

RESEARCH ARTICLE

Sparse Coding and Counting for Robust Visual Tracking

Risheng Liu^{1,2*}, Jing Wang³, Xiaoke Shang⁴, Yiyang Wang³, Zhixun Su³, Yu Cai³

1 School of Software Technology, Dalian University of Technology, Dalian City, Liaoning Province, China, **2** Key Laboratory for Ubiquitous Network and Service Software of Liaoning Province, Dalian City, Liaoning Province, China, **3** School of Mathematic Sciences, Dalian University of Technology, Dalian City, Liaoning Province, China, **4** Dalian Campus, Luxun Academy of Fine Arts, Dalian City, Liaoning Province, China

* rsliu@dlut.edu.cn



Abstract

In this paper, we propose a novel sparse coding and counting method under Bayesian framework for visual tracking. In contrast to existing methods, the proposed method employs the combination of L_0 and L_1 norm to regularize the linear coefficients of incrementally updated linear basis. The sparsity constraint enables the tracker to effectively handle difficult challenges, such as occlusion or image corruption. To achieve real-time processing, we propose a fast and efficient numerical algorithm for solving the proposed model. Although it is an NP-hard problem, the proposed accelerated proximal gradient (APG) approach is guaranteed to converge to a solution quickly. Besides, we provide a closed solution of combining L_0 and L_1 regularized representation to obtain better sparsity. Experimental results on challenging video sequences demonstrate that the proposed method achieves state-of-the-art results both in accuracy and speed.

OPEN ACCESS

Citation: Liu R, Wang J, Shang X, Wang Y, Su Z, Cai Y (2016) Sparse Coding and Counting for Robust Visual Tracking. PLoS ONE 11(12): e0168093. doi:10.1371/journal.pone.0168093

Editor: Quan Zou, Tianjin University, CHINA

Received: January 10, 2016

Accepted: November 24, 2016

Published: December 16, 2016

Copyright: © 2016 Liu et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All relevant data are within the paper.

Funding: This work is partially supported by the National Natural Science Foundation of China (Nos. 61672125, 61300086, 61432003), the Fundamental Research Funds for the Central Universities (DUT15QY15), the Hong Kong Scholar Program (No. XJ2015008), and National Science and Technology Major Project (Nos. 2013ZX04005-021, 2014ZX001011). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Introduction

Visual tracking plays an important role in computer vision and has many applications such as video surveillance, robotics, motion analysis and human computer interaction. Even though various algorithms have come out, it is still a challenge problem due to complex object motion, heavy occlusion, illumination change and background clutter.

Visual tracking algorithms can be roughly categorized into two major categories: discriminative methods and generative methods. Discriminative methods (e.g., [1–3]) view object tracking as a binary classification problem in which the goal is to separate the target object from the background. Generative methods (e.g., [4–8]) employ a generative appearance model to represent the target's appearance.

We focus on the generative one and will briefly review the relevant work below. Recently, sparse representation has been successfully applied to visual tracking (e.g., [9–12]). The trackers based on sparse representation are under the assumption that the appearance of a tracked object can be sparsely represented by a over-complete dictionary which can be dynamically updated to maintain holistic appearance information. Traditionally, the over-complete dictionary is a series of redundant object templates, however, a set of basis vectors from target

Competing Interests: The authors have declared that no competing interests exist.

subspace as dictionary is also used because an orthogonal dictionary performs as efficient as the redundant one. In visual tracking, we will call the L_1 regularized object representation “sparse coding” (e.g., [9]), and the L_0 regularized object representation “sparse counting” (e.g., [13]). [9] has been shown to be robust against partial occlusions, which improves the tracking performance. However, because of using redundant dictionary, heavy computational overhead in L_1 minimization hampers the tracking speed. Very recent efforts have been made to improve this method in terms of both speed and accuracy by using accelerated proximal gradient (APG) algorithm [14] or modeling the similarity between different candidates [11]. Different from [9], IVT [5] incrementally learns a low-dimensional PCA subspace representation, which adapts online to the appearance changes of the target. To get rid of image noise, Lu *et al.* [15] introduce L_1 noise regularization into the PCA reconstruction, which is able to handle partial occlusion and other challenging factors. Pan *et al.* [13] employs L_0 norm to regularize the linear coefficients of incrementally updated linear basis (sparse counting) to remove the redundant features of the basis vectors. However, sparse counting will cause unstable solutions because of its nonconvexity and discontinuity. Although the sparse coding has good performance, it may cause biased estimation since it penalizes true large coefficients more, and produce over-penalization. Consequently, it is necessary to find a way to overcome the disadvantages of sparse coding and sparse counting.

From the viewpoint of statistics, sparse representation are similar to variable selection when the dictionary is fixed. Besides, it is a bonus that Bayesian framework has been successfully applied to select variables by enforcing appropriate priors. Laplace priors were used to avoid overfitting and enforce sparsity in sparse linear model, which derives sparse coding problem. To further enforce sparsity and reduce over-penalization of sparse coding, each coefficient is assigned with a Bernoulli variable. Therefore, a novel model interpreted from a Bayesian perspective by carrying maximum a posteriori (MAP) is proposed, which turns out to be a combination of sparse coding and counting model. In paper [16], Lu *et al.* also consider L_0 and L_1 norm under a Bayesian perspective. However, considering that there will be occlusion, illumination change and background clutter in tracking, we restraint the noise with L_1 norm. Besides, We use an orthogonal dictionary to replace the redundant object templates as similar atoms of redundant templates may cause mistake of coefficients and huge computational complexity. Lastly, We propose closed solution of regularization which is the combination of the L_0 norm and L_1 norm. However Lu *et al.* obtain the approximate solution by using the Greedy Coordinate Descent.

Tracking results by using unconstrained regularization, sparse counting, sparse coding and our model under the same dictionary D are shown in Fig 1, respectively. As shown in Fig 1, one can see that the coefficients of unconstrained regularization and sparse coding are actually not sparse and the target object is not tracked well. Similarly, sparse counting with sparsity coefficients sometimes cannot obtain appropriate linear combination of the orthogonal basis vectors, which will interfere with the tracking accuracy. However, we note that our method is able to reconstruct the object well and find the good candidate, then facilitating the tracking results. We also compare our model with unconstrained regularization, sparse counting, sparse coding over all 50 sequences in benchmark, the precision and success plots are shown in Fig 2. One can see the parameter setting in the section Experimental Results.

Contributions: The contributions of this work are threefold.

- We propose a sparse coding and counting model from a novel Bayesian perspective for visual tracking. Compared to the state-of-the-art algorithms, the proposed method achieves more reliable tracking results.

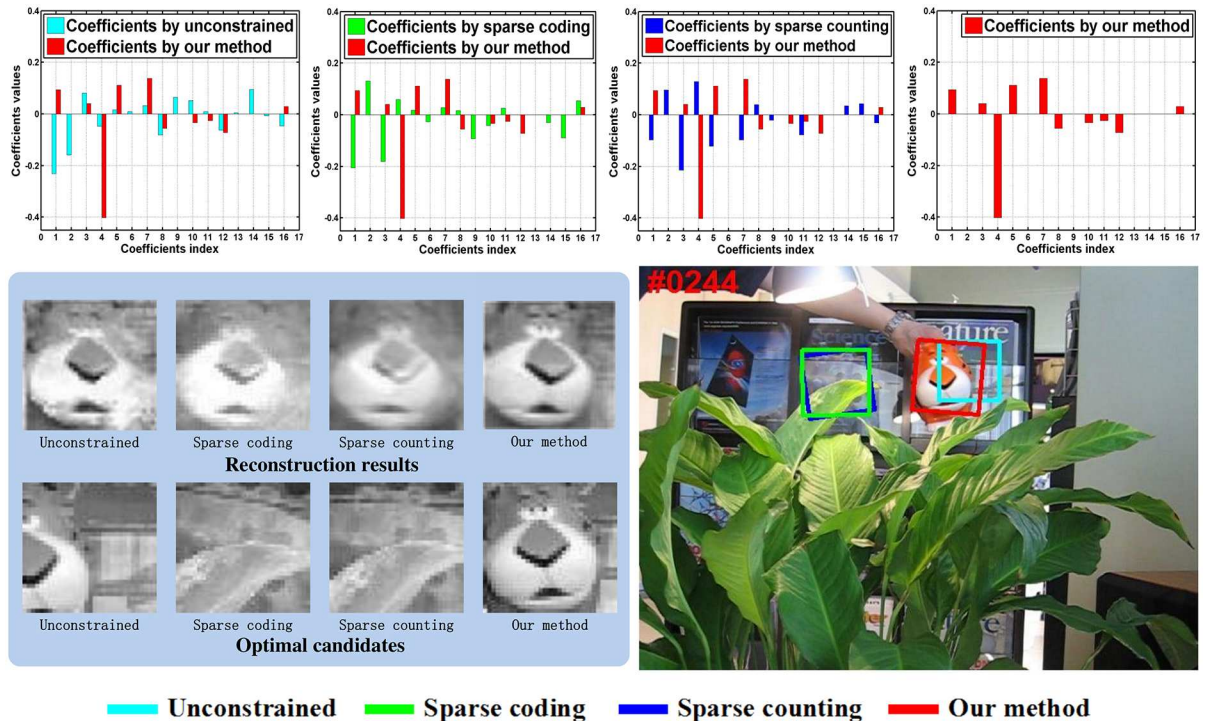


Fig 1. The comparison of coefficients, optimal candidates and reconstruction. The top is the coefficients of our method versus unconstrained, sparse coding and sparse counting regularization, respectively. The bottom is the optimal candidates and reconstruction results by using unconstrained, sparse coding, sparse counting and our method under same dictionary, respectively.

doi:10.1371/journal.pone.0168093.g001

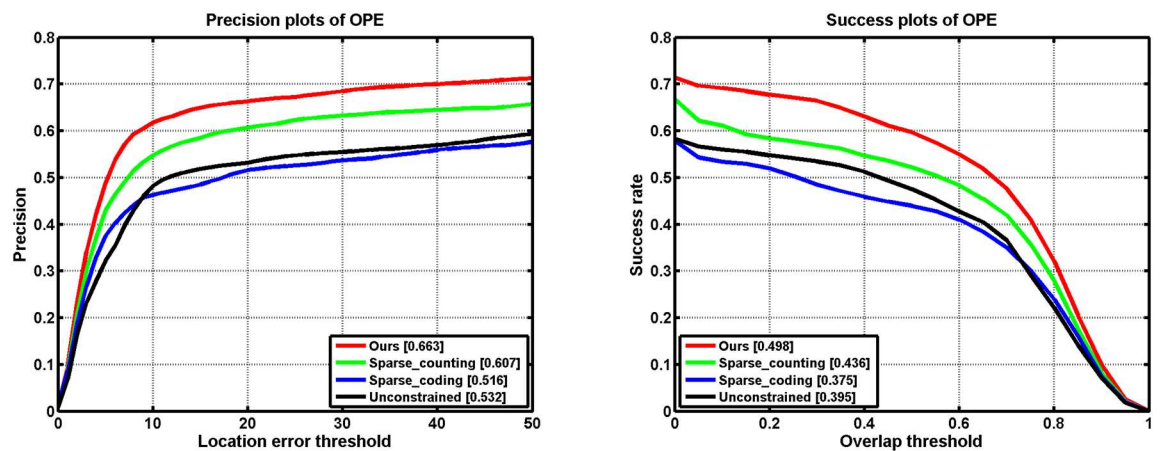


Fig 2. Precision and success plots of overall performance comparison among unconstrained regularization, sparse counting, sparse coding and ours for the 22 videos in the benchmark. The mean precision scores are reported in the legends.

doi:10.1371/journal.pone.0168093.g002

- We propose closed solution of combining the L_0 norm and L_1 norm based regularization in a unique one.
- Although the sparse coding and counting related minimization is an NP-hard problem, we show that the proposed model can be efficiently estimated by the proposed APG method. This makes our tracking method computationally attractive in general and comparable in speed with SP method [15] and the accelerated L_1 tracker [14].

Visual Tracking based on the Particle Filter

In this paper, we employ a particle filter to track the target object. The particle filter provides an estimate of posterior distribution of random variables related to Markov chain. Given a set of observed image vectors $\mathbf{Y}_t = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t\}$ up to the t -th frame and target state variable \mathbf{x}_t that describes the six affine motion parameters, the posterior distribution $p(\mathbf{x}_t|\mathbf{Y}_t)$ based on the Bayesian theorem is estimated by:

$$p(\mathbf{x}_t|\mathbf{Y}_t) \propto p(\mathbf{y}_t|\mathbf{x}_t) \int p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{Y}_{t-1})d\mathbf{x}_{t-1}, \tag{1}$$

where $p(\mathbf{y}_t|\mathbf{x}_t)$ is the observation model that estimates the likelihood of an observed image patch \mathbf{y}_t belonging to the object class, and $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ is the motion model that describes the state transition between consecutive frames.

The Motion Model: The motion model $p(\mathbf{x}_t|\mathbf{x}_{t-1}) = N(\mathbf{x}_t; \mathbf{x}_{t-1}, \Sigma)$ models the parameters by independent Gaussian distribution around the counterpart in \mathbf{x}_{t-1} , where Σ is a diagonal covariance matrix whose elements are the variances of the affine parameters. In the tracking framework, the optimal target state $\hat{\mathbf{x}}_t$ is obtained by the maximal approximate posterior (MAP) probability: $\hat{\mathbf{x}}_t = \arg \max_{\mathbf{x}_t^i} p(\mathbf{x}_t^i|\mathbf{Y}_t)$, where \mathbf{x}_t^i indicates the i -th sample of the state \mathbf{x}_t .

The observation model: In this paper, we assume that the tracked target object is generated by a subspace (spanned by \mathbf{D} and centered at $\boldsymbol{\mu}$) with corruption (i.i.d Gaussian Laplacian noise),

$$\mathbf{y} = \mathbf{D}\boldsymbol{\alpha} + \epsilon + \mathbf{e}, \tag{2}$$

where $\mathbf{y} \in \mathbb{R}^N$ denotes an observation vector centered at $\boldsymbol{\mu}$, the columns of $\mathbf{D} = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_K\} \in \mathbb{R}^{N \times K}$ are orthogonal basis vectors of the subspace, $\boldsymbol{\alpha}$ indicates the coefficients of basis vectors, ϵ and \mathbf{e} stand for the Gaussian noise and Laplacian noise vector respectively. the Gaussian component models small dense noise and the Laplacian one aims to handle outliers. As proposed by [17], under the i.i.d Gaussian-Laplacian noise assumption, the distance between the vector \mathbf{y} and the subspace $(\mathbf{D}, \boldsymbol{\mu})$ is the least soft threshold squares distance:

$$d(\boldsymbol{\alpha}, \mathbf{e}) = \min_{\boldsymbol{\alpha}, \mathbf{e}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\boldsymbol{\alpha} - \mathbf{e}\|_2^2 + \lambda \|\mathbf{e}\|_1. \tag{3}$$

Thus, for each observation \mathbf{y}_t corresponding to a predicted state \mathbf{x}_t , the observation model $p(\mathbf{y}_t|\mathbf{x}_t)$ that is set to be

$$p(\mathbf{y}_t|\mathbf{x}_t) = \exp(-\tau d(\boldsymbol{\alpha}^*, \mathbf{e}^*)), \tag{4}$$

where $\boldsymbol{\alpha}^*$ and \mathbf{e}^* are the optimal solution of Eq (18) which will be introduced in detail in next section, and τ is a constant controlling the shape of the Gaussian kernel.

Model Update: It is essential to update the observation model for handling appearance change of the target in visual tracking. Since the error term \mathbf{e} can be used to identify some outliers (e.g., Laplacian noise, illumination), we adopt the strategy proposed by [17] to update the

appearance model using the incremental PCA with mean update [5] as follows,

$$y_i = \begin{cases} y_i, & e_i = 0, \\ \mu_i, & \text{otherwise,} \end{cases} \quad (5)$$

where y_i , e_i , and μ_i are the i -th elements of \mathbf{y} , \mathbf{e} , and $\boldsymbol{\mu}$, respectively, $\boldsymbol{\mu}$ is the mean vector computed the same as [5].

Object Representation by Bayesian Framework

Motivation

Considering \mathbf{y} as the vectorized target object region, it can be represented by an feature subspace with both sparse corruptions and dense errors, i.e.,

$$\mathbf{y} = \mathbf{D}\boldsymbol{\alpha} + \boldsymbol{\epsilon} + \mathbf{e}. \quad (6)$$

Most existing sparsity based trackers aim to directly utilize L_1 regularization on $\boldsymbol{\alpha}$ to suppress small coefficients for subspace reconstruction. However, by carefully investigating the soft-thresholding operator corresponding to L_1 minimization subproblem, it can be observed that such simple regularization will consistently suppress the values of the coefficients, thus destroy the discriminative property of the learned feature subspace.

To address this limitation in existing work, we here incorporate two different sparse regularization techniques within the Bayesian perspective, which has the capacity to encode prior knowledge and to make valid estimation of uncertainty. In other words, our goal is to propose a Bayesian inference framework to incorporate both the coefficients thresholding and selection to improve the discrimination of our feature subspace learning formulation. Specifically, by defining an index vector $\mathbf{r} = [r_1, r_2, \dots, r_K]$ ($r_l = \mathbb{I}(\boldsymbol{\alpha}_l \neq 0)$, $l = 1, 2, \dots, K$), Eq (6) can be rewritten as

$$y_j = \sum_{l=1}^K d_{jl} r_l \boldsymbol{\alpha}_l + \epsilon_j + e_j, \quad j = 1, 2, \dots, N. \quad (7)$$

Here the additional index vector \mathbf{r} can be considered as a dictionary selection operator and we will enforce particular prior distribution on it to enhance the discriminative power of our model for subspace reconstruction. To further enhance the representative ability of our model, we will also develop a novel dictionary learning framework to build orthogonal subspace dictionary for Eq (6). Please notice that the orthogonality of the learned dictionary will also significantly simplify the numerical optimization process. Please see the following sections for more details.

Bayesian Formulation

Now we will introduce our model under Bayesian framework in detail. The joint posterior distribution of $\boldsymbol{\alpha}$, \mathbf{r} , \mathbf{e} and σ^2 based on the Bayesian theorem can be written as

$$p(\boldsymbol{\alpha}, \mathbf{r}, \mathbf{e}, \sigma^2 | \mathbf{D}, \mathbf{y}, \tilde{\boldsymbol{\mu}}, \tau_1, \tau_2, \kappa, \hat{\sigma}) \propto p(\mathbf{y} | \mathbf{D}, \boldsymbol{\alpha}, \mathbf{r}, \mathbf{e}, \sigma^2) p(\boldsymbol{\alpha} | \sigma^2, \tilde{\boldsymbol{\mu}}) p(\mathbf{r} | \kappa) p(\mathbf{e} | \hat{\sigma}) p(\sigma^2 | \tau_1, \tau_2), \quad (8)$$

where $p(\mathbf{y} | \mathbf{D}, \boldsymbol{\alpha}, \mathbf{r}, \mathbf{e}, \sigma^2)$, $p(\boldsymbol{\alpha} | \sigma^2, \tilde{\boldsymbol{\mu}})$, $p(\mathbf{r} | \kappa)$, $p(\mathbf{e} | \hat{\sigma})$, $p(\sigma^2 | \tau_1, \tau_2)$, denote the priors on the noisy vectorized target region, the coefficient vector $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_K]$, the index vector $\mathbf{r} = [r_1, r_2, \dots, r_K]$ ($r_l = \mathbb{I}(\boldsymbol{\alpha}_l \neq 0)$, $l = 1, 2, \dots, K$), the Laplacian noise, and the noise level, respectively. In Eq (8), the parameters $\tilde{\boldsymbol{\mu}}$, τ_1 , τ_2 , $\hat{\sigma}$, and κ are the relevant constant parameters of the priors.

We generally assume that the noise ϵ_j follows the Gaussian distribution, *i.e.*, $p(\epsilon_j) = N(0, \sigma^2)$. We treat the Laplacian noise term e_j as missing values with the same Laplacian prior. Therefore, the Prior $p(\mathbf{y}|\mathbf{D}, \mathbf{a}, \mathbf{r}, \mathbf{e}, \sigma^2)$ has the follow distribution:

$$p(\mathbf{y}|\mathbf{D}, \mathbf{a}, \mathbf{r}, \mathbf{e}, \sigma^2) = \prod_{j=1}^N P(y_j|\mathbf{d}_j, \mathbf{a}, \mathbf{r}, e_j, \sigma^2) = \prod_{j=1}^N N\left(\sum_{l=1}^K d_{jl}r_l\mathbf{a}_l + e_j, \sigma^2\right). \tag{9}$$

To enforce sparsity, the coefficients \mathbf{a} are assumed to follow Laplace distribution.

$$p(\mathbf{a}|\sigma^2, \tilde{\mu}) = \prod_{l=1}^K p(\mathbf{a}_l|\sigma^2, \tilde{\mu}) = \prod_{l=1}^K \frac{1}{2\sigma^2\tilde{\mu}^{-1}} \exp\left(-\frac{|\mathbf{a}_l|}{\sigma^2\tilde{\mu}^{-1}}\right). \tag{10}$$

Our goal is to remove redundant features while preserving the useful parts in the dictionary. As Laplace priors resulting sparse coding may lead to over penalization on the large coefficients, we assume the index variable r_l of each coefficient \mathbf{a}_l to be a Bernoulli variable to enforce sparsity and reduce over penalization.

$$p(\mathbf{r}|\kappa) = \prod_{l=1}^K \kappa^{r_l} (1 - \kappa)^{1-r_l}, \tag{11}$$

where $\kappa \leq 1/2$. Here, the Bernoulli prior on r_l means that r_l will have probability κ to be 1 and $1 - \kappa$ to be 0, if the prior information is known.

The noise e_j is aims at handling outliers, so it follows Laplace distribution:

$$p(\mathbf{e}|\hat{\sigma}) = \prod_{j=1}^N p(e_j|\hat{\sigma}) = \prod_{j=1}^N \frac{1}{2\hat{\sigma}} \exp\left(-\frac{|e_j|}{\hat{\sigma}}\right). \tag{12}$$

The variances of noises are assigned with Inverse Gamma prior as follow:

$$p(\sigma^2|\tau_1, \tau_2) = \frac{\tau_2^{\tau_1}}{\Gamma(\tau_1)} \sigma^{-2(\tau_1+1)} \exp\left(-\frac{\tau_2}{\sigma^2}\right), \tag{13}$$

where $\Gamma(\cdot)$ denotes the gamma function.

Then, the optimal $\mathbf{a}, \mathbf{r}, \mathbf{e}, \sigma^2$ are obtained by the MAP probability. After taking the negative logarithm, the formula is

$$(\mathbf{a}^*, \mathbf{r}^*, \mathbf{e}^*, \sigma^{*2}) = \arg \min_{\mathbf{a}, \mathbf{r}, \mathbf{e}, \sigma^2} \{-2 \log p(\mathbf{a}, \mathbf{r}, \mathbf{e}, \sigma^2|\mathbf{D}, \mathbf{y}, \tilde{\mu}, \tau_1, \tau_2, \kappa, \hat{\sigma})\}. \tag{14}$$

Combining the aforementioned Eqs (8)–(13), we have

$$\begin{aligned} & -2 \log p(\mathbf{a}, \mathbf{r}, \mathbf{e}, \sigma^2|\mathbf{D}, \mathbf{y}, \tilde{\mu}, \tau_1, \tau_2, \kappa, \hat{\sigma}) \\ &= \frac{1}{\sigma^2} \sum_{j=1}^N \left(y_j - \sum_{l=1}^K d_{jl}r_l\mathbf{a}_l - e_j\right)^2 + \frac{1}{\sigma^2} \frac{2\sigma^2}{\hat{\sigma}} \sum_{j=1}^N |e_j| + \frac{2\tilde{\mu}}{\sigma^2} \sum_{l=1}^K |\mathbf{a}_l| + \\ & (2N + 2K + 2\tau_1 + 2) \log \sigma^2 + \frac{2\tau_2}{\sigma^2} + \sum_{l=1}^K r_l \log \frac{(1 - \kappa)^2}{\kappa^2} + const. \end{aligned} \tag{15}$$

With fixing $\sigma^2 = 1$, Eq (15) can be rewritten as

$$\|\mathbf{y} - \mathbf{D}\mathbf{a} - \mathbf{e}\|_2^2 + 2\beta\|\mathbf{e}\|_1 + 2\tilde{\mu}\|\mathbf{a}\|_1 + \rho_\kappa\|\mathbf{a}\|_0 + const, \tag{16}$$

where $\rho_\kappa = \log(1 - \kappa)^2/\kappa^2$, $\beta = \sigma^2/\hat{\sigma}$. With $\gamma \in [0, 1]$, $\lambda = \tilde{\mu} + 1/2\rho_\kappa$ and

$\gamma = 4\tilde{\mu}/(2\tilde{\mu} + \rho_\kappa)$, Eq (16) can be rewritten as

$$\frac{1}{2} \| \mathbf{y} - \mathbf{D}\boldsymbol{\alpha} - \mathbf{e} \|_2^2 + \beta \| \mathbf{e} \|_1 + \lambda (\gamma \| \boldsymbol{\alpha} \|_1 + (1 - \gamma) \| \boldsymbol{\alpha} \|_0) + \text{const.} \tag{17}$$

Final Optimization Model

By observing the objective function in Eq (17), it can be found that the essential regularization in Eq (17) is a combination of the sparse coding and the sparse counting. With a fixed appropriate orthogonal dictionary \mathbf{D} , Eq (17) can be written as the following optimization problem

$$\min_{\boldsymbol{\alpha}, \mathbf{e}} \frac{1}{2} \| \mathbf{y} - \mathbf{D}\boldsymbol{\alpha} - \mathbf{e} \|_2^2 + \beta \| \mathbf{e} \|_1 + \lambda (\gamma \| \boldsymbol{\alpha} \|_1 + (1 - \gamma) \| \boldsymbol{\alpha} \|_0), \tag{18}$$

where $\| \cdot \|_0$ denotes the L_0 norm which counts the number of non-zero elements, γ , λ and β are regularization parameters, and $\| \cdot \|_2$ and $\| \cdot \|_1$ denote L_2 and L_1 norms, respectively. The term $\| \mathbf{e} \|_1$ is used to reject outliers (e.g., occlusions), while $\| \boldsymbol{\alpha} \|_0$ and $\| \boldsymbol{\alpha} \|_1$ are used to select the most discriminative subspace features. Notice that we also implicitly assume that $\mathbf{D}^\top \mathbf{D} = \mathbf{I}$, where \mathbf{I} is an identity matrix.

Theory of Fast Numerical Algorithm

It is known that APG is an excellent algorithm for convex programming [18, 19] and has been used in visual tracking. In this section, we propose a fast numerical algorithm for solving the proposed nonconvex and nonsmooth model by using APG approach. The experimental results show that it can converge to a solution quickly and achieve attractive performance. Besides, the closed solution of the combining L_0 and L_1 based regularization is provided.

APG Algorithm for Solving Eq (19)

Eq (18) contains two subproblem: one is solving $\boldsymbol{\alpha}$ given fixed \mathbf{e} , the other one is solving \mathbf{e} given fixed $\boldsymbol{\alpha}$, the formula is shown as follow

$$\begin{cases} \boldsymbol{\alpha} = \arg \min_{\boldsymbol{\alpha}} \frac{1}{2} \| \mathbf{y} - \mathbf{D}\boldsymbol{\alpha} - \mathbf{e} \|_2^2 + \lambda \gamma \| \boldsymbol{\alpha} \|_1 + \lambda (1 - \gamma) \| \boldsymbol{\alpha} \|_0, \\ \mathbf{e} = \arg \min_{\mathbf{e}} \frac{1}{2} \| \mathbf{y} - \mathbf{D}\boldsymbol{\alpha} - \mathbf{e} \|_2^2 + \beta \| \mathbf{e} \|_1. \end{cases} \tag{19}$$

Solving Eq (19) is an NP-hard problem because it involves a discrete counting metric. We adopt a special optimization strategy based on the APG approach [18], which ensures each step be solved easily. In APG Algorithm, we need to solve

$$\begin{cases} \boldsymbol{\alpha}_{k+1}^* = \arg \min_{\boldsymbol{\alpha}} \lambda \gamma \| \boldsymbol{\alpha} \|_1 + \lambda (1 - \gamma) \| \boldsymbol{\alpha} \|_0 + \frac{L}{2} \| \boldsymbol{\alpha} - \mathbf{z}_{k+1}^\alpha + \frac{\nabla_{\boldsymbol{\alpha}} F(\mathbf{z}_{k+1})}{L} \|_2^2, \\ \mathbf{e}_{k+1}^* = \arg \min_{\mathbf{e}} \beta \| \mathbf{e} \|_1 + \frac{L}{2} \| \mathbf{e} - \mathbf{z}_{k+1}^e + \frac{\nabla_{\mathbf{e}} F(\mathbf{z}_{k+1})}{L} \|_2^2, \end{cases} \tag{20}$$

where $\mathbf{z}_{k+1} = (\mathbf{z}_{k+1}^\alpha, \mathbf{z}_{k+1}^e)$, $\nabla_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}, \mathbf{e}) = \mathbf{D}^\top (\mathbf{D}\boldsymbol{\alpha} + \mathbf{e} - \mathbf{y})$, $\nabla_{\mathbf{e}} F(\boldsymbol{\alpha}, \mathbf{e}) = \mathbf{e} - (\mathbf{y} - \mathbf{D}\boldsymbol{\alpha})$, and L is a Lipschitz constant.

The solutions of Eq (20) can be obtained by

$$\begin{cases} \alpha_{k+1}^* = \mathcal{E}_{(\lambda\gamma/L, \lambda(1-\gamma)/L)}\left(\mathbf{z}_{k+1}^a - \frac{\nabla_{\alpha} F(\mathbf{z}_{k+1})}{L}\right), \\ \mathbf{e}_{k+1}^* = \mathcal{S}_{\beta/L}\left(\mathbf{z}_{k+1}^e - \frac{\nabla_{\mathbf{e}} F(\mathbf{z}_{k+1})}{L}\right), \end{cases} \quad (21)$$

where $\mathcal{S}_{\theta}(y) = \text{sign}(y) \max(|y| - \theta, 0)$, and $\mathcal{E}_{(\delta, \eta)}(y)$ is defined as

$$\mathcal{E}_{(\delta, \eta)}(y) = \begin{cases} y - \delta, & y > \delta + \sqrt{2\eta}, \\ y + \delta, & y < -\delta - \sqrt{2\eta}, \\ 0, & \text{otherwise.} \end{cases} \quad (22)$$

The numerical algorithm for solving Eq (19) is summarized in Algorithm 1. Due to the orthogonality of \mathbf{D} , Algorithm 1 converges fast, and its computation cost does not increase compared to the solver of L_1 regularized model.

Algorithm 1 Fast Numerical Algorithm for Solving Eq (19)

Initialize: Set initial guesses $\alpha_0 = \alpha_{-1} = \mathbf{0}$, $\mathbf{e}_0 = \mathbf{e}_{-1} = \mathbf{0}$, and $t_0 = t_{-1} = 1$.
while not convergence or termination **do**
Step 1: $\mathbf{z}_{k+1}^a := \alpha_k + \frac{t_{k-1}-1}{t_k}(\alpha_k - \alpha_{k-1})$;
Step 2: $\mathbf{z}_{k+1}^e := \mathbf{e}_k + \frac{t_{k-1}-1}{t_k}(\mathbf{e}_k - \mathbf{e}_{k-1})$;
Step 3: $\alpha_{k+1} = \mathcal{E}_{(\lambda\gamma/L, \lambda(1-\gamma)/L)}\left(\mathbf{z}_{k+1}^a - \frac{\nabla_{\alpha} F(\mathbf{z}_{k+1})}{L}\right)$;
Step 4: $\mathbf{e}_{k+1} = \mathcal{S}_{\beta/L}\left(\mathbf{z}_{k+1}^e - \frac{\nabla_{\mathbf{e}} F(\mathbf{z}_{k+1})}{L}\right)$;
Step 5: $t_{k+1} := \frac{1 + \sqrt{1 + 4t_k^2}}{2}$, $k \leftarrow k + 1$.
end while

Closed-form Solution for Combining L_1 and L_0 Regularization

This subsection mainly focus on a sparse combinatory model which combines L_0 and L_1 norm together as the regularizer term

$$\min_x \frac{1}{2}(x - y)^2 + \delta|x| + \eta|x|_0, \quad (23)$$

where $x, y \in \mathbb{R}^1$, and $|x|$ denotes L_0 norm: if $x = 0$, then $|x|_0 = 0$, and $|x|_0 = 1$, otherwise.

Proposition 1. The optimal solution x^* of the Eq (23) is defined as

$$x^* = \begin{cases} y - \delta, & y > \delta + \sqrt{2\eta}, \\ y + \delta, & y < -\delta - \sqrt{2\eta}, \\ 0, & \text{otherwise.} \end{cases} \quad (24)$$

Proof. First, we denote $E(x) = \frac{1}{2}(x - y)^2 + \delta|x| + \eta|x|_0$. It is obvious that if $x = 0$, then $E(0) = \frac{1}{2}y^2$. Then we need to discuss the case that $x \neq 0$:

1. if $x > 0$, then $E(x) = \frac{1}{2}(x - y)^2 + \delta x + \eta$. Writing its K.K.T condition, we get $x = y - \delta$, and the objective value is $E(y - \delta) = -\frac{1}{2}\delta^2 + \delta y + \eta$.

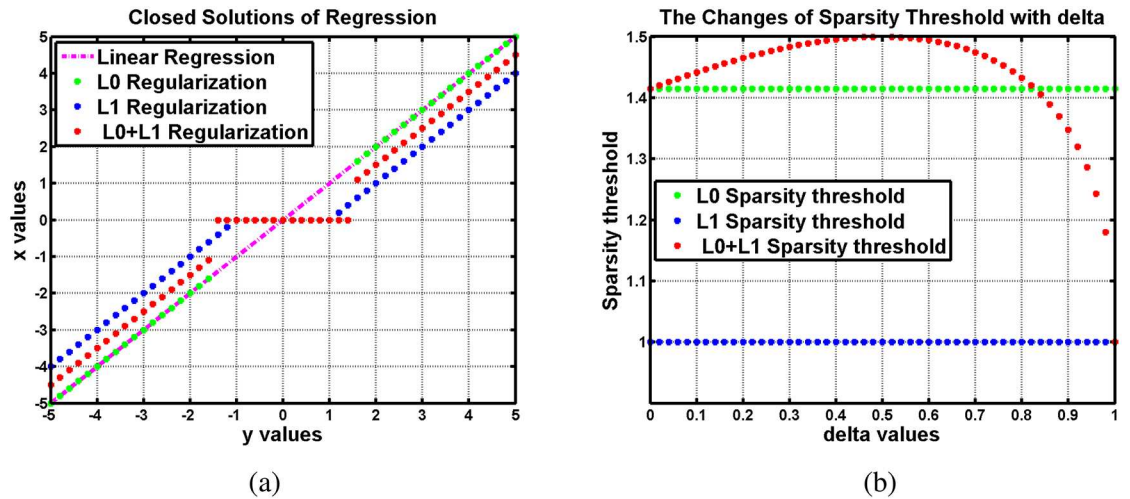


Fig 3. Analysis about combination of L_1 and L_0 regularization. (a) shows the closed solutions of linear regression, L_0 , L_1 , $L_0 + L_1$ regularized regression, respectively. (b) shows the sparsity threshold changes of L_0 , L_1 and $L_0 + L_1$ regularized regression, respectively.

doi:10.1371/journal.pone.0168093.g003

- if $x < 0$, then $E(x) = \frac{1}{2}(x - y)^2 - \delta x + \eta$. It is easy to get $x = y + \delta$, and the objective value is $E(y + \delta) = -\frac{1}{2}\delta^2 - \delta y + \eta$.

Then, we need to compare these three cases, if $E(0) > E(x - \delta)$, we have $(\delta - y)^2 > 2\eta$. Combining with $x = y - \delta > 0$, we have $y > \delta + \sqrt{2\eta}$. Similarly, if $E(0) > E(x + \delta)$, then we have $y < -\delta - \sqrt{2\eta}$. And $x = 0$, otherwise.

If $x \in \mathbb{R}^N$, the Eq (23) changes into

$$\min_x \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \delta \|\mathbf{x}\|_1 + \eta \|\mathbf{x}\|_0, \tag{25}$$

where $\|\mathbf{x}\|_1 = \sum_{i=1}^N |x_i|$ and $\|\mathbf{x}\|_0 = \sum_{i=1}^N |x_i|_0$. It is obvious that Eq (23) can be turned into

$$\min_{x_i} \sum_{i=1}^N \frac{1}{2} (x_i - y_i)^2 + \delta |x_i| + \eta |x_i|_0. \tag{26}$$

So it can be seen as a sequence of optimization of x_i , $i = 1, \dots, n$, and each can be solved by proposition.

In Eq (25), if we set $\delta = 0$ and $\eta = 0$, the model degenerates to the linear regression. If we set $\delta = 0$, Eq (25) reduces to L_0 regularized regression, while becoming L_1 regularized regression when $\eta = 0$. Fig 3(a) shows the closed solutions of these four cases. We set $\delta = \eta = 0.5$ in Eq (25) ($L_0 + L_1$ regularized regression), $\eta = 1$ in L_0 regularized regression, and $\delta = 1$ in L_1 regularized regression. We note that $L_0 + L_1$ regularized regression has the same sparsity as L_0 regularized regression, while causing little over penalization than L_1 regularized regression. In Fig 3(b), sparsity threshold changes of L_0 , L_1 and $L_0 + L_1$ regularized regression are shown, respectively. When $\delta = 1 - \eta$ changes from 0 to 1, the sparsity threshold of $L_0 + L_1$ varies from that of L_0 to the threshold of L_1 . Besides, it is obvious that the threshold of $L_0 + L_1$ is larger than those of L_0 and L_1 in interval $(0, 0.8]$.

Orthogonal Dictionary Learning for Visual Tracking

In this section, we demonstrate dictionary learning in detail through three parts: dictionary initialization, orthogonal dictionary update and dictionary reinitialization.

Dictionary Initialization: There are two schemes to initialize the orthogonal dictionary, one is doing PCA for the set of initial first k frames \mathbf{Y}_k , the other is doing RPCA for \mathbf{Y}_k . When initial frames do not undergo corruption (e.g., occlusion or illumination), we do PCA for \mathbf{Y}_k instead of RPCA. The whole process of PCA is doing skinny SVD for \mathbf{Y}_k and get the basis vectors of column space as the initial dictionary. However, when initial frames have large sparse noise, RPCA is selected to get the intrinsic low-rank features \mathbf{Z}_k , which can be obtained by solving [7]:

$$\min_{\mathbf{Z}_k, \mathbf{E}_k} \|\mathbf{Z}_k\|_* + \lambda \|\mathbf{E}_k\|_1, \text{ s.t. } \mathbf{Y}_k = \mathbf{Z}_k + \mathbf{E}_k. \quad (27)$$

When solving Eq (27), the skinny SVD of \mathbf{Z}_k is readily available: $\mathbf{Z}_k = \mathbf{U}_k \Sigma_k \mathbf{V}_k^T$, and $\mathbf{D} = \mathbf{U}_k$ is the initial orthogonal dictionary. As the analysis in [6], the skinny SVD of \mathbf{Z}_k is readily available when solving Eq (27): Fig 4(a) shows that PCA initialization and RPCA initialization both perform well when the initial first k frames have little noise. The initial frames is generally clean, therefore, we choose PCA initialization as the default.

Orthogonal Dictionary Update: As the appearance of a target may change drastically, it is necessary to update the orthogonal dictionary \mathbf{D} . Here we adopt an incremental PCA algorithm [21] to update the dictionary.

Dictionary Reinitialization: When the tracker is prone to drift, dynamically reinitializing dictionary to obtain the intrinsic subspace features is needed. We adopt the strategy proposed by [7]. The reinitialization is performed at t -th frame if $\sigma = \|\mathbf{e}_t\|_0 / \text{len}(\mathbf{e}_t) > \text{thr}$, where \mathbf{e}_t is the noise item at t -th frame, $\text{len}(\cdot)$ is the length of vector, and $\text{thr} > 0$ is a threshold parameter (generally 0.5). If $\sigma > \text{thr}$, we reinitialize the dictionary in the same way as initialization of dictionary by doing RPCA, but \mathbf{Y}_t in Eq (27) is different. Here, \mathbf{Y}_t consists of optimal candidate observations respectively from the initial n (generally 10) frames and the latest $t - n$ frames (we set $t = 30$). Fig 4(b) compares the tracking performance within and without RPCA reinitialization when the object undergoes variable illumination. After reinitializing dictionary, our tracker retracks the object, so reinitializing dictionary is efficient to improve the reconstruction ability. In Algorithm 2, we summarize the overall tracking process for frame t .

Experimental Results

In this section, we compare the performance of our proposed tracker with several state-of-the-art tracking algorithms, such as TLD [22], IVT [5], ASLA [23], L_1 APG [14], MTT [11], SP [15], SPOT [24], FOT [25], SST [26], SCM [27], MIL [2], and Struck [3], on twenty-two video sequences from the popular benchmark [20] including basketball, bolt, boy, car4, carDark, carScale, crossing, david, david2, david3, deer, faceocc1, faceocc2, fish, football, mountainBike, shaking, skating1, trellis, walking, walking2 and woman. These sequences are publicly available online at http://cvlab.hanyang.ac.kr/tracker_benchmark/datasets.html. Representative videos including Tiger1 and Singer1 have been downloaded from the open video data-sets of the paper [28]. Our tracker is implemented in MATLAB and runs at 4.2 fps on an Intel 2.53 GHz Dual-Core CPU with 8GB memory, running Windows 7 and Matlab (R2013b). We empirically set $\eta = 0.1$, $\lambda = 0.5$, $\gamma = 0.1$, $\tau = 0.05$ and the Lipschitz constant $L = 2$. Before solving Eq (18), all the candidates \mathbf{y} are centralized. Considering the efficiency, the updated orthogonal dictionary \mathbf{D} is taken columns corresponding to the 16 largest eigenvalues of PCA or RPCA, 600 particles are adopted, and the model is incrementally updated every 5 frames. In the

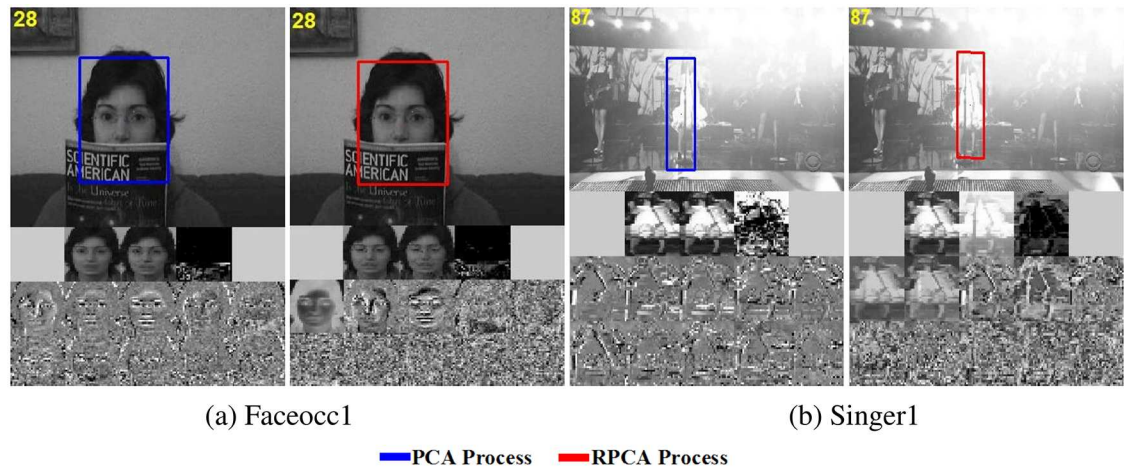


Fig 4. Comparison of PCA and RPCA. The upper portion of the image is the tracking frame. the middle of the image consists of three sub-pictures, the left is the mean image, the middle is the reconstruction result, and the right is the Laplace noise. the bottom of the image is the top ten basis vectors of dictionary. (a) shows the tracking results of PCA and RPCA dictionary initialization. The tracking performance with and without RPCA reinitialization is shown in (b). Reprinted from [20] under a CC BY license, with permission from Yi Wu, original copyright 2013.

doi:10.1371/journal.pone.0168093.g004

following, we present both qualitative and quantitative comparisons of above mentioned methods.

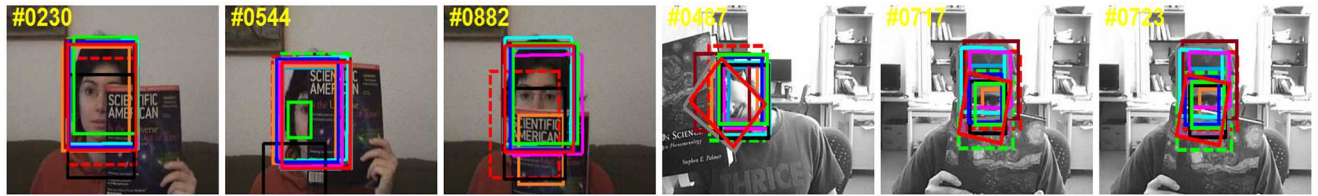
Algorithm 2 Our Robust Visual Tracking Algorithm

Initialization: Initialize orthogonal dictionary \mathbf{D} by performing PCA on \mathbf{Y}_k .
Input: State \mathbf{x}_{t-1} ($t > k$) and orthogonal dictionary \mathbf{D} .
Step 1: Draw new samples \mathbf{x}_i^t from \mathbf{x}_{t-1} and obtain corresponding candidates \mathbf{y}_i^t .
Step 2: Obtain α_i^t and \mathbf{e}_i^t using Eq (19).
Step 3: For each candidate, calculate the observation probability $p(\mathbf{y}_i^t | \mathbf{x}_i^t)$ using Eq (4).
Step 4: Find the tracking result patch \mathbf{y}_i^* with the maximal observation likelihood and its corresponding noise \mathbf{e}_i^* .
Step 5: perform an incremental PCA algorithm to update the orthogonal dictionary \mathbf{D} every five frames. If $\sigma > thr$, reinitializing Dictionary at t -th frame using Eq (27).
Output: State \mathbf{x}_t^* and corresponding image patch; orthogonal dictionary \mathbf{D} .

Qualitative Evaluation

We choose some examples from part of 22 sequences to illustrate the effectiveness of our method. Fig 5 shows the visualization results.

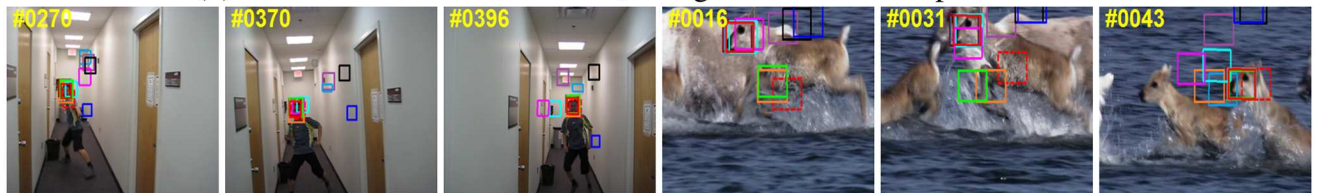
Heavy Occlusion: Fig 5(a) and 5(b) show three challenging sequences with heavy occlusion. In *Faceocc1* and *Faceocc2*, the targets undergo with heavy occlusion and in-plane rotation, it can be seen that our method outperforms the other tracking algorithms. *David3* demonstrates that the proposed method can capture the accurate location of objects in terms of position, and scale when the target undergoes severe occlusion (e.g., *David3* #0085). However, IVT, L_1 APG, MIL, SP, SCM, ASLA, TLD, SPOT, FOT, SST, MTT, and Struck methods drift away from the target object when occlusion occurs. For these four sequences, the IVT method performs poorly since conventional PCA is not robust to occlusions. Although L_1 APG



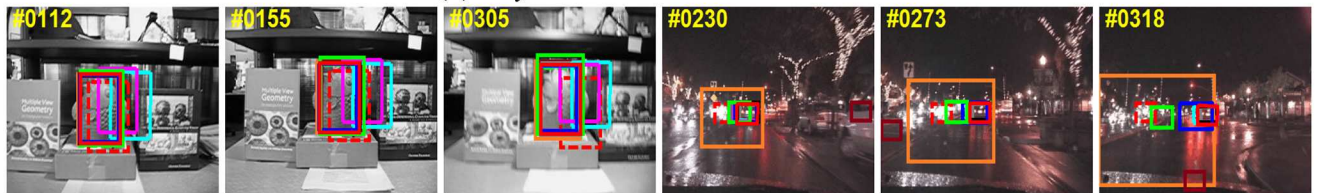
(a) *Faceocc1* and *Faceocc2* with heavy occlusion and in-plane rotation.



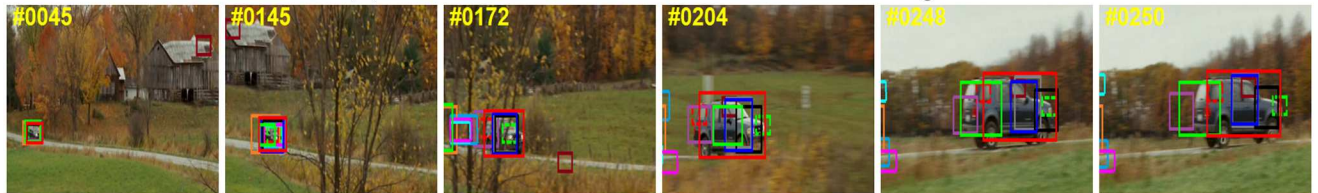
(b) *David3* with severe occlusion, background clutter and pose variation.



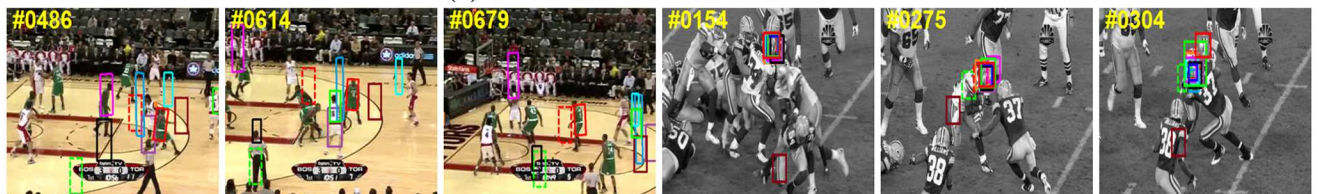
(c) *Boy* and *Deer* with fast motion.



(d) *Fish* and *CarDark* with illumination changes.



(e) *CarScale* with scale variation.



(f) *Basketball* and *Football* with background clutter.

— Ours — TLD — IVT — ALSA — LIAPG — MTT — SP — SPOT — FOT — SST — SCM — MIL — Struck

Fig 5. Sampled tracking results of evaluated algorithms on ten challenging image sequences.

doi:10.1371/journal.pone.0168093.g005

and SP utilize sparsity to model outliers, it is observed that their occlusion detection are not stable when drastic change of appearance happens. In contrast, our method is robust to heavy occlusion. This is because our combination of L_0 and L_1 regularized appearance model can exactly reconstruct the object.

Fast Motion: Fig 5(c) show the sequences *Boy* and *Deer* with fast motion. It is difficult to predict the locations of the tracked objects when they undergo abrupt motion. In *Boy*, the captured images are blurred seriously, but Struck and our method track the target faithfully throughout the images. IVT, MTT, ALSA, SCM and SST methods drift away seriously. We note that most of the other trackers have drift problem due to the abrupt motion and background clutter in sequence *Deer*. In contrast, the SST and our method successfully track the target for whole video.

Illumination Changes and Scale Variation: In Fig 5(d) and 5(e), we test three challenging sequences with illumination changes and scale variation. *Fish* chips contain significant illumination variation. We can see that the L_1 APG, MTT, and MIL methods are less effective in these cases (e.g., *Fish* #0305). In *CarDark*, our method still performs well, but TLD, FOT, and MIL fail. Our method also achieves good performance in *CarScale* with scale variation (e.g., *CarScale* #0204). For subspace-based approaches, they may fail to update the appearance model as the calculation of coefficients in their models may have redundant background features. Our method can successfully adapt to variable drastic changes since the combination of sparse coding and sparse counting is not merely stable but also applicable to obtain the intrinsic features of the subspace.

Background Clutters: Fig 5(f) demonstrates the tracking results in *Basketball* and *Football* with background clutter. *Basketball* is a difficult sequence because it contains cluttered background, illumination change, heavy occlusion and non-rigid pose variation. Unless our tracker, none of the compared algorithms can work well on it (e.g., *Basketball* #0486 and #0614). As shown in *Football*, our tracker performs relatively well (e.g., *Football* #304) as it has excluded background clutters in the sparse errors, but TLD, FOT, and MIL fail.

Quantitative Evaluation

We use two metrics to evaluate the proposed algorithm with other state-of-the-art methods. The first metric is the center location error measured with manually labeled ground truth data. The second one is the overlap rate, i.e., $score = \frac{area(R_T \cap R_G)}{area(R_T \cup R_G)}$, where R_T is the tracking bounding box and R_G is the ground truth bounding box. The larger average scores mean more accurate results.

Table 1 shows the average overlap rates. Table 2 reports the average center location errors (in pixels) where a smaller average error means a more accurate result. Notice that the results are calculated by averaging 5 runs of these algorithms. As can be seen from the table, the most sequences generated by our method have lower average error and higher overlap rate values. We provide the precision and success plots in Fig 6 to evaluate our performance over all the 22 sequences. The evaluation parameters are set as default in [20]. We note that the our algorithm performs well for the videos with occlusion, low resolutionn, in plane rotation, and background clutter based on the precision metric and the success rate metric as shown in Figs 7 and 8 respectively. Both table and figures show that our method achieves favorable performance against other state-of-the-art methods.

To further compare the running time of four subspace-based tracking algorithms (i.e. IVT, L_1 APG, SP and our method), we calculated the average Frames Per Second (FPS) for 32×32 image patch (see the last row of Table 1). For L_1 APG, we reported FPS for its APG acceleration.

Table 1. Average overlap rate and average frame per second (FPS). The best and the second results are shown in **BOLD** fonts and **BOLD** fonts, respectively.

	TLD	IVT	ASLA	L_1 APG	MTT	SP	SPOT	FOT	SST	SCM	MIL	Struck	Ours
Faceocc1	0.58	0.73	0.32	0.76	0.70	0.79	0.74	0.60	0.79	0.79	0.60	0.73	0.80
Faceocc2	0.62	0.73	0.65	0.69	0.75	0.59	0.69	0.64	0.63	0.73	0.67	0.79	0.69
David3	0.10	0.48	0.43	0.38	0.10	0.46	0.77	0.41	0.30	0.41	0.54	0.29	0.73
Boy	0.66	0.26	0.37	0.73	0.50	0.36	0.57	0.64	0.36	0.38	0.49	0.76	0.81
Deer	0.60	0.03	0.03	0.60	0.61	0.72	0.72	0.16	0.62	0.07	0.12	0.74	0.82
Fish	0.81	0.77	0.85	0.34	0.16	0.83	0.83	0.78	0.86	0.75	0.45	0.85	0.87
CarDark	0.45	0.66	0.85	0.88	0.83	0.77	0.00	0.26	0.86	0.84	0.20	0.89	0.85
Jogging-2	0.66	0.14	0.14	0.15	0.13	0.73	0.20	0.12	0.12	0.73	0.14	0.20	0.74
CarScale	0.45	0.63	0.61	0.50	0.49	0.60	0.01	0.35	0.55	0.59	0.41	0.41	0.81
Basketball	0.02	0.11	0.39	0.23	0.19	0.23	0.01	0.17	0.20	0.46	0.22	0.20	0.63
Football	0.49	0.56	0.53	0.55	0.58	0.69	0.01	0.55	0.40	0.49	0.59	0.53	0.59
Average	0.46	0.34	0.36	0.41	0.34	0.50	0.43	0.37	0.39	0.44	0.39	0.41	0.70
FPS	21.74	27.83	7.48	2.47	0.99	2.35	–	376.48	2.12	0.37	28.06	10.01	4.27

doi:10.1371/journal.pone.0168093.t001

Table 2. Average center location error and average frame per second (FPS). The best and the second results are shown in **BOLD** fonts and **BOLD** fonts, respectively.

	TLD	IVT	ASLA	L_1 APG	MTT	SP	SPOT	FOT	SST	SCM	MIL	Struck	Ours
Faceocc1	27.37	18.42	78.06	17.33	21.00	14.14	17.17	29.00	13.00	13.04	29.86	18.78	12.88
Faceocc2	12.28	7.42	19.35	12.76	9.836	10.43	11.78	11.94	12.82	5.96	9.02	13.60	5.50
David3	208.00	51.95	87.76	90.00	341.33	8.74	6.27	33.40	104.50	73.09	29.68	106.50	5.79
Boy	4.49	91.25	106.07	7.03	12.77	58.09	8.93	5.79	66.97	51.02	12.83	3.84	2.57
Deer	30.93	182.69	160.06	24.19	18.91	6.84	13.95	80.30	13.81	103.54	100.73	5.27	4.59
Fish	6.54	5.67	3.85	29.43	45.50	3.99	4.52	6.50	3.14	8.54	24.14	3.40	3.08
CarDark	27.47	8.43	1.54	1.04	1.57	1.35	121.58	34.43	1.19	1.30	43.48	0.95	1.31
CarScale	22.60	11.90	24.64	79.78	87.61	13.36	207.01	106.20	87.05	33.38	33.47	36.43	7.66
Basketball	213.86	107.11	82.64	137.53	106.80	39.79	169.86	118.02	105.93	52.90	91.92	118.6	7.92
Football	14.26	14.34	15.00	15.11	13.67	5.22	202.03	13.36	17.21	16.30	12.09	17.31	7.28
Average	50.48	72.54	78.20	64.58	85.48	39.38	69.46	55.66	88.42	48.26	48.92	49.17	7.97
FPS	21.74	27.83	7.48	2.47	0.99	2.35	–	376.48	2.12	0.37	28.06	10.01	4.27

doi:10.1371/journal.pone.0168093.t002

It can be seen that IVT is quite faster than other trackers as its computation only involves matrix-vector multiplication. Both SP and our method are faster than L_1 APG. It is also observed that our method is much faster than SP. This is due to the different choices of the optimization scheme. SP adopts a naive altering minimization strategy, in contrast, our method is efficiently solved by APG.

Conclusion

In this paper, we propose sparse coding and counting method under Bayesian framework for robust visual tracking. The proposed method combines L_0 regularization and L_1 regularized sparse representation in a unique one, therefore, it has better ability to sparsely represent an object and the reconstruction result are also better. Besides, to solve the proposed model, we develop a fast and efficient APG algorithm. Moreover, the closed solution of the combination

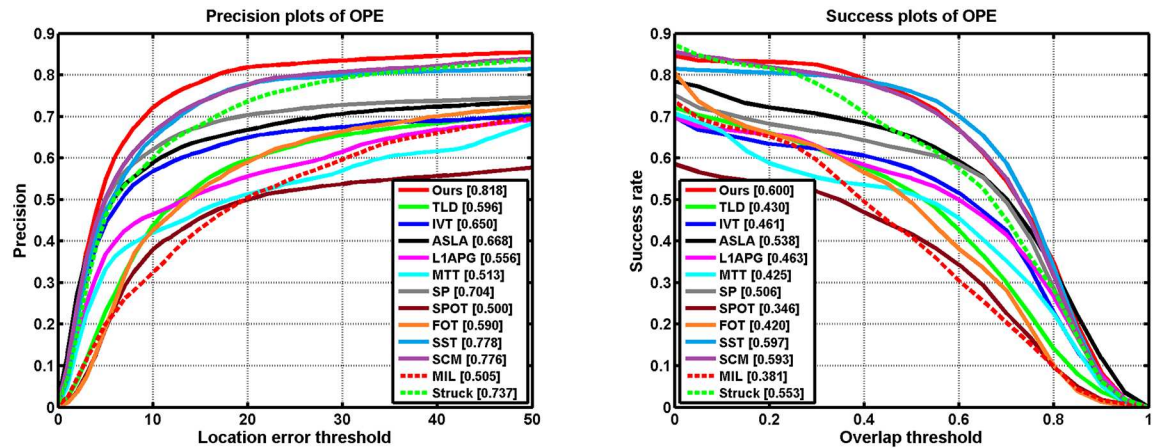


Fig 6. Precision and success plots over all the 22 sequences. The mean precision scores are reported in the legends.

doi:10.1371/journal.pone.0168093.g006

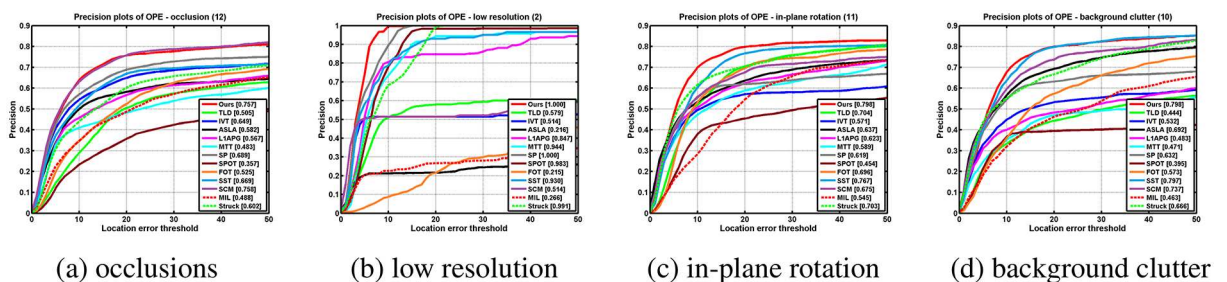


Fig 7. The plots of OPE with attributes based on the precision metric.

doi:10.1371/journal.pone.0168093.g007

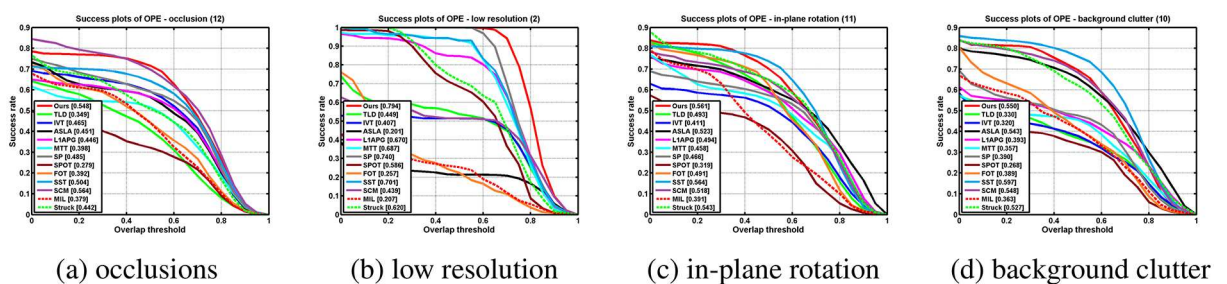


Fig 8. The plots of OPE with attributes using the success rate metric.

doi:10.1371/journal.pone.0168093.g008

of L_0 norm and L_1 norm regularization is provided. Extensive experiments testify to the superiority of our method over state-of-the-art methods, both qualitatively and quantitatively.

Acknowledgments

This work is partially supported by the National Natural Science Foundation of China (Nos. 61300086, 61432003, 61672125), the Fundamental Research Funds for the Central Universities (DUT15QY15), the Hong Kong Scholar Program (No. XJ2015008), and National Science and Technology Major Project (Nos. 2013ZX04005-021, 2014ZX001011).

Author Contributions

Conceptualization: RL JW ZS.

Data curation: JW.

Formal analysis: RL JW.

Funding acquisition: RL ZS.

Investigation: RL JW.

Methodology: RL JW.

Project administration: RL ZS.

Software: JW YC.

Supervision: ZS.

Validation: JW XS.

Visualization: JW YW.

Writing – original draft: RL JW.

Writing – review & editing: RL JW XS.

References

1. Liu R, Cheng J, Lu H. A robust boosting tracker with minimum error bound in a co-training framework. In: ICCV; 2009. p. 1459–1466.
2. Babenko B, Yang MH, Belongie SJ. Visual tracking with online Multiple Instance Learning. In: CVPR; 2009. p. 983–990.
3. Hare S, Saffari A, Torr PHS. Struck: Structured output tracking with kernels. In: ICCV; 2011. p. 263–270.
4. Jepson AD, Fleet DJ, El-Maraghi TF. Robust Online Appearance Models for Visual Tracking. IEEE TPAMI. 2003; 25(10):1296–1311. doi: [10.1109/TPAMI.2003.1233903](https://doi.org/10.1109/TPAMI.2003.1233903)
5. Ross DA, Lim J, Lin RS, Yang MH. Incremental Learning for Robust Visual Tracking. IJCV. 2008; 77(1-3):125–141. doi: [10.1007/s11263-007-0075-7](https://doi.org/10.1007/s11263-007-0075-7)
6. Liu R, Lin Z, Su Z, Gao J. Linear time Principal Component Pursuit and its extensions using ℓ_1 filtering. Neurocomputing. 2014; 142:529–541. doi: [10.1016/j.neucom.2014.03.046](https://doi.org/10.1016/j.neucom.2014.03.046)
7. Zhang C, Liu R, Qiu T, Su Z. Robust visual tracking via incremental low-rank features learning. Neurocomputing. 2014; 131:237–247. doi: [10.1016/j.neucom.2013.10.020](https://doi.org/10.1016/j.neucom.2013.10.020)
8. Liu R, Jin W, Su Z, Zhang C. Latent Subspace Projection Pursuit with Online Optimization for Robust Visual Tracking. IEEE MultiMedia. 2014; 21:47–55. doi: [10.1109/MMUL.2014.49](https://doi.org/10.1109/MMUL.2014.49)
9. Mei X, Ling H. Robust visual tracking using ℓ_1 minimization. In: ICCV; 2009. p. 1436–1443.
10. Liu B, Yang L, Huang J, Meer P, Gong L, Kulikowski C. Robust and fast collaborative tracking with two stage sparse optimization. In: ECCV; 2010. p. 624–637.
11. Zhang T, Ghanem B, Liu S, Ahuja N. Robust Visual Tracking via Structured Multi-Task Sparse Learning. IJCV. 2013; 101(2):367–383. doi: [10.1007/s11263-012-0582-z](https://doi.org/10.1007/s11263-012-0582-z)
12. Jin W, Liu R, Su Z, Zhang C, Bai S. Robust visual tracking using latent subspace projection pursuit. In: ICME; 2014. p. 1–6.
13. Pan J, Lim J, Su Z, Yang MH. ℓ_0 -Regularized Object Representation for Visual Tracking. BMVC. 2013;.
14. Bao C, Wu Y, Ling H, Ji H. Real time robust ℓ_1 tracker using accelerated proximal gradient approach. In: CVPR; 2012. p. 1830–1837.
15. Wang D, Lu H, Yang MH. Online Object Tracking With Sparse Prototypes. IEEE TIP. 2013; 22(1):314–325.

16. Lu X, Wang Y, Yuan Y. Sparse Coding From a Bayesian Perspective. *IEEE Transactions on Neural Networks and Learning Systems*. 2013; 24(6):929–939. doi: [10.1109/TNNLS.2013.2245914](https://doi.org/10.1109/TNNLS.2013.2245914) PMID: [24808474](https://pubmed.ncbi.nlm.nih.gov/24808474/)
17. Wang D, Lu H, Yang MH. Least Soft-threshold Squares Tracking. In: *CVPR*; 2013. p. 2371–2378.
18. Lin Z, Ganesh A, Wright J, Wu L, Chen M, Ma Y. Fast convex optimization algorithms for exact recovery of a corrupted low-rank matrix. *UIUC*; 2009.
19. Tseng P. On accelerated proximal gradient methods for convex-concave optimization; 2008. submitted to *SIAM J. Optimiz.*
20. Wu Y, Lim J, Yang MH. Online Object Tracking: A Benchmark. In: *CVPR*; 2013. p. 2411–2418.
21. Levey A, Lindenbaum M. Sequential Karhunen-Loeve basis extraction and its application to images. *IEEE Trans on IP*. 2000; 9(8):1371–1374.
22. Kalal Z, Mikolajczyk K, Matas J. Tracking-Learning-Detection. *IEEE TPAMI*. 2012; 34(7):1409–1422. doi: [10.1109/TPAMI.2011.239](https://doi.org/10.1109/TPAMI.2011.239)
23. Jia X, Lu H, Yang MH. Visual tracking via adaptive structural local sparse appearance model. In: *CVPR*; 2012. p. 1822–1829.
24. Zhang L, Maaten L. Structure preserving object tracking. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2013. p. 1838–1845.
25. Vojř T, Matas J. The enhanced flock of trackers. In: *Registration and Recognition in Images and Videos*. Springer; 2014. p. 113–136.
26. Zhang T, Liu S, Xu C, Yan S, Ghanem B, Ahuja N, et al. Structural sparse tracking. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2015. p. 150–158.
27. Zhong W, Lu H, Yang MH. Robust object tracking via sparsity-based collaborative model. In: *CVPR*; 2012. p. 1838–1845.
28. Bai Q, Wu Z, Sclaroff S, Betke M, Monnier C. Randomized ensemble tracking. In: *Proceedings of the IEEE International Conference on Computer Vision*; 2013. p. 2040–2047.