambda_1 \mathcal{E}_1(\mathbf{R}) + \lambda_2 \mathcal{E}_2(\mathbf{R}), \end{aligned} \quad (12)$$

where $\mathbf{i}_i(\mathbf{x})$ is an indicator function describing whether the pixel \mathbf{x} of the image \mathbf{h}_i is in shadow, $\alpha_j(\mathbf{x})$ is the integrability constraint term based only on surface normals as defined in (8), and $\beta_j(\mathbf{x})$ is the symmetry constraint term given as

$$\beta_j^2(\mathbf{x}) = \{a_j(\mathbf{x}) + a_j(\bar{\mathbf{x}})\}^2 + \{b_j(\mathbf{x}) - b_j(\bar{\mathbf{x}})\}^2 + \{c_j(\mathbf{x}) - c_j(\bar{\mathbf{x}})\}^2; \quad j = 1, \dots, m. \quad (13)$$

One approach could be to directly minimize the cost function over \mathbf{W} , \mathbf{F} , and \mathbf{S} . This is in principle possible but numerically difficult as the number of unknowns depends on the image size, which can be quite large in practice.

As shown in [17], the surface normals can be recovered up to a generalized bas-relief (GBR) ambiguity. To resolve the GBR ambiguity, we normalize the matrix R by keeping $\|R\|_2 = 1$. Another ambiguity between f_j and s_j is a nonzero scale, which can be removed by normalizing f : $f_j^T \mathbf{1} = 1$, where $\mathbf{1}_{m \times 1}$ is a vector of 1's.

To summarize, we perform the following task:

$$\min_{R, F, S} \mathcal{E}(R, F, S) \quad \text{subject to } \|R\|_2 = 1, F^T \mathbf{1} = 1. \quad (14)$$

An iterative algorithm can be designed to solve (14). While solving for F and S with R fixed is quite easy, solving for R given F and S is very difficult because the integrability constraint terms require partial derivatives of the surface normals that are nonlinear in R . Regular algorithms such as steepest descent are inefficient. One main contribution in this paper is that we propose a linearized algorithm to solve for R . Appendix-I presents the details of the complete algorithm.

To demonstrate how the algorithm works, we design the following scenario with $m = 2$ so that the rank of interest is $2 \times 3 = 6$. To defeat the photometric stereo algorithm, which requires one object illuminated by at least three sources, and the bilinear analysis which requires two fixed objects illuminated by at least the same three lighting sources, we synthesize eight images by taking random linear combinations of two basis objects illuminated by eight different lighting sources. Fig. 1 displays the synthesized images and the recovered class-specific albedo-shape matrix, which clearly shows the two basis objects. Our algorithm converges within 100 iterations.

One notes that the special case $m = 1$ of our algorithm can be readily applied to photometric stereo (with the symmetry constraint removed) to robustly recover the albedos and surface normals for one object.

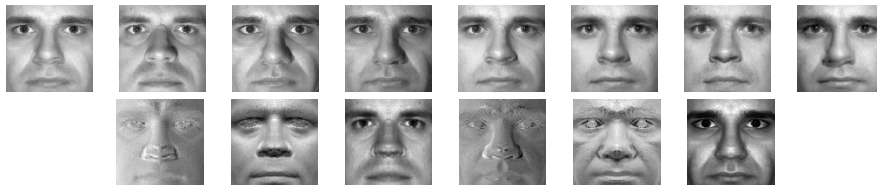


Fig. 1. Row 1: Eight Synthesized images that are random linear combinations of two basis objects illuminated by eight different lighting sources. Row 2: Recovered class-specific albedo-shape matrix showing the two basis objects (i.e. the three columns of T_1 and T_2).

4 Recognition experiments

We study an extreme recognition setting with the following features: there is no identity overlap between the training set and the gallery and probe sets; only

one image for one object is stored in the gallery set; the lighting conditions for the training, gallery and probe sets are completely unknown. The only known fact is that the face is in frontal view. Our strategy is to: (i) Learn W from the training set using the recovery algorithm described in Section 3; (ii) With W given, learn the identity signature f 's for both the gallery and probe sets using the recovery algorithm described in Section 4.1 and no knowledge of illumination directions; and (iii) Perform recognition using the nearest correlation coefficient. Suppose that one gallery image g has its signature f_g and one probe image p has its signature f_p , their correlation coefficient is $cc(p, g) = (f_p, f_g) / \sqrt{(f_p, f_p)(f_g, f_g)}$, where (x, y) is an inner-product such as $(x, y) = x^T \Sigma y$ with Σ learned or given.

4.1 Separating illumination

With the class-specific albedo-shape matrix W available, we proceed to solve the problem of separating illumination, *v.i.z.*, for an arbitrary image h , find the illuminant vector s and the identity signature f . For convenience of recognition, we normalize f to the same range: $f^T \mathbf{1} = 1$. Appendix-II presents the recovery algorithm which infers the shadow pixels as well.

4.2 PIE dataset

We use the Pose and Illumination and Expression (PIE) dataset [12] in our experiment. Fig. 2 shows the distribution of all 21 flashes used in PIE and their estimated positions using our algorithm. Since the flashes are almost symmetrically distributed about the head position, we only use 12 of them distributed on the right half of the unit sphere in Fig. 2. In total, we used $68 \times 12 = 816$ images in a fixed view as there are 68 subjects in the PIE database. Fig. 2 also displays one PIE object under the selected 12 illuminants. Registration is performed by aligning the eyes and mouth to canonical positions. No flow computation is carried out for further alignment. We use cropped face regions. After the pre-processing step, the actual image size is 50×50 . Also, we only study gray images by taking the average of the red, green, and blue channels. We use all 68 images under one illumination to form a gallery set and under another illumination to form a probe set. The training set is taken from sources other than the PIE dataset. Thus, we have 132 tests, with each test giving rise to a recognition score.

4.3 Recognition performance

The training set is first taken from the Yale illumination database [8]. There are only 10 subjects (i.e. $m = 10$) in this database and each subject has 64 images in frontal view illuminated by 64 different lights. We only use 9 lights and Fig. 2 shows one Yale object under 9 lights.

Table 1 lists the recognition rate for all test protocols for the PIE database using the Yale database as the training set. Even with $m = 10$, we obtain quite good results. One observation is that when the flashes become separated, the

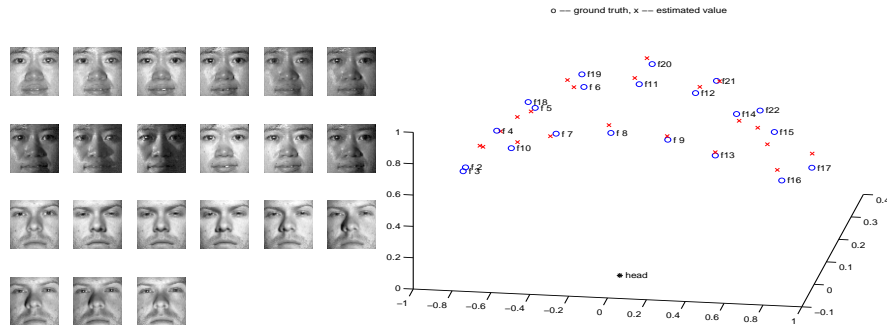


Fig. 2. Left: Rows 1-2 display one PIE object under the selected 12 illuminants and Rows 3-4 one Yale object under 9 lights used in the training set. Right: Flash distribution in the PIE database. For illustrative purposes, we move their positions on a unit sphere as only the illuminant directions matter. ‘o’ means the ground truth and ‘x’ the estimated values. It is quite accurate for estimation of directions of flashes near frontal pose. But when the flashes are very off-frontal, accuracy goes down slightly.

Gallery	f08	f09	f11	f12	f13	f14	f15	f16	f17	f20	f21	f22	Average
Probe													
f08	-	96/F	96/99	87/99	66/97	60/97	46/79	29/72	22/43	85/99	78/97	53/93	65/88
f09	94/F	-	96/99	96/99	90/99	87/99	56/97	40/91	24/60	84/97	96/97	68/97	75/94
f11	94/99	91/99	-	97/F	72/F	72/F	38/90	28/76	16/65	F/F	94/F	51/99	69/93
f12	88/99	94/99	97/F	-	88/F	93/F	57/F	41/93	28/76	94/F	F/F	76/F	78/97
f13	56/99	87/99	59/F	85/F	-	F/F	90/F	71/F	50/88	54/99	87/F	F/F	76/99
f14	51/99	85/99	63/F	93/F	F/F	-	90/F	66/F	49/96	59/99	91/F	99/F	77/99
f15	33/84	40/94	37/93	49/F	85/F	88/F	-	93/F	78/F	32/88	49/F	97/F	62/96
f16	19/69	26/87	26/78	32/90	59/F	44/F	84/F	-	93/F	26/69	31/F	63/F	46/89
f17	14/44	28/60	19/51	26/71	50/84	41/91	68/99	94/F	-	19/56	26/75	44/94	39/75
f20	90/97	85/97	99/F	97/F	65/F	69/F	38/90	26/74	21/68	-	93/F	53/F	67/93
f21	79/97	94/97	93/F	F/F	88/F	94/F	62/F	49/97	28/82	91/F	-	76/F	78/98
f22	43/90	65/97	46/96	75/F	99/F	99/F	97/F	76/F	59/99	43/97	74/F	-	70/98
Average	60/89	72/93	66/92	76/96	78/98	77/99	66/96	56/91	42/80	63/91	74/96	71/98	67/93

Table 1. Recognition rate obtained by our rank constrained approach using the Yale database (the left number in each cell) and Vetter database (the right number in each cell) as the training set. ‘F’ means 100 and ‘fnn’ flash number *nn*.

recognition rate decreases. Also, using images under frontal or near-frontal illuminants as galleries produces good results. For comparison, we also implemented the ‘Eigenface’ approach [15] by training the projection directions from the same training set. Its average recognition rate is only 35% while ours is 67%.

Generalization capacity with $m = 10$ is rather restrictive. We now increase m from 10 to 100 by using the Vetter’s 3D face database [18]. As this is a 3D database, we actually have W available. However, we believe that using a training set of $m = 100$ from other sources can yield similar performance. Table 1 presents the recognition rates. Significant improvements have been achieved due to the increase in m . The average rate is 93%. This seems to suggest that a moderate rank of 100 is enough to span the entire face space under a fixed view.

As a comparison, Romdhani et. al. [18] reported the recognition rates only with ‘f12’ being the gallery set and their average is 98% while ours is 96%. Our approach is very different from [18]. In [18] depths and texture maps of explicit 3D face models are used, while our approach is image-based and recovers

albedos and surface normals. 3D models can be then be reconstructed. In the experiments, (i) we use the ‘illum’ part of the PIE database that is close to the Lambertian model and they use the ‘light’ part that includes an ambient light; (ii) we use gray-valued images and they use color images; (iii) we assume known pose but unknown illumination but they assume unknown pose but known illumination; and (iv) compared to [18], our alignment is rather crude and can be improved using flow computations. We believe that our recognition rates can be boosted using the color images and finer alignment.

5 Conclusions

We presented an approach that naturally combines the rank-constraint for identity with illumination modeling. By using the integrability and symmetry constraints, we then achieved a linear algorithm that recovers the albedos and surface normals for a class of face images under the most general setting, i.e., the observation matrix consists of different objects under different illuminations. Further, after separating the illuminations, we obtained illumination-invariant identity signatures which produced good recognition performances under illumination variations. We still need to investigate pose variations and extreme lighting conditions that cause more shadows.

References

1. R. Basri and D. Jacobs, “Lambertian reflectance and linear subspaces,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 2, pp. 218 - 233, 2003.
2. R. Basri and D. Jacobs, “Photometric stereo with general, unknown lighting,” *IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 374-381, 2001.
3. D. Jacobs, “Linear fitting with missing data for structure-from-motion,” *Computer Vision and Image Understanding*, vol. 82, pp. 57 - 81, 2001.
4. A. Shashua, “On photometric issues in 3D visual recognition from a single 2D image,” *International Journal of Computer Vision*, vol. 21, no. 1, pp. 99-122, 1997.
5. A. Shashua and T. R. Raviv, “The quotient image: Class based re-rendering and recognition with varying illuminations,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 23, no. 2, pp. 129-139, 2001.
6. A. L. Yuille, D. Snow, Epstein R., and P. N. Belhumeur, “Determining generative models of objects under varying illumination: Shape and albedo from multiple images using SVD and integrability,” *International Journal of Computer Vision*, vol. 35, no. 3, pp. 203-222, 1999.
7. W. Zhao and R. Chellappa, “Symmetric shape from shading using self-ratio image,” *International Journal of Computer Vision*, vol. 45, pp. 55-752, 2001.
8. A. Georghiades, P. Belhumeur, and D. Kriegman, “From few to many: illumination cone models for face recognition under variable lighting and pose,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 23, pp. 643 -660, 2001.
9. R. T. Frankot and R. Chellappa, “A method for enforcing integrability in shape from shading problem,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. PAMI-10, pp. 439-451, 1987.

10. D. Forsyth, "Shape from texture and integrability," *International Conference on Computer Vision*, pp. 447-453, 2001.
11. I. Shimshoni, Y. Moses, and M. Lindenbaum., "Shape reconstruction of 3D bilaterally symmetric surfaces," *International Journal of Computer Vision*, vol. 39, pp. 97-100, 2000.
12. T. Sim, S. Baker, and M. Bsat, "The CMU pose, illuminatin, and expression (PIE) database," *Prof. Face and Gesture Recognition*, pp. 53-58, 2002.
13. P. J. Philipps, H. Moon, P. Rauss, and S. Rivzi, "The FERET evaluation methodology for face-recognition algorithms," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, pp. 1090-1104, 2000.
14. W. T. Freeman and J. B. Tenenbaum, "Learning bilinear models for two-factor problems in vision," *IEEE Conf. on Computer Vision and Pattern Recognition*, 1997.
15. M. Turk and A. Pentland, "Eigenfaces for recognition," *Journal of Cognitive Neuroscience*, vol. 3, pp. 72-86, 1991.
16. M.E. Brand, "Incremental singular value decomposition of uncertain data with missing values," *European Conference on Computer Vision*, pp. 707-720, 2002.
17. P. Belhumeur, D. Kriegman, and A. Yuille, "The bas-relief ambiguity," *International Journal of Computer Vision*, vol. 35, pp. 33-44, 1999.
18. V. Blanz and T. Vetter, "Face recognition based on fitting a 3D morphable model," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, pp. 1063-1074, 2003.
19. Q. F. Zheng and R. Chellappa, "Estimation of illuminant direction, albedo and shape from shading," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. PAMI-13, pp. 680-702, 1991.
20. A. Shashua and A. Levin. "Linear Image Coding for Regression and Classification using the Tensor-rank Principle." *IEEE Conf. on Computer Vision and Pattern Recognition*, 2001.
21. C. Bregler, A. Hertzmann, and H. Biermann "Recovering Non-Rigid 3D Shape from Image Streams." *IEEE Conf. on Computer Vision and Pattern Recognition*, 2000.
22. M.A.O. Vasilescu and D. Terzopoulos, "Multilinear Analysis of Image Ensembles: TensorFaces," *European Conference on Computer Vision*, 2002.

Appendix-I: Recovering **R**, **F**, and **S** from **H**

This appendix presents an iterative algorithm that recovers **R**, **F** and **S** from the observation matrix **H**. In fact, we also infer $\mathbf{I} = [\downarrow_{\mathbf{x}} [\Rightarrow_i \mathbf{i}_i(\mathbf{x})]]$, which is an indication matrix for **H**.

We first concentrate on the most difficult part of updating **R** with **F**, **S**, and **I** fixed. We take vector derivatives of \mathcal{E} with respect to $\{r_{ij}; i = a, b, c; j = 1, \dots, m\}$ and treat the three terms in \mathcal{E} separately.

[About \mathcal{E}_0 .] With $\mathbf{f}_{j'} \doteq [\downarrow_{j=1}^m f_{j'j}]$ and $\mathbf{s}_{j'} \doteq [s_{j'a}, s_{j'b}, s_{j'c}]^T$,

$$\frac{\partial \mathcal{E}_0}{\partial r_{ij}} = \sum_{j'=1}^n \sum_{\mathbf{x}=1}^d \mathbf{i}_{j'(\mathbf{x})} \{ \hat{\mathbf{w}}(\mathbf{x})^T \mathbf{R}(\mathbf{f}_{j'} \otimes \mathbf{s}_{j'}) - h_{j'(\mathbf{x})} \} \hat{\mathbf{w}}(\mathbf{x}) f_{j'j} s_{j'i}$$

$$\begin{aligned}
&= \sum_{j'=1}^n \sum_{\mathcal{X}=1}^d \hat{\mathbf{i}}_{j'}(\mathbf{x}) \left\{ \sum_{l=a,b,c} \sum_{k=1}^m \hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{lk} f_{j'k} s_{j'l} - h_{j'}(\mathbf{x}) \right\} \hat{\mathbf{w}}(\mathbf{x}) f_{j'j} s_{j'i} \\
&= \sum_{l=a,b,c} \sum_{k=1}^m \left\{ \sum_{j'=1}^n \sum_{\mathcal{X}=1}^d \hat{\mathbf{i}}_{j'}(\mathbf{x}) f_{j'k} s_{j'l} f_{j'j} s_{j'i} \hat{\mathbf{w}}(\mathbf{x}) \hat{\mathbf{w}}(\mathbf{x})^T \right\} \mathbf{r}_{lk} - \sum_{j'=1}^n \sum_{\mathcal{X}=1}^d \hat{\mathbf{i}}_{j'}(\mathbf{x}) h_{j'}(\mathbf{x}) f_{j'j} s_{j'i} \hat{\mathbf{w}}(\mathbf{x}) \\
&= \sum_{l=a,b,c} \sum_{k=1}^m \mathbf{O}_{ij}^{lk} \mathbf{r}_{lk} - \gamma_{ij}, \tag{15}
\end{aligned}$$

where $\{\mathbf{O}_{ij}^{lk}; l = a, b, c; k = 1, \dots, m\}$ are properly defined $3m \times 3m$ matrices, and γ_{ij} is a properly defined $3m \times 1$ vector.

[About \mathcal{E}_1 .] Using forward differences to approximate partial derivatives ³,

$$\begin{aligned}
\frac{\partial \hat{a}_{j(x,y)}}{\partial y} &\simeq \hat{a}_{j(x,y+1)} - \hat{a}_{j(x,y)}; & \frac{\partial \hat{b}_{j(x,y)}}{\partial x} &\simeq \hat{b}_{j(x+1,y)} - \hat{b}_{j(x,y)}; \\
\frac{\partial \hat{c}_{j(x,y)}}{\partial x} &\simeq \hat{c}_{j(x+1,y)} - \hat{c}_{j(x,y)}; & \frac{\partial \hat{c}_{j(x,y)}}{\partial y} &\simeq \hat{c}_{j(x,y+1)} - \hat{c}_{j(x,y)},
\end{aligned} \tag{16}$$

we have

$$\alpha_{j(x,y)} \approx \hat{b}_{j(x+1,y)} \hat{c}_{j(x,y)} - \hat{b}_{j(x,y)} \hat{c}_{j(x+1,y)} + \hat{a}_{j(x,y)} \hat{c}_{j(x,y+1)} - \hat{a}_{j(x,y+1)} \hat{c}_{j(x,y)}. \tag{17}$$

Suppose we are given the product of albedo and surface normal as in (11), we can derive the albedo $p_j(\mathbf{x})$ and surface normals $\hat{a}_j(\mathbf{x})$, $\hat{b}_j(\mathbf{x})$, and $\hat{c}_j(\mathbf{x})$ as follows:

$$p_j(\mathbf{x}) = \sqrt{(\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{aj})^2 + (\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{bj})^2 + (\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{cj})^2}, \tag{18}$$

$$\hat{a}_j(\mathbf{x}) = \frac{\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{aj}}{p_j(\mathbf{x})}, \quad \hat{b}_j(\mathbf{x}) = \frac{\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{bj}}{p_j(\mathbf{x})}, \quad \hat{c}_j(\mathbf{x}) = \frac{\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{cj}}{p_j(\mathbf{x})}. \tag{19}$$

So, their partial derivatives with respect to \mathbf{r}_{aj} are

$$\frac{\partial \hat{a}_j(\mathbf{x})}{\partial \mathbf{r}_{aj}} = \frac{\hat{\mathbf{w}}(\mathbf{x})}{p_j(\mathbf{x})} - \hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{aj} \frac{\hat{\mathbf{w}}(\mathbf{x}) \hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{aj}}{p_j^3(\mathbf{x})} = \frac{1 - \hat{a}_j^2(\mathbf{x})}{p_j(\mathbf{x})} \hat{\mathbf{w}}(\mathbf{x}), \tag{20}$$

$$\frac{\partial \hat{a}_j(\mathbf{x})}{\partial \mathbf{r}_{bj}} = -\hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{aj} \frac{\hat{\mathbf{w}}(\mathbf{x}) \hat{\mathbf{w}}(\mathbf{x})^T \mathbf{r}_{bj}}{p_j^3(\mathbf{x})} = \frac{-\hat{a}_j(\mathbf{x}) \hat{b}_j(\mathbf{x})}{p_j(\mathbf{x})} \hat{\mathbf{w}}(\mathbf{x}), \quad \frac{\partial \hat{a}_j(\mathbf{x})}{\partial \mathbf{r}_{cj}} = \frac{-\hat{a}_j(\mathbf{x}) \hat{c}_j(\mathbf{x})}{p_j(\mathbf{x})} \hat{\mathbf{w}}(\mathbf{x}). \tag{21}$$

Similarly, we can derive their partial derivatives with respect to \mathbf{r}_{bj} and \mathbf{r}_{cj} , which are summarized as follows:

$$\frac{\partial \hat{k}_{j(x)}^l}{\partial \mathbf{r}_{lj}} = \frac{-\hat{k}_{j(x)}^l \hat{l}_{j(x)}}{p_j(\mathbf{x})} \hat{\mathbf{w}}(\mathbf{x}), \quad \frac{\partial \hat{k}_{j(x)}^k}{\partial \mathbf{r}_{kj}} = \frac{1 - \hat{k}_{j(x)}^2}{p_j(\mathbf{x})} \hat{\mathbf{w}}(\mathbf{x}), \quad k, l \in \{a, b, c\}, \quad k \neq l. \tag{22}$$

³ Partial derivatives of boundary pixels require different approximations. But, similar derivations (skipped here due to space limitation) can be derived.

Notice that $\frac{\partial \hat{a}_j(\mathbf{x})}{\partial \mathbf{r}_{bj}} = \frac{\partial \hat{b}_j(\mathbf{x})}{\partial \mathbf{r}_{aj}}$, $\frac{\partial \hat{a}_j(\mathbf{x})}{\partial \mathbf{r}_{cj}} = \frac{\partial \hat{c}_j(\mathbf{x})}{\partial \mathbf{r}_{aj}}$, and $\frac{\partial \hat{b}_j(\mathbf{x})}{\partial \mathbf{r}_{cj}} = \frac{\partial \hat{c}_j(\mathbf{x})}{\partial \mathbf{r}_{bj}}$, which implies saving in computations.

We now compute the partial derivative of $\alpha_{j(x,y)}$ with respect to \mathbf{r}_{aj} :

$$\begin{aligned}
 \frac{\partial \alpha_{j(x,y)}}{\partial \mathbf{r}_{aj}} &= \frac{\partial}{\partial \mathbf{r}_{aj}} \{ \hat{b}_{j(x+1,y)} \hat{c}_{j(x,y)} - \hat{b}_{j(x,y)} \hat{c}_{j(x+1,y)} + \hat{a}_{j(x,y)} \hat{c}_{j(x,y+1)} - \hat{a}_{j(x,y+1)} \hat{c}_{j(x,y)} \} \\
 &= \left\{ \frac{\hat{a}_{j(x,y)} \hat{c}_{j(x,y)}}{p_{j(x,y)} p_{j(x,y+1)}} \hat{\mathbf{w}}_{(x,y)} \hat{\mathbf{w}}_{(x,y+1)}^{\mathbf{T}} - \frac{\hat{a}_{j(x,y+1)} \hat{c}_{j(x,y+1)}}{p_{j(x,y)} p_{j(x,y+1)}} \hat{\mathbf{w}}_{(x,y+1)} \hat{\mathbf{w}}_{(x,y)}^{\mathbf{T}} \right\} \mathbf{r}_{aj} + \\
 &\quad \left\{ \frac{\hat{a}_{j(x+1,y)} \hat{c}_{j(x+1,y)}}{p_{j(x,y)} p_{j(x+1,y)}} \hat{\mathbf{w}}_{(x+1,y)} \hat{\mathbf{w}}_{(x,y)}^{\mathbf{T}} - \frac{\hat{a}_{j(x,y)} \hat{c}_{j(x,y)}}{p_{j(x,y)} p_{j(x+1,y)}} \hat{\mathbf{w}}_{(x,y)} \hat{\mathbf{w}}_{(x+1,y)}^{\mathbf{T}} \right\} \mathbf{r}_{bj} + \\
 &\quad \left\{ \frac{\hat{a}_{j(x,y)} \hat{b}_{j(x,y)}}{p_{j(x,y)} p_{j(x+1,y)}} \hat{\mathbf{w}}_{(x,y)} \hat{\mathbf{w}}_{(x+1,y)}^{\mathbf{T}} - \frac{\hat{a}_{j(x+1,y)} \hat{b}_{j(x+1,y)}}{p_{j(x,y)} p_{j(x+1,y)}} \hat{\mathbf{w}}_{(x+1,y)} \hat{\mathbf{w}}_{(x,y)}^{\mathbf{T}} + \right. \\
 &\quad \left. \frac{1 - \hat{a}_{j(x,y)}^2}{p_{j(x,y)} p_{j(x,y+1)}} \hat{\mathbf{w}}_{(x,y)} \hat{\mathbf{w}}_{(x,y+1)}^{\mathbf{T}} - \frac{1 - \hat{a}_{j(x,y+1)}^2}{p_{j(x,y)} p_{j(x+1,y)}} \hat{\mathbf{w}}_{(x,y+1)} \hat{\mathbf{w}}_{(x,y)}^{\mathbf{T}} \right\} \mathbf{r}_{cj} \\
 &= \mathbf{P}_{aj(x,y)}^a \mathbf{r}_{aj} + \mathbf{P}_{aj(x,y)}^b \mathbf{r}_{bj} + \mathbf{P}_{aj(x,y)}^c \mathbf{r}_{cj} = \sum_{l=a,b,c} \mathbf{P}_{aj(x,y)}^l \mathbf{r}_{lj}, \tag{23}
 \end{aligned}$$

where $\mathbf{P}_{aj(x,y)}^a$, $\mathbf{P}_{aj(x,y)}^b$, and $\mathbf{P}_{aj(x,y)}^c$ are properly defined matrices of dimension $3m \times 3m$. By the same token, using properly defined $\mathbf{P}_{bj(x,y)}^a$, $\mathbf{P}_{bj(x,y)}^b$, $\mathbf{P}_{bj(x,y)}^c$, $\mathbf{P}_{cj(x,y)}^a$, $\mathbf{P}_{cj(x,y)}^b$, and $\mathbf{P}_{cj(x,y)}^c$, we can calculate $\frac{\partial \alpha_{ij(x,y)}}{\partial \mathbf{r}_{ij}} = \sum_{l=a,b,c} \mathbf{P}_{ij(x,y)}^l \mathbf{r}_{lj}$ for $i = a, b, c$, and, finally,

$$\frac{\partial \mathcal{E}_1}{\partial \mathbf{r}_{ij}} = \sum_{\mathbf{x}=1}^d \alpha_{j(\mathbf{x})} \sum_{l=a,b,c} \mathbf{P}_{ij(\mathbf{x})}^l \mathbf{r}_{lj} = \sum_{l=a,b,c} \mathbf{P}_{ij}^l \mathbf{r}_{lj}; \quad \mathbf{P}_{ij}^l \doteq \sum_{\mathbf{x}=1}^d \alpha_{j(\mathbf{x})} \mathbf{P}_{ij(\mathbf{x})}^l. \tag{24}$$

[About \mathcal{E}_2 .] The symmetry constraint term $\beta_{j(\mathbf{x})}$ defined as in (13) can be expressed as

$$\beta_{j(\mathbf{x})}^2 = \mathbf{r}_{aj}^{\mathbf{T}} \mathbf{Q}_{(\mathbf{x})}^a \mathbf{r}_{aj} + \mathbf{r}_{bj}^{\mathbf{T}} \mathbf{Q}_{(\mathbf{x})}^b \mathbf{r}_{bj} + \mathbf{r}_{cj}^{\mathbf{T}} \mathbf{Q}_{(\mathbf{x})}^c \mathbf{r}_{cj}, \tag{25}$$

where $\mathbf{Q}_{(\mathbf{x})}^a$, $\mathbf{Q}_{(\mathbf{x})}^b$, and $\mathbf{Q}_{(\mathbf{x})}^c$ are symmetric matrices with size $3m \times 3m$:

$$\mathbf{Q}_{(\mathbf{x})}^a = (\hat{\mathbf{w}}_{(\mathbf{x})} + \hat{\mathbf{w}}_{(\bar{\mathbf{x}})})(\hat{\mathbf{w}}_{(\mathbf{x})} + \hat{\mathbf{w}}_{(\bar{\mathbf{x}})})^{\mathbf{T}}, \quad \mathbf{Q}_{(\mathbf{x})}^b = (\hat{\mathbf{w}}_{(\mathbf{x})} - \hat{\mathbf{w}}_{(\bar{\mathbf{x}})})(\hat{\mathbf{w}}_{(\mathbf{x})} - \hat{\mathbf{w}}_{(\bar{\mathbf{x}})})^{\mathbf{T}}, \quad \mathbf{Q}_{(\mathbf{x})}^c = \mathbf{Q}_{(\mathbf{x})}^b. \tag{26}$$

The derivatives of $\beta_{j(\mathbf{x})}^2/2$ and \mathcal{E}_2 with respect to \mathbf{r}_{aj} , \mathbf{r}_{bj} , and \mathbf{r}_{cj} are

$$\frac{\partial \{\beta_{j(\mathbf{x})}^2/2\}}{\partial \mathbf{r}_{ij}} = \mathbf{Q}_{(\mathbf{x})}^i \mathbf{r}_{ij}; \quad \frac{\partial \mathcal{E}_2}{\partial \mathbf{r}_{ij}} = \sum_{\mathbf{x}=1}^d \mathbf{Q}_{(\mathbf{x})}^i \mathbf{r}_{ij} = \mathbf{Q}^i \mathbf{r}_{ij}; \quad \mathbf{Q}^i = \sum_{\mathbf{x}=1}^d \mathbf{Q}_{(\mathbf{x})}^i. \tag{27}$$

Putting the above derivations together and using $\frac{\partial \mathcal{E}}{\partial \mathbf{r}_{ij}} = 0$, we have

$$\sum_{l=a,b,c} \sum_{k=1}^m \mathbf{O}_{ij}^{lk} \mathbf{r}_{lk} + \lambda_1 \sum_{l=a,b,c} \mathbf{P}_{ij}^l \mathbf{r}_{lj} + \lambda_2 \mathbf{Q}^i \mathbf{r}_{ij} = \gamma_{ij}; \quad i = a, b, c; \quad j = 1, \dots, m. \tag{28}$$

We therefore arrive at a set of equations linear in $\{r_{ij}; i = a, b, c; j = 1, \dots, m\}$ that can be solved easily. After finding the new \mathbf{R} , we normalize it using $\mathbf{R} = \mathbf{R} / \|\mathbf{R}\|_2$.

We now illustrate how to update $\mathbf{F} = [\Rightarrow_i \mathbf{f}_i]$, $\mathbf{S} = [\Rightarrow_i \mathbf{s}_i]$, and $\mathbf{I} = [\Rightarrow_i \mathbf{I}_i]$ with \mathbf{R} fixed (or \mathbf{W} fixed). First notice that they are only involved in \mathcal{E}_0 . Moreover, \mathbf{f}_i , \mathbf{s}_i and \mathbf{I}_i are related with only the image \mathbf{h}_i . This becomes the same as the illumination separation problem defined in Section 4 and Appendix-II presents such a recovery algorithm, which also is iterative in nature. After running one iterative step to obtain the updated \mathbf{F} , \mathbf{S} , and \mathbf{I} , we proceed to update \mathbf{R} again and this process carries on until convergence.

Appendix-II: Recovering \mathbf{f} and \mathbf{s} from \mathbf{h} given \mathbf{W}

This recovery task is equivalent to minimizing the cost function defined as

$$\mathcal{E}_{\mathbf{h}}(\mathbf{f}, \mathbf{s}) \doteq \|\mathbf{I} \circ (\mathbf{h} - \mathbf{W}(\mathbf{f} \otimes \mathbf{s}))\|^2 + (\mathbf{1}^T \mathbf{f} - 1)^2, \quad (29)$$

where $\mathbf{I}_{d \times 1}$ indicates the inclusion or exclusion of the pixels of the image \mathbf{h} and \circ denotes the Hadamard (element-wise) product. Notice that (29) actually can be easily generalized as a cost function for robust estimation if the L_2 norm $\|\cdot\|$ is replaced by a robust function, and \mathbf{I} by an appropriate weight function.

The following algorithm is an extension of bilinear analysis, with occlusion embedded. Firstly, we solve the least square (LS) solution \mathbf{f} , given \mathbf{s} and \mathbf{I} .

$$\mathbf{f} = \begin{bmatrix} \mathbf{W}_{\mathbf{f}} \\ \mathbf{1}^T \end{bmatrix}^{\dagger} \begin{bmatrix} \mathbf{I} \circ \mathbf{h} \\ 1 \end{bmatrix}; \quad \mathbf{W}_{\mathbf{f}} \doteq [\Rightarrow_i (\mathbf{T}_i \mathbf{s})]_{d \times m}. \quad (30)$$

where $[\cdot]^{\dagger}$ denotes the pseudo-inverse. Secondly, we solve the LS solution \mathbf{s} , given \mathbf{f} and \mathbf{I} :

$$\mathbf{s} = \mathbf{W}_{\mathbf{S}}^{\dagger} (\mathbf{I} \circ \mathbf{h}); \quad \mathbf{W}_{\mathbf{S}} \doteq [[\Rightarrow_i \mathbf{a}_i] \mathbf{f}, [\Rightarrow_i \mathbf{b}_i] \mathbf{f}, [\Rightarrow_i \mathbf{c}_i] \mathbf{f}]_{d \times 3} \doteq [\mathbf{A} \mathbf{f}, \mathbf{B} \mathbf{f}, \mathbf{C} \mathbf{f}], \quad (31)$$

where $\mathbf{A}_{d \times m} \doteq [\Rightarrow_i \mathbf{a}_i]$, $\mathbf{B}_{d \times m} \doteq [\Rightarrow_i \mathbf{b}_i]$, and $\mathbf{C}_{d \times m} \doteq [\Rightarrow_i \mathbf{c}_i]$ contain the information on the product of albedos and x , y , and z directions of the surface normals, respectively. In the third step, given \mathbf{f} and \mathbf{s} we update \mathbf{I} as follows⁴:

$$\mathbf{I} = [|\mathbf{h} - \mathbf{W}(\mathbf{f} \otimes \mathbf{s})| < \eta], \quad (32)$$

where η is a pre-defined threshold.

Note that in Eqs. (30) and (31), additional saving in computation is possible. We can form matrices $\mathbf{W}'_{\mathbf{f}}$ and $\mathbf{W}'_{\mathbf{S}}$ and vector \mathbf{h}' , with a reduced dimension, from $\mathbf{W}_{\mathbf{f}}$, $\mathbf{W}_{\mathbf{S}}$, and \mathbf{h} , respectively, by discarding those rows corresponding to the excluded pixels and applying the primed version in (30) and (31).

For fast convergence, we use the following initial values in our implementation. We estimate \mathbf{s} using the algorithm presented in [19] and set \mathbf{I} to exclude those pixels whose intensities are smaller than a certain threshold.

⁴ This is a Matlab operation which performs an element-wise comparison.