

Real-time Robust Principal Components' Pursuit

Chenlu Qiu, Namrata Vaswani

Dept. of Electrical and Computer Engineering, Iowa State University
Ames, IA, 50010
{chenlu, namrata}@iastate.edu

Abstract—In the recent work of Candes et al, the problem of recovering low rank matrix corrupted by i.i.d. sparse outliers is studied and a very elegant solution, principal component pursuit, is proposed. It is motivated as a tool for video surveillance applications with the background image sequence forming the low rank part and the moving objects/persons/abnormalities forming the sparse part. Each image frame is treated as a column vector of the data matrix made up of a low rank matrix and a sparse corruption matrix. Principal component pursuit solves the problem under the assumptions that the singular vectors of the low rank matrix are spread out and the sparsity pattern of the sparse matrix is uniformly random. However, in practice, usually the sparsity pattern and the signal values of the sparse part (moving persons/objects) change in a correlated fashion over time, for e.g., the object moves slowly and/or with roughly constant velocity. This will often result in a low rank sparse matrix.

For video surveillance applications, it would be much more useful to have a real-time solution. In this work, we study the online version of the above problem and propose a solution that automatically handles correlated sparse outliers. In fact we also discuss how we can potentially use the correlation to our advantage in future work. The key idea of this work is as follows. Given an initial estimate of the principal directions of the low rank part, we causally keep estimating the sparse part at each time by solving a noisy compressive sensing type problem. The principal directions of the low rank part are updated every-so-often. In between two update times, if new Principal Components' directions appear, the “noise” seen by the Compressive Sensing step may increase. This problem is solved, in part, by utilizing the time correlation model of the low rank part. We call the proposed solution “Real-time Robust Principal Components' Pursuit”. It still requires the singular vectors of the low rank part to be spread out, but it does not require i.i.d.-ness of either the sparse part or the low rank part.

I. INTRODUCTION

Principal Components' Analysis (PCA) tries to find the “principal components' space” with the smallest dimension that spans a given dataset. In practice, data is noisy and in this case PCA finds the smallest subspace to represent the dataset with a given mean squared error (MSE) tolerance.

Given a low rank data matrix $M \in \mathbb{R}^{m \times n}$ (each column of M is one data vector), PCA finds its principal components (PCs) as the left singular vectors of M that have nonzero singular values. This is the same as first estimating the data covariance as $(1/n)MM^T$, computing its eigenvalue decomposition (EVD) and retaining eigenvectors corresponding to nonzero eigenvalues. When data is noisy, this is replaced by

arranging the eigenvectors in decreasing order of eigenvalues, and retaining the smallest number of eigenvectors so that the sum of the remaining eigenvalues (which is equal to the residual MSE) is less than the MSE tolerance.

When the noise is small, the above approach works well. However, covariance matrix estimation, and hence the corresponding EVD, are sensitive to even a few large outliers in the data. Unfortunately, in practice these do occur, e.g. when trying to compute the principal components' subspace for a video sequence, parts of it may get occluded by other moving objects. There has been a large amount of work in literature on “Robust PCA”, e.g. [1], [2], [3], [4], [5], [6], [7], [8], [9], most of which either assumes the locations of the missing/corrupted data points are known [3], which is not a practical assumption, or (ii) first tries to detect the corrupted pixels and then either fills in the corrupted location using some heuristics or (iii) often just removes the entire outlier vector. In a series of recent works [10], [11], [12], a very elegant solution to this problem was provided that treats the outlier as a sparse vector. In [10], the data matrix M consists of a low rank matrix that is corrupted by sparse outliers, i.e.

$$M = L + S$$

where L is a low rank matrix having a singular value decomposition (SVD) $L \stackrel{SVD}{=} UDV^T$ and S is sparse and can have arbitrary large magnitude. Let $\|L\|_*$ denotes the nuclear norm of L , i.e., the sum of singular values of L . It is shown in [10] that L and S can be recovered with high probability by solving a convex optimization problem, named as Principal Component Pursuit (PCP),

$$\begin{aligned} \min_{L,S} \quad & \|L\|_* + \lambda \|S\|_1 \\ \text{subject to} \quad & L + S = M \end{aligned} \quad (1)$$

provided the singular vectors of L are spread out (not sparse), the support and signs of S are uniformly random (thus not low rank), the rank of L and the fraction of corrupted entries in S are both sufficient small. A more recent work, [12], extends the result of [10] showing that, with a proper weighting parameter λ , PCP can recover L and S with high probability even if the size of support set of S is large, as long as the rank of L is small enough. But it requires that S has random support and random signs.

PCP [10] is motivated as a tool for surveillance applications, with the background variations approximately lying in a low dimension subspace, and the sparse part being the “moving persons” or “abnormalities” to be detected. It is an offline

method which treats each image frame as a column vector of the data matrix M . While this is a very elegant and novel idea, there are certain limitations.

- 1) In surveillance, it would be more useful to obtain the estimates of the sparse part on-the-fly rather than offline.
- 2) The sparsity pattern (support and signs) of the sparse part may change slowly or in a correlated fashion, which may result in a low rank sparse matrix. In this case, PCP assumption will not get satisfied and as a result it will not work, e.g. see Fig.2.
- 3) The principal directions (set of eigenvectors corresponding to nonzero eigenvalues) can change over time. So the rank of the matrix L will keep increasing over time thus making PCP impossible to do after sometime.

This last issue may get resolved by not using all frames of M , but only the latest image frames. But the first two issues still remain.

In this paper, we propose an online approach to solve this problem. Our goal is to causally keep estimating the sparse part S_t at each time, and to keep updating the principal directions every-so-often. The t -th column of M , M_t , is the data acquired at time t . It can be split as

$$M_t = L_t + S_t = [U \ I] \begin{bmatrix} x_t \\ S_t \end{bmatrix} \quad (2)$$

where $x_t := U^T L_t$ and the matrix U is an *unknown* $m \times m$ orthonormal matrix. The support of the vector S_t changes slowly over time. Given an initial estimate of $P_t := (U)_{N_t}$, denoted \hat{P}_t , we solve for the sparse vector S_t by first finding the orthogonal complement matrix $\hat{P}_{t,\perp}$ and then using the projection of M_t onto $\hat{P}_{t,\perp}$, denoted by y_t ,

$$y_t := \hat{P}_{t,\perp}^T M_t = \hat{P}_{t,\perp}^T L_t + \hat{P}_{t,\perp}^T S_t.$$

to solve for S_t . Notice that if $\hat{P}_t \approx P_t$ the first term will be close to zero and can be treated as “noise”. When $\hat{P}_t \neq P_t$ (new directions added), the “noise” can be reduced by using the time correlation model on L_t . Furthermore, recent estimates of $L_t := M_t - S_t$ are stored and used to periodically update P_t as described in Sec. III-D. There are also some limitations of our method.

- 1) We need an approximately accurate initial estimate of the PCs’ basis, \hat{P}_0 , which is easy to get using training data without sparse corruptions.
- 2) The orthogonal complement $\hat{P}_{t,\perp}$ needs to satisfy some conditions for Compressive Sensing to succeed.
- 3) An appropriate choice of constraint parameter ϵ is needed for estimating S_t .

The above idea is somewhat related to that of [13] in that both try to cancel the “message” signal and only solve for the sparse “error” signal, but with the big difference that in [13], P_t is *known*. Other related work which also uses P_t known is [14], [15]. However in our problem P_t is *unknown* and can change with time. Our method requires the columns of $\hat{P}_{t,\perp}$ be spread out (not sparse), but it does not require S_t to have independent nonzero entries. In fact, we can utilize the correlated support change of S_t over time, t , to our advantage

in future work. Since we update the principal directions on-the-fly, the dimension of the principal subspace remains bounded.

A model similar to (2) but for a static problem and with U being a *known* matrix, was introduced in [16], [17]. A method, termed as pursuit of justice (PJ), is introduced to solve for the sparse vector $u = [x_t, S_t]^T$ which solving the following ℓ_1 minimization problem

$$\begin{aligned} \min_u \quad & \|u\|_1 \\ \text{subject to} \quad & M_t = Au \end{aligned} \quad (3)$$

where $A := [U \ I]$. Notice that in our problem U is *unknown*, and thus we cannot use sparse reconstruction techniques to find x_t . Given an estimate \hat{P}_t , the above can be modified to $A = [\hat{P} \ \hat{P}_{t,\perp} \ I]$. However, this does not work as shown in Fig. 3.

A. Notations

The set operations \cup , \cap and \setminus have the usual meanings. For any set $T \subset \{1, \dots, m\}$, T^c denotes the complement set of T , i.e., $T^c := \{1, \dots, m\} \setminus T$.

For a non diagonal matrix A , we let A_i denote the i th column of A and we let A_T denote a matrix composed of the columns of A indexed by T . For two set T_1 and T_2 , we let A_{T_1, T_2} denote a submatrix of A consisting of the rows indexed by T_1 and columns indexed by T_2 . For a diagonal matrix Q , Q_T denotes a submatrix of Q consisting of the rows and columns indexed by T . In other words, Q_T is a diagonal matrix with $(Q_T)_{j,j} = (Q)_{T_j, T_j}$.

For vector v , v_i denotes the i th entry of v and v_T denotes a vector consisting of the entries of v indexed by T . $\|v\|_k$ denotes the ℓ_k norm of v . The support of v , $\text{supp}(v)$, is the set of indices at which v has nonzero value, $\text{supp}(v) := \{i : v_i \neq 0\}$.

We use \emptyset to denote an empty set or an empty matrix.

II. PROBLEM DEFINITION AND SIGNAL MODEL

The t th column of M , $M_t \in \mathbb{R}^{m \times 1}$, is the data at time t which can be split as

$$\begin{aligned} M_t &= L_t + S_t \\ L_t &= U x_t = P_t a_t \end{aligned}$$

where $x_t := U^T L_t$ and S_t are sparse vectors with slowly changing support $N_t := \text{supp}(x_t)$ and $T_t := \text{supp}(S_t)$, respectively. N_t is modeled as being piecewise constant with time. The vector $a_t := (x_t)_{N_t}$ is the none-zero part of x_t . The principal components’ basis at each time t , $P_t := U_{N_t}$, is a submatrix of U whose columns span the principal components’ subspace at time t . It is unknown and can change over time.

Since the matrix U does not change with time (in this work), the only way P_t changes is when the set N_t changes. This happens every d frames. We assume that x_t and hence $L_t = U x_t$ follows a *piecewise stationary model with nonstationary transients when switching pieces*. For every d frames, there are some supporting indices get added or deleted from N_t . Specifically when an element i gets added into the support, it gets added with an initial small variance $\theta \sigma_i^2$ (with $0 < \theta < 1$)

and then at future times follows a first order autoregressive (AR-1) model with AR parameter f and stable variance σ_i^2 . Recall that an AR-1 model is asymptotically stationary. Thus, after the initial transient period, x_t is stationary until the next support change time. Before an element i gets deleted, it starts decaying as $(x_t)_i = f_d(x_{t-1})_i$, with $0 < f_d < f < 1$, and soon decays to zero.

A. Mathematical description of signal model for x_t (and hence for L_t)

The support set of x_t , N_t , is a union of three disjoint sets Δ_t , D_t , and E_t , i.e., $N_t = \Delta_t \cup D_t \cup E_t$. The addition set $\Delta_t := N_t \setminus N_{t-1}$ is the set of indices for the new appearing eigenvectors $(U)_{\Delta_t}$. The set $D_t \subset (N_t \cap N_{t-1})$ is the set of indices of those eigenvectors whose eigenvalues are decreasing at time t . The set $E_t := N_t \cap N_{t-1} \setminus D_t$ is the set of indices for existing eigenvectors with non-decreasing eigenvalues. The sets D_t and Δ_t can be empty. For any time τ with ‘‘decreasing’’ set D_τ , we assume that D_τ will not get added to N_t for any $t > \tau$.

Let $\Sigma = \text{diag}(\sigma_i^2)$, $i = 1, \dots, m$, be a diagonal matrix with non-increasing positive diagonal elements, i.e. σ_i^2 satisfying $\sigma_i^2 \geq \sigma_{i+1}^2$. We model x_t as

$$\begin{aligned} x_0 &= 0, \quad N_0 = \emptyset \\ x_t &= F_t x_{t-1} + \nu_t, \quad \nu_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, Q_t) \end{aligned} \quad (4)$$

where F_t and Q_t are two diagonal matrices defined as below

$$\begin{aligned} (F_t)_{\Delta_t} &= 0, \quad (Q_t)_{\Delta_t} = \theta(\Sigma)_{\Delta_t}, \\ (F_t)_{E_t} &= fI, \quad (Q_t)_{E_t} = (1 - f^2)(\Sigma)_{E_t}, \\ (F_t)_{D_t} &= f_d I, \quad (Q_t)_{D_t} = 0, \\ (F_t)_{N_t^c} &= 0, \quad (Q_t)_{N_t^c} = 0. \end{aligned}$$

where f , f_d , and θ are scalars satisfying $0 < f_d < f < 1$ and $0 < \theta < 1$.

From the model on x_t , we notice the following:

- a) At time $t = \tau$, $(x_\tau)_{\Delta_\tau}$ starts with

$$(x_\tau)_{\Delta_\tau} \sim \mathcal{N}(0, \theta(\Sigma)_{\Delta_\tau}).$$

Small θ ensures that new directions gets added at a small value and increase slowly. $(x_\tau)_{D_\tau}$ decays as

$$(x_\tau)_{D_\tau} = f_d(x_{\tau-1})_{D_\tau}$$

$(x_\tau)_{E_\tau}$ follows an AR-1 model with parameter f :

$$(x_\tau)_{E_\tau} = f(x_{\tau-1})_{E_\tau} + (v_\tau)_{E_\tau}$$

- b) At time $t > \tau$, if Δ_τ is not removed from the support set, the variance of $(x_t)_{\Delta_\tau}$ gradually increases as

$$(x_t)_i \sim \mathcal{N}(0, (1 - (1 - \theta)f^{2(t-\tau)})_{\Sigma_{i,i}}), \quad i \in \Delta_\tau$$

Eventually, the variance of $(x_t)_{\Delta_\tau}$ converges to $(\Sigma)_{\Delta_\tau}$. For example, with $f = 0.9$ and $\theta = 0.4$, the variance of $(x_t)_{\Delta_\tau}$ gets to $0.9(\Sigma)_{\Delta_\tau}$ in 18 frames.

- c) At time $t > \tau$, the variance of $(x_t)_{D_\tau}$ decays as

$$(x_t)_{D_\tau} \sim \mathcal{N}(0, f_d^{2(t-\tau)}(\Sigma)_{D_\tau})$$

Eventually, $(x_t)_{D_\tau}$ decays to zero. For example, with $f_d = 0.1$, the variance of $(x_t)_{D_\tau}$ decrease to $0.0001(\Sigma)_{D_\tau}$ in 2 frames.

B. Model for S_t

Recall that in video surveillance applications, the data matrix M is obtained by stacking each image frame as a column vector, whose low rank component L corresponds to background variation lying in a low rank subspace and sparse component S captures the moving objects in the foreground. In this work, we use a simple model for the sparse component S modeling the activity of the moving objects as described below.

We assume that, in each image frame, there are k ($k \geq 1$) objects in the foreground. Each object occupies a 3×3 pixel block which has nonzero pixel values. All other pixels in the foreground have zero values. Let CG_t^i denote the coordinate of the center of gravity of the i th object at time t , $i = 1, \dots, k$. For the next image frame, each CG_t^i can either be static with probability p or move one step to the left/right/top/bottom with probability $(1 - p)/4$ each, i.e., for $i = 1, \dots, k$,

$$CG_t^i = \begin{cases} CG_{t-1}^i & \text{with probability } p \\ CG_{t-1}^i + (1, 0), & \text{with probability } (1 - p)/4 \\ CG_{t-1}^i + (-1, 0), & \text{with probability } (1 - p)/4 \\ CG_{t-1}^i + (0, 1), & \text{with probability } (1 - p)/4 \\ CG_{t-1}^i + (0, -1), & \text{with probability } (1 - p)/4 \end{cases}$$

with $p = 0.8$. The pixels in each block move accordingly. Except if the objects move very fast or if they are very small, there will be overlap between their regions from frame to frame. We then stack the resulting foreground image frame as columns of S . Clearly, the support of S_t , t th column of S , is time correlated and the signs of these nonzero entries are fixed. This is quite different from [10] and [12] where random support and random signs are assumed on the sparse part S .

III. REAL-TIME ROBUST PCP

An overview of our method, real-time robust PCP (RR-PCP), is shown in Fig.1. We first discuss the approach to recursively reconstruct the sparse component S_t . Next, we discussed the way we track the changes of the principal directions. Finally, a complete algorithm is given in Algorithm 2.

A. RR-PCP: recursively reconstruction of the sparse part S_t

Using \hat{P}_t , which is an estimate of principal components P_t at time t , we can rewrite L_t and M_t as

$$\begin{aligned} L_t &= \hat{P}_t \alpha_t + \hat{P}_{t,\perp} \beta_t \\ M_t &= \hat{P}_t \alpha_t + \hat{P}_{t,\perp} \beta_t + S_t \end{aligned}$$

where $\hat{P}_{t,\perp}$ is an orthogonal complement of \hat{P}_t ; $\alpha_t := \hat{P}_t^T L_t$ is the projection of L_t onto the subspace spanned by \hat{P}_t ; and $\beta_t := \hat{P}_{t,\perp}^T L_t$ is the projection of L_t onto the subspace spanned by $\hat{P}_{t,\perp}$. Notice that \hat{P}_t is an estimate of P_t . It is either just a slight rotation of P_t with $\text{span}(P_t) = \text{span}(\hat{P}_t)$ or there may be some missing and extra principal directions. The column vectors of $\hat{P}_{t,\perp}$ are the eigenvectors spanning the null space of \hat{P}_t^T . The orthogonal complement $\hat{P}_{t,\perp}$ is not unique.

Let

$$y_t := \hat{P}_{t,\perp}^T M_t = \hat{P}_{t,\perp}^T S_t + \beta_t \quad (5)$$

If there is no missing principal direction, i.e., $\text{span}(P_t) \subseteq \text{span}(\hat{P}_t)$, $\beta_t = 0$. If there are missing principal directions, $\text{span}(P_t) \not\subseteq \text{span}(\hat{P}_t)$ and $\beta_t \neq 0$. In this case, β_t in (5) is the “noise” resulting from the estimation error of current principal directions. This now becomes a noisy sparse reconstruction problem, with “noise” β_t . When $\|\beta_t\|_2^2$ is not very large, we can causally recover S_t by solving

$$\min_s \|s\|_1 \text{ s.t. } \|\hat{P}_{t,\perp}^T(M_t - s)\|_2^2 \leq \epsilon \quad (6)$$

and hence estimate L_t as

$$\hat{L}_t = M_t - \hat{S}_t$$

where ϵ is a parameter with some small positive value.

For the case of missing principal directions, $\text{span}(P_t) \not\subseteq \text{span}(\hat{P}_t)$. Let $\mathcal{S}_{t,\text{miss}} := \text{span}(P_t) \setminus \text{span}(\hat{P}_t)$ denote the “missing” subspace and let $P_{t,\text{miss}}$ be its orthonormal basis matrix. Thus $\text{span}(P_{t,\text{miss}}) = \mathcal{S}_{t,\text{miss}}$ and

$$\text{span}(P_t) = \text{span}(\hat{P}_t) \oplus \text{span}(P_{t,\text{miss}}) \quad (7)$$

$$\text{span}(\hat{P}_{t,\perp}) = \text{span}(P_{t,\perp}) \oplus \text{span}(P_{t,\text{miss}}) \quad (8)$$

Therefore,

$$\beta_t = \hat{P}_{t,\perp}^T L_t = \hat{P}_{t,\perp}^T P_{t,\text{miss}} P_{t,\text{miss}}^T L_t \quad (9)$$

with $P_{t,\text{miss}}^T L_t$ being the projection of L_t onto the subspace $\mathcal{S}_{t,\text{miss}}$. If $P_{t,\text{miss}}^T L_t$ starts with small values, $\|\beta_t\|_2^2$ shall be small and it can increase over time. When $\|\beta_t\|_2^2$ is getting too large, (6) may give incorrect estimate \hat{S}_t . Thus, we need to update \hat{P}_t and get those missing directions detected.

B. Canceling the “noise” using the time correlation of x_t

It is expensive to update \hat{P}_t and $\hat{P}_{t,\perp}$ very frequently, especially for some real-time applications. But notice that we can cancel out some of β_t by using the model on x_t from Sec.II-A. We modify (6) as

$$\min_s \|s\|_1 \text{ s.t. } \|\hat{P}_{t,\perp}^T(M_t - s - f\hat{L}_{t-1})\|_2^2 \leq \epsilon \quad (10)$$

Let $\hat{\beta}_{t-1} := \hat{P}_{t-1,\perp}^T \hat{L}_{t-1}$. Note that in (10), the “noise” is $\beta_t - f\hat{\beta}_{t-1}$ while in (6), the “noise” is β_t .

Next, we discuss an example showing that the “noise” $\beta_t - f\hat{\beta}_{t-1}$ in (10) is smaller than the “noise” β_t in (6).

Suppose at time $t = \tau - 1$, we have an exact estimate of all principal directions, $\hat{P}_{\tau-1} = P_{\tau-1}$. At time $t = \tau$, a support change occurs with $N_\tau = N_{\tau-1} \cup \Delta_\tau$ and $D_\tau = \emptyset$. The principal directions at time $t = \tau$ are $P_\tau = [P_{\tau-1}, P_{\Delta_\tau}]$ where $P_{\Delta_\tau} = U_{\Delta_\tau}$ are the new added principal directions. However, this change is unknown to us and we just use $\hat{P}_\tau = \hat{P}_{\tau-1} = P_{\tau-1}$. Therefore,

$$\text{span}(P_\tau) = \text{span}(P_{\tau-1}) \oplus \text{span}(P_{\Delta_\tau})$$

Thus, at time $t \geq \tau$,

$$\begin{aligned} \beta_t &= \hat{P}_{\tau,\perp}^T L_t = \hat{P}_{\tau,\perp}^T (P_{\tau-1} P_{\tau-1}^T L_t + P_{t,\Delta_\tau} P_{t,\Delta_\tau}^T L_t) \\ &= \hat{P}_{\tau,\perp}^T P_{t,\Delta_\tau} P_{t,\Delta_\tau}^T U x_t \end{aligned} \quad (11)$$

Assuming L_{t-1} is correctly recovered, i.e., $\hat{L}_{t-1} = L_{t-1}$, so $\hat{\beta}_{t-1} = \beta_{t-1}$, then

$$\begin{aligned} \beta_t - f\hat{\beta}_{t-1} &= \hat{P}_{\tau,\perp}^T P_{t,\Delta_\tau} P_{t,\Delta_\tau}^T U \begin{bmatrix} (\nu_t)_{\Delta_\tau} \\ (\nu_t)_{E_\tau} \\ 0 \end{bmatrix} \\ &= \hat{P}_{\tau,\perp}^T P_{t,\Delta_\tau} (\nu_t)_{\Delta_\tau} \end{aligned}$$

Let $B = (\hat{P}_{\tau,\perp}^T P_{t,\Delta_\tau})^T (\hat{P}_{\tau,\perp}^T P_{t,\Delta_\tau})$, then at time $t \geq \tau$,

$$\begin{aligned} \mathbb{E}(\|\beta_t\|_2^2) &= \sum_{i \in \Delta_\tau} B_{i,i} (1 - (1 - \theta)f^{2(t-\tau)}) \sigma_i^2 \\ \mathbb{E}(\|\beta_t - f\hat{\beta}_{t-1}\|_2^2) &= \sum_{i \in \Delta_\tau} B_{i,i} (1 - f^2) \sigma_i^2 \end{aligned} \quad (12)$$

Clearly, $\mathbb{E}(\|\beta_t\|_2^2)$ shall be much larger than $\mathbb{E}(\|\beta_t - f\hat{\beta}_{t-1}\|_2^2)$. For example, with $f = 0.9$, $\theta = 0.4$, at $t = \tau + 1$, $\mathbb{E}(\|\beta_t\|_2^2)$ is $0.514 \sum_{i \in \Delta_\tau} B_{i,i} \sigma_i^2$ while $\mathbb{E}(\|\beta_t - f\hat{\beta}_{t-1}\|_2^2)$ is $0.19 \sum_{i \in \Delta_\tau} B_{i,i} \sigma_i^2$; and the ratio of $\mathbb{E}(\|\beta_t\|_2^2)$ and $\mathbb{E}(\|\beta_t - f\hat{\beta}_{t-1}\|_2^2)$ keeps increasing over time. Recall that β_t and $\beta_t - f\hat{\beta}_{t-1}$ are the “noise” in (6) and (10), therefore, assuming we have an accurate estimate of principal basis at time $\tau - 1$, (10) shall give more accurate estimate of \hat{S}_t than (6). We show a plot of the expectations of $\|\beta_t\|_2^2$ and $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ in Fig. 4.

In (6) and (10), we need an appropriate parameter ϵ which should be proportional to the “noise” term $\|\beta_t\|_2^2$ in (6) or $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ in (10). The “noise” $\|\beta_t\|_2^2$ and $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ also changes over time. Thus, we shall set the ϵ adaptively. In our work, we use ϵ proportionally to $2\|\hat{\beta}_{t-1}\|_2^2$ for (6) and ϵ proportional to $2\|\hat{\beta}_{t-1} - f\hat{\beta}_{t-2}\|_2^2$ for (10).

If the constraint is too tight (ϵ is too small), (6) or (10) may give solutions with some small nonzero values outside the true support set T_t (also verified by our numerical experiments). As first done in [18], we can also do support thresholding followed by least square estimation to reduce these errors, i.e., we can solve

$$\hat{T}_t = \{i : (\hat{S}_t)_i \geq \gamma\} \quad (13)$$

$$(\hat{S}_t)_{\hat{T}_t} = ((\hat{P}_{t,\perp}^T)_{\hat{T}_t})^\dagger (y_t - f\hat{P}_{t,\perp}^T \hat{L}_{t-1}) \quad (14)$$

$$(\hat{S}_t)_{\hat{T}_t^c} = 0$$

$$\hat{L}_t = M_t - \hat{S}_t \quad (15)$$

C. Using model on S_t

Currently, we do not use the model on S_t . A simple way to use it is to do modified-CS [19] as

$$\min_s \|s_{T_{\text{pred}}}\|_1 \text{ s.t. } \|\hat{P}_{t,\perp}^T(M_t - s - f\hat{L}_{t-1})\|_2^2 \leq \epsilon \quad (16)$$

with T_{pred} being an estimate of the support of S_t . As [19] shows, if T_{pred} is an approximately correct estimate of currently support, (16) should improve the performance of (10).

In previous work [20], [21], we used the previous support estimate, \hat{T}_{t-1} , as T_{pred} . This is sufficient for the problems considered in [20], where support changes very slowly over time, e.g. in case of wavelet coefficients of a medical image sequences. But for our current problem, even with one or two pixel motion between frames, the support change will be significant and \hat{T}_{t-1} will have large error w.r.t. T_t . A better

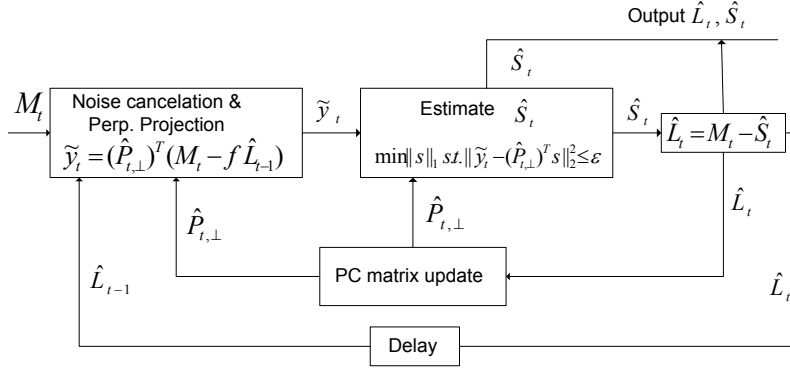


Fig. 1: Real-time Robust PCP

solution, is to use the motion model to predict the object(s)' location in the next frame and use this prediction to obtain T_{pred} at time t . The details of how to do this, especially for multiple objects, will be worked out in future work.

D. RR-PCP: Recursively estimating the low rank part

When some new principal directions appear, we need to detect these directions timely before the “noise” gets large. Now, $\mathbb{E}(\|\beta_t - f\hat{\beta}\|_2^2)$ in (12) seems not increase with time. But this assumes $\hat{L}_{t-1} = L_{t-1}$ which is not true. When some existing directions vanish, they also need to be removed from \hat{P}_t . Otherwise, the number of estimated principal directions keeps increasing and thus the number of columns in $\hat{P}_{t,\perp}$, which is the number of measurements for (10), keeps decreasing.

At initial time, we have the training data $L^0 := [L_1, \dots, L_{t_0}]$, which contains no sparse component. According to our signal model (4), the data sequence L_t is time correlated. Thus, we need a long sequence's data to get an accurate estimate of its covariance. But notice in our model, the sequence $L_t - fL_{t-1}$ is time independent and has same eigenvectors as L_t . Thus, we estimate principal directions of L^0 by estimating the covariance of $L_t - fL_{t-1}$ and computing its EVD. Let P_0 and G_0 be the eigenvectors and (non-zero) eigenvalues of the covariance matrix of $L_t - fL_{t-1}$, $t = 2, \dots, t_0$, i.e. $P_0 G_0 P_0^T = L^0 (L^0)^T$. Let \hat{P}_{stable} be an orthonormal matrix whose columns are the correctly estimated eigenvectors and let \hat{G}_{stable} be a diagonal matrix whose diagonal elements are the correspondingly correctly estimated eigenvalues. Let $\hat{P}_t = \hat{P}_{\text{stable}} = P_0$ and let $\hat{G}_t = \hat{G}_{\text{stable}} = G_0$.

Our PCs update procedure is designed to estimate the current principal directions for data generated according to the piecewise stationary model on x_t (and hence on L_t) that was described in Sec. II-A¹. Assume that every d frames, k new directions get added or removed or both from the PCs' subspace. The newly added directions' variance starts at a small value and slowly stabilizes to the stable value. For deleted directions, we set $(\nu_t)_i = 0$ immediately and we replace f by $f_d < f$ (ensures quicker decay).

¹In future work, we will analyze real data and study existing literature to come up with more realistic models and the corresponding PCs update algorithms. For example, in practice U may not be fixed, but may also rotate gradually over time.

Consider a change time, $t = \tau$. Let $P_{\text{new}} := (U)_{N_\tau \setminus N_{\tau-1}}$ be the matrix containing the k newly added directions. Our PCs update algorithm assumes the following,

- 1) The previous additions are detected and correctly estimated before a new set gets added.
- 2) Let the data matrix D contain τ_d frames of $\hat{L}_t - f\hat{L}_{t-1}$ after the new directions have been added. Then $\text{span}(P_{\text{new}})$ is contained in the span of the data, $\text{span}(D)$.

Assumption 1) holds approximately if d is large enough. Assumption 2) holds with high probability if $\tau_d \gg k$.

We split the estimate of PCs' basis, \hat{P}_t , into two parts, $\hat{P}_t = [\hat{P}_{\text{stable}}, \hat{P}_{\text{new}}]$ where \hat{P}_{stable} is the “stable” (correctly estimated) set of principal directions and \hat{P}_{new} are the new ones which are still being rotated and corrected. We would like to compute an initial estimate of P_{new} as soon as possible (using only a few frames after $\|\hat{\beta}_t\|^2$ exceeds a threshold). Say we use τ_d frames and let the matrix D contains $\hat{L}_{t-1} - f\hat{L}_{t-1}$ for these frames. We can compute an initial estimate of the new directions, \hat{P}_{new} , by computing the principal directions of the sample covariance matrix of $(I - \hat{P}_{\text{stable}}\hat{P}_{\text{stable}}^T)D$. This is done by step 1.b) of Algorithm 1. By assumption 2), if $\tau_d > k$, then we would have found the correct span, i.e. $\text{span}(\hat{P}_{\text{new}}) \supseteq \text{span}(P_{\text{new}})$. But notice that without enough data, even though $\text{span}(\hat{P}_{\text{new}})$ contains $\text{span}(P_{\text{new}})$, it will typically contain many extra directions. As more data comes in, we keep rotating \hat{P}_{new} every-so-often until variances along some directions become approximately zero and these get thresholded out. Once this has happened, the estimated rotation matrix P along the existing directions becomes close to identity and remains this way. This is the time we can add \hat{P}_{new} into \hat{P}_{stable} . This is done by step 1.c) of Algorithm 1.

When the variances along some directions in \hat{P}_{stable} begin to decrease and eventually decay to zero, we compute the variance of last τ_{del} frames along \hat{P}_{stable} and then remove directions with small variance from \hat{P}_{stable} . This is done by step 2) of Algorithm 1.

The above PCs update procedure is summarized in Algorithm 1. In Algorithm 1, D and D_{del} are data matrix to store the data difference $\hat{L}_t - f\hat{L}_{t-1}$. The parameters, τ_d , τ_r , and τ_{del} , are the length of each data piece we use to detect new directions, to rotate and correct newly added directions, and to remove decayed directions, respectively. We use two small

Algorithm 1 Updating \hat{P}_t

1) Detect new appearing directions

a) If status = stable, compute $\|\hat{\beta}_{t-1}\|_2^2 := \|\hat{P}_{t-1,\perp}^T L_{t-1}\|_2^2$. If $\|\hat{\beta}_{t-1}\|_2^2 > \delta$, set status \leftarrow detection and store data in D , i.e., $D \leftarrow [D, \hat{L}_{t-1} - f\hat{L}_{t-2}]$. If not, keep status = stable. go to step 2).

b) If status = detection,

- If there are less than τ_d frames in D , keep storing data difference in D , i.e., $D = [D, \hat{L}_{t-1} - f\hat{L}_{t-2}]$.
- If there are τ_d data frames in D , compute $K = (I - \hat{P}_{\text{stable}}\hat{P}_{\text{stable}}^T)D$.

* Estimate \hat{P}_{new} by computing the EVD of KK^T , i.e.,

$$\frac{1}{\tau_d} KK^T \stackrel{\text{EVD}}{=} PGP^T$$

$$T_d = \{i : G_{i,i} > \xi_d\}$$

$$\hat{P}_{\text{new}} = P_{T_d}, \hat{G}_{\text{new}} = G_{T_d}$$

where G is a square matrix and G_{T_d} is a submatrix of G consisting the rows and columns indexed by T_d .

* Let $D = \emptyset$.

* If $\hat{P}_{\text{new}} = \emptyset$, set status \leftarrow stable and set $l = 0$. If $\hat{P}_{\text{new}} \neq \emptyset$, set status \leftarrow rotation and set $l = \tau_d$.

- Go to step 2)

c) If status = rotation,

- If there are less than τ_r frames in D , keep storing data difference in D , i.e., $D = [D, \hat{L}_{t-1} - f\hat{L}_{t-2}]$.
- If there are τ_r data frames in D , let $K = \hat{P}_{\text{new}}^T D$.

* Rotate \hat{P}_{new} and \hat{G}_{new} using K , i.e.,

$$\frac{1}{l + \tau_r} (l\hat{G}_{\text{new}} + KK^T) \stackrel{\text{EVD}}{=} PGP^T$$

$$T_r = \{i : G_{i,i} > \xi_r\}$$

$$\hat{P}_{\text{new}} = (\hat{P}_{\text{new}}P)_{T_r}, \hat{G}_{\text{new}} = (G)_{T_r}$$

If P is approximately an identity matrix, let $\hat{P}_{\text{stable}} \leftarrow [\hat{P}_{\text{stable}}, \hat{P}_{\text{new}}]$, $\hat{G}_{\text{stable}} \leftarrow [\hat{G}_{\text{stable}}, \hat{G}_{\text{new}}]$. Set status \leftarrow stable, and let $\hat{P}_{\text{new}} = \emptyset$, $\hat{G}_{\text{new}} = \emptyset$, $l = 0$. If not keep status \leftarrow rotation and let $l = l + \tau_d$.

* $D = \emptyset$.

- Go to step 2).

2) Remove decayed directions from \hat{P}_{stable} .

- If there are less than τ_{del} data in D_{del} , keep store data difference in D_{del} , i.e., $D_{\text{del}} = [D_{\text{del}}, \hat{L}_{t-1} - f\hat{L}_{t-2}]$.
- If there τ_{del} data in D_{del} , detect decayed directions as follows

- Find $T_{\text{del}} := \{i : \frac{1}{\tau_{\text{del}}} (\hat{P}_{\text{stable}})_i^T D_{\text{del}} D_{\text{del}}^T (\hat{P}_{\text{stable}})_i < 0.05 (\hat{G}_{\text{stable}})_{i,i}\}$.

- Remove $(\hat{P}_{\text{stable}})_{T_{\text{del}}}$ from \hat{P}_{stable} and remove $(\hat{G}_{\text{stable}})_{T_{\text{del}}}$ from \hat{G}_{stable} .

- Set $D_{\text{del}} = \emptyset$.

3) Let $\hat{P}_t = [\hat{P}_{\text{stable}}, \hat{P}_{\text{new}}]$.

thresholds, ξ_d and ξ_r , to detect new directions \hat{P}_{new} in step 1b) and to threshold extra directions out from \hat{P}_{new} in step 1c). They are proportional to the total variance along all existing stable directions.

E. A complete algorithm

The complete algorithm of real-time robust PCP is given in Algorithm 2.

In Algorithm 2, we compute $\hat{P}_{t,\perp}$, orthogonal complement of \hat{P}_t or equivalently the null space basis of \hat{P}_t^T , using the QR decomposition of \hat{P}_t [22]. Suppose \hat{P}_t is $m \times r$ matrix with $m \gg r$. We find an $m \times m$ orthonormal matrix H such that

$$\hat{P}_t \stackrel{\text{QR}}{=} HJ = [H_1 \ H_2] \begin{bmatrix} J_1 \\ 0 \end{bmatrix}$$

where J_1 is an $r \times r$ upper triangular matrix. H_1 consists of the first r columns of H and H_2 is made up of the last $m - r$ columns. The columns of H_2 span the null space of \hat{P}_t^T and we let $\hat{P}_{t,\perp} = H_2$.

Algorithm 2 Real-time Robust PCP (noise canceled)

Training: Given training data $L_0 = [L_1, \dots, L_{t_0}]$, estimate principal components of L_0 by computing the eigen-pairs of the sample covariance of $L_t - fL_{t-1}$. Let P_0 and G_0 denote the eigenvectors and eigenvalues. Set $\hat{P}_{\text{stable}} = P_0$, $\hat{G}_{\text{stable}} = G_0$.

At time $t = t_0$,

- Set status = stable. Let $\hat{P}_t = \hat{P}_{\text{stable}}$, $\hat{G}_t = \hat{G}_{\text{stable}}$, $\hat{P}_{\text{new}} = \emptyset$, $\hat{G}_{\text{new}} = \emptyset$.
- Let $D = \emptyset$, $l = 0$, $D_{\text{del}} = \emptyset$.

FOR $t > t_0$, do the following:

- 1) Estimate PCs' subspace of low rank part using Algorithm 1 and compute $\hat{P}_{t,\perp}$.
 - 2) Estimate sparse part, S_t , by solving (10) with $\epsilon = 2\|\hat{\beta}_{t-1} - f\hat{\beta}_{t-2}\|_2^2$.
 - 3) Support thresholding and least square estimation: do (13), (14), and (15).
 - 4) Increment t by 1 and go to step 1).
-

IV. EXPERIMENT RESULTS

We simulated $L_t \in \mathbb{R}^{128 \times 1}$ using the model described in Sec.II-A. The first $t_0 = 5 \times 10^3$ frames contains no sparse part and we use it as training data. The sparse vector, x_t , follows a AR-1 model with parameter $f = 0.9$. There are 32 principal directions with variances ranging from 1×10^4 to 9. Recall that in our model, L_t is time correlated and the sequence $L_t - fL_{t-1}$ has same eigenvectors as L_t , we get initial estimate of PCs' subspace by estimating the covariance of $L_t - fL_{t-1}$ and computing its EVD.

The sparse component $S_t \in \mathbb{R}^{128 \times 1}$ first arises at time $t = t_0 + 1$. The nonzero entries of S_t has positive magnitude 5, which is usually much smaller than magnitude of the nonzero entries of x_t . For $t > t_0 + 1$, the support of S_t changes following the model described in Sec.II-B, resulting in a low rank matrix S .

At time $t = t_0 + 5$, we add one new direction P_{new} with variance 50 to PCs' basis and let it starts at a small value

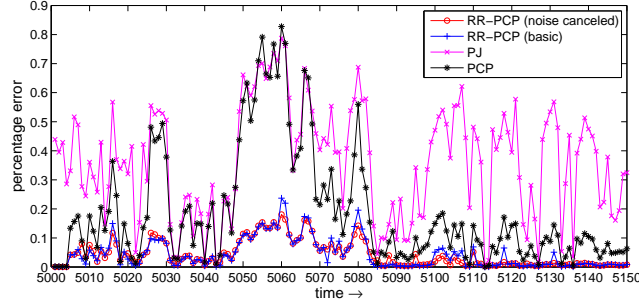


Fig. 2: Comparison of RR-PCP (noise canceled), RR-PCP (basic), PJ and PCP

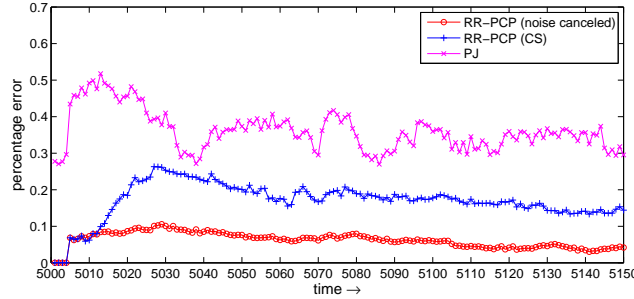


Fig. 3: Monte Carlo averaging for three on-line methods, RR-PCP (noise canceled), RR-PCP (basic), and PJ

with $\theta = 0.4$. It slowly stabilizes to the variance 50. At time $t = t_0 + 100$, one existing direction (not P_{new}) begins to decay with $f_d = 0.1$.

We do RR-PCP (noise canceled), RR-PCP (basic), PJ and PCP using the data generated as described above and plot the percentage error in Fig. 2. The percentage error is defined as

$$\text{percentage error} := \frac{\|S_t - \hat{S}_t\|_2}{\|S_t\|_2}$$

In Fig. 2, $S_{1:t}$ is $1024 \times t$ dimensional at time t , but its rank ranges from 0 to 51.

For RR-PCP (noise canceled), we do algorithm 2 with $\tau_d = \tau_r = \tau_{\text{del}} = 20$. At $t = t_0 + 6$, it detects the appearance of new directions and set status = detection. At $t = t_0 + 26$, a new piece of data containing τ_d frames are available, RR-PCP do step 1b) of Algorithm 1 and get \hat{P}_{new} , an estimate of the new direction P_{new} . There are 7 new directions in \hat{P}_{new} , and the coherence between these new estimated directions and the true one P_{new} ranges from 0.9393 to 0.0051. So, with τ_d frames of data, it approximately finds a subspace containing P_{new} and some extra directions. For $t > t_0 + 26$, we do step 1c) of Algorithm 1 to rotate \hat{P}_{new} closer to the true one and threshold out those extra directions for every $\tau_r = 20$ new frames of data. For example, at $t = t_0 + 47$ when a new piece of data is available, we do step 1c) which rotates \hat{P}_{new} closer to the true P_{new} and get 2 directions thresholded out. The maximum coherence of \hat{P}_{new} and P_{new} goes up to 0.9505. At time $t = t_0 + 68$, another two directions are thresholded out and the maximum coherence of \hat{P}_{new} and P_{new} goes up to 0.9526; at time $t = t_0 + 110$, the rotation matrix P is close to an identity matrix (on-diagonal elements larger than 0.9999 and

off-diagonal elements smaller than 0.01). Only one direction is left in \hat{P}_{new} , with coherence 0.9553 to P_{new} . It sets status = stable and adds \hat{P}_{new} to the stable set of principal directions. At time $t = t_0 + 126$, it removes the deleted direction from the estimated PCs' basis successfully.

For RR-PCP (basic), we do same thing as RR-PCP (noise canceled) but replace (10) with (6) and replace (14) by doing LS on y_t . We see that error of RR-PCP (basic) is larger than RR-PCP (noise canceled) because it does not use the information contained in \hat{L}_{t-1} , i.e. $\|\beta_t\|^2$ is larger than $\|\beta_t - f\hat{\beta}_{t-1}\|^2$ (see Fig.4).

For PJ, we solve (3) with $A = [\hat{P}_t, \hat{P}_{t,\perp}, I]$. PJ recovers x_t and S_t while RR-PCP (noise canceled) and RR-PCP (basic) cancel the term x_t by $\hat{P}_{t,\perp}^T$. Recall that x_t has variance ranging from 1×10^4 to 9, the magnitude of x_t is much larger than S_t . PJ recovers the significant part x_t and cannot get S_t recovered correctly.

For the off-line method PCP, at each time t , we solve (1) using all available data frames², $[M_1, \dots, M_t]$, and plot the error for current frame S_t . The error of PCP is large because the support of S_t is time correlated and S_t does not has random signs. To implement PCP in a causal fashion, it requires about 200 - 300 seconds at every time t , while RR-PCP takes about 1.7 seconds at every time t .

We do 50 times Monte Carlo simulation for three on-line methods, average the percentage error and plot them in Fig. 3. As can be seen from Fig.3, our method RR-PCP (noise canceled), gives the smallest error. In Fig.4, we plot the

²We use Accelerated Proximal Gradient method with code available at http://perception.csl.illinois.edu/matrix-rank/sample_code.html.

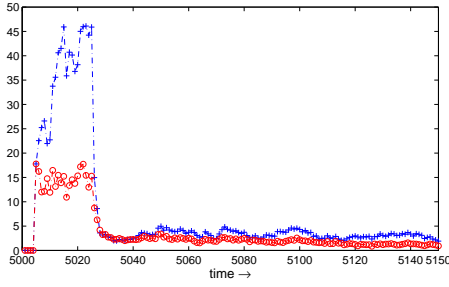


Fig. 4: $\mathbb{E}\|\beta_t\|_2^2$ v.s. $\mathbb{E}\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$

expectation of $\|\beta_t\|_2^2 := \|\hat{P}_{t,\perp}^T L_t\|_2^2$ and $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2 := \|\hat{P}_{t,\perp}^T (L_t - f\hat{L}_{t-1})\|_2^2$ for RR-PCP (noise canceled). It shows that, for $t < t_0 + 26$ when there are some missing principal directions, $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ is much smaller than $\|\beta_t\|_2^2$. For $t > t_0 + 26$, RR-PCP (noise canceled) gets an estimate \hat{P}_{new} and adds it to the estimate of current PCs' basis, thus, both $\|\beta_t\|_2^2$ and $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ decreases. However, $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ is still slightly less than $\|\beta_t\|_2^2$ due to the time correlated model on L_t . Recall that $\|\beta_t\|_2^2$ is the noise in (6) and $\|\beta_t - f\hat{\beta}_{t-1}\|_2^2$ is the noise in (10), that is the reason why RR-PCP (noise canceled) is better than RR-PCP (basic).

V. DISCUSSION AND FUTURE WORK

In this work, we used a simple motion model on the sparse vector S_t as explained in Sec. II-B. Under this model, the support of S_t changes slowly over time, resulting in a low rank matrix S . Because of this, PCP is unable to distinguish S from the low rank L . But our method, RR-PCP, works because it does not require the sparse matrix S to be uniformly random. In this work, we have not utilized the correlated support change of S_t to our advantage. But, in fact, RR-PCP can be improved significantly by using this knowledge and by adapting the modified-CS idea of [19], [20] to incorporate motion prediction. We can use the knowledge of the object's motion model and the previous support estimate to obtain the current support prediction T_{pred} of the sparse part. If this prediction is accurate and is used in (16), with an appropriately chosen ϵ , the reconstruction error should reduce significantly, especially when the support size of S_t is large. In future work, we will develop realistic motion models and corresponding motion prediction algorithms to get reliable support predictions of the sparse part. We will also analyze their performance, first assuming P_t is perfectly known and later for the practical case of P_t unknown.

Our PCs updating procedure is designed for the data generated according to the piecewise stationary model on x_t while U is a constant but unknown orthonormal matrix. In future work, we will analyze real data and study existing literature to come up with more realistic models and the corresponding PC update algorithms.

For very-large scale data, it is computationally and memory intensive to compute $\hat{P}_{t,\perp}$. In future work, we will develop computational efficient alternatives. For example we can use the fact that $\|\hat{P}_{t,\perp}^T z\|_2 = \|\hat{P}_{t,\perp} \hat{P}_{t,\perp}^T z\|_2 = \|(I - \hat{P}_t \hat{P}_t^T) z\|_2$.

A somewhat related work is Jin-Rao's approach [14] which solves

$$\min_{\alpha, s} \|s\|_1 \text{ s.t. } \|M_t - P_t \alpha - s\|_2^2 \leq \epsilon \quad (17)$$

In [14], the matrix P_t is a known and fixed regression coefficients' matrix, which is no longer true in our problem. We can use the time correlated model on x_t (and hence on α_t) and the motion model on S_t to modify (17) following a similar way of RR-PCP.

REFERENCES

- [1] P. Huber, *Robust statistics*, Wiley and Sons, 1981.
- [2] Huan Xu, C. Caramanis, and S. Mannor, "High dimensional principal component analysis with contaminated data", in *Networking and Information Theory, 2009. ITW 2009. IEEE Information Theory Workshop on*, jun. 2009, pp. 246–250.
- [3] Sam Roweis, "Em algorithms for pca and spca", in *Advances in Neural Information Processing Systems*, 1998, pp. 626–632, MIT Press.
- [4] F. De La Torre and M. Black, "A framework for robust subspace learning", in *International Journal on Computer vision*, 2003, pp. 54:117–142.
- [5] Matthew Brand, "Incremental singular value decomposition of uncertain data with missing values", in *In ECCV*, 2002, pp. 707–720.
- [6] Yongmin Li, Li qun Xu, Jason Morphet, and Richard Jacobs, "An integrated algorithm of incremental and robust pca", in *Proceedings of IEEE International Conference on Image Processing*, 2003, pp. 245–248.
- [7] Lei Xu and A. Yuille, "Robust principal component analysis by self-organizing rules based on statistical physics approach", in *IEEE Trans. Neural Networks*, 1995, pp. 131–143.
- [8] D. Skocaj and A. Leonardis, "Weighted and robust incremental method for subspace learning", oct. 2003, pp. 1494–1501 vol.2.
- [9] F. De la Torre and M.J. Black, "Robust principal component analysis for computer vision", 2001, vol. 1, pp. 362–369 vol.1.
- [10] Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright, "Robust principal component analysis?", *Submitted for publication*.
- [11] Venkat Chandrasekaran, Sujay Sanghavi, Pablo A. Parrilo, and Alan S. Willsky, "Sparse and low-rank matrix decompositions", in *Allerton'09*, 2009.
- [12] A. Ganesh, J. Wright, X. Li, E. J. Candes, and Y. Ma, "Dense Error Correction for Low-Rank Matrices via Principal Component Pursuit", *ArXiv e-prints*, January 2010.
- [13] E. Candes and T. Tao, "Decoding by linear programming", *IEEE Trans. Info. Th.*, vol. 51(12), pp. 4203–4215, Dec. 2005.
- [14] Yuzhe Jin and Bhaskar Rao, "Algorithms for robust linear regression by exploiting the connection to sparse signal recovery", in *IEEE Intl. Conf. Acous. Speech. Sig.Proc.(ICASSP)*, 2010.
- [15] Kaushik Mitra, Ashok Veeraraghavan, and Rama Chellappa, "A robust regression using sparse learning for high dimensional parameter estimation problems", in *IEEE Intl. Conf. Acous. Speech. Sig.Proc.(ICASSP)*, 2010.
- [16] John Wright and Yi Ma, "Dense error correction via l1-minimization", *IEEE Transactions on Information Theory*, 2009.
- [17] J.N. Laska, M.A. Davenport, and R.G. Baraniuk, "Exact signal recovery from sparsely corrupted measurements through the pursuit of justice", in *Signals, Systems and Computers, 2009 Conference Record of the Forty-Third Asilomar Conference on*, Nov 2009, pp. 1556–1560.
- [18] E. Candes and T. Tao, "The dantzig selector: statistical estimation when p is much larger than n", *Annals of Statistics*, 2006.
- [19] Wei Lu and Namrata Vaswani, "Modified basis pursuit denoising (modified-bpdn) for noisy compressive sensing with partially known support", in *IEEE Intl. Conf. Acous. Speech. Sig.Proc.(ICASSP)*, 2010.
- [20] N. Vaswani and W. Lu, "Modified-cs: Modifying compressive sensing for problems with partially known support", *IEEE Trans. Signal Processing*, September 2010.
- [21] N. Vaswani and Wei Lu, "Modified-cs: Modifying compressive sensing for problems with partially known support", in *ISIT 2009*, June 2009, pp. 488–492.
- [22] R. A. Horn and C. R. Johnson, *Matrix Analysis*, Cambridge University Press, February 1990.