

pedestrian is heavily affected by camera perspective and crowd density, also it is observed that different features can be more important given various crowdedness levels. In addition, their evaluations show that the actual performance of a regression model can be quite different from what one may anticipate, subject to the nature of data, especially when it is applied to unseen crowd density.

Unlike regression techniques, our proposed method based on sparse representation, does not need to select either the optimal feature set or the regression model. The main idea behind sparse representation is, if a collection of representative samples are found, we should expect that a typical sample has a very sparse representation with respect to such a learned basis. In other words, given sufficient diversity in the training images, the new test image can be well represented as a sparse linear combination of the training set. This sparse representation would naturally encode the semantic information of the image [9]. In order to reduce the time complexity of finding the sparse representation, random projection is utilized as our choice of dimensionality reduction method.

It is commonly believed that the Sparse Representation-based Classification (SRC) requires a rich set of training images of every class that can span the variation under testing conditions. To fulfill this requirement, we use a semi-supervised learning framework to avoid exhaustive manual image annotation. Extensive experimental results suggest that our proposed method is fast, accurate and scalable to large-scale datasets.

The remainder of the paper is organized as follows: the theory of sparse representation is summarized in Section 2. Section 3 shows how to apply general classification framework to people counting task. In Section 4 we discuss how we exploit semi-supervised regression to deal with few labelled training samples effectively. Experimental setup is explained in Section 5 and results and discussion are presented in Section 6, followed by conclusion remarks in Section 7.

2. Sparse representation

Sparse representation (SR) has proven to be an extremely powerful tool for acquiring, representing, and compressing high-dimensional signals. This success is mainly due to the fact that important classes of signals such as audio and images have naturally sparse representations with respect to fixed bases e.g. Fourier and Wavelet. Moreover, in recent years, efficient and fast algorithms have been proposed for computing such representations [9]. The problem solved by sparse representation is to search for the most compact representation of a signal (image) in terms of a linear combination of relatively few base elements in a basis or over-complete dictionary. If the optimal representation is sufficiently sparse, it can be efficiently computed by greedy methods or convex optimization. Typically, the sparse representation technique is cast into an l_1 -minimization problem, which is equivalent to the l_0 -minimization under some conditions. This l_0 - l_1 equivalence has provided computational convenience as evidenced by Compressed Sensing (CS) [10].

In the recent years, variations and extensions of l_1 -minimization have been applied to many computer vision tasks, including face recognition [11], background modelling [12] and image classification [13]. In almost all of these applications, using sparsity as a prior leads to the state-of-the-art results [9]. The ability of sparse representation to uncover semantic information derives in part from a simple but important property of the data: although the images (or their features) are naturally very high dimensional, in many applications images belonging to the same class exhibit degenerate structure. That is, they lie on or near low-dimensional subspaces or submanifolds [9]. So, if a collection of representative samples are found, we

should expect that a typical sample has a very sparse representation with respect to such a (possibly learned) basis. Such a sparse representation, if computed correctly, could naturally encode the semantic information of the image [9]. SRC seeks a sparse representation of the query image in terms of the over-complete dictionary and then performs the recognition by checking which class yields the least representation error. SRC can be considered as a generalization of Nearest Neighbor (NN) and Nearest Feature Subspace (NFS). Generally speaking, Nearest Feature based Classifiers (NFCs) aim to find a representation of the query image, and classify it to the best representor. According to the mechanism of representing the query image, NFCs include Nearest Neighbor, Nearest Feature Line (NFL), Nearest Feature Plane (NFP) and Nearest Feature Subspace. More specifically, NN is the simplest one with no parameters, which classifies the query image to its nearest neighbor. NN, NFL and NFP all use a subset of the training samples with the same label to represent the query image, while NFS represents the query image by all the training samples of the same class. In general, the larger samples lead to better stability of a method. The most generalized classifier is SRC, which considers all possible supports (within each class or across multiple classes) and adaptively chooses the minimal number of training samples needed to represent each test sample. In the next section, we show how this sparse representation can be used in people counting application.

3. People counting based on sparse representation and random projection

3.1. People counting as sparse representation

Suppose that we have a set of labelled (annotated) training images from a pedestrian dataset where the number of people present in each image is given. We assume these labelled training images $\{x_i, l_i\}$ are from C different classes. Here, class (label) l_i is equal to the count, i.e. number of people in the image $x_i \in R^m$, where x_i is the vector representation of the image, which could be its raw pixels or features computed from the raw pixels. Given sufficient training samples from the i th class, any new test sample $x_{test} \in R^m$ from the same class, will approximately lie in the linear span of the training samples associated with class i .

$$x_{test} \approx \sum_{j|l_j=i} x_j \alpha_j = X_i \alpha_i \quad (1)$$

where $X_i \in R^{m \times n_i}$ concatenates all of the images of class i . Since the class label of the test image is initially unknown, we would form a linear representation similar to Eq. (1), now in terms of all training samples. We define a new matrix (dictionary) $\Psi \in R^{m \times n}$ for the entire training set as the concatenation of all $n = \sum_i n_i$ training samples of all C classes:

$$x_{test} = [X_1, X_2, \dots, X_C] \alpha = \Psi \alpha \in R^m \quad (2)$$

where

$$\alpha = [\dots, 0^T, \alpha_i^T, 0^T, \dots]^T \in R^n \quad (3)$$

α is a coefficient vector whose entries are zero except those associated with the i th class. We notice that α is a highly sparse vector and on average, only a fraction of $1/C$ coefficients are nonzero and the dominant nonzero coefficients in the sparse representation α reveal the true class of test image. Indeed, in the test phase, we wish to represent a new unlabelled image in a Ψ -dependent space in which the image has a sparse representation. In general, this vector is the sparsest solution to the system of equations $x_{test} = \Psi \alpha$ which is found by solving the following optimization problem:

$$\alpha^* = \operatorname{argmin} \|\alpha\|_0 \quad \text{s.t. } \Psi \alpha = x_{test} \quad (4)$$

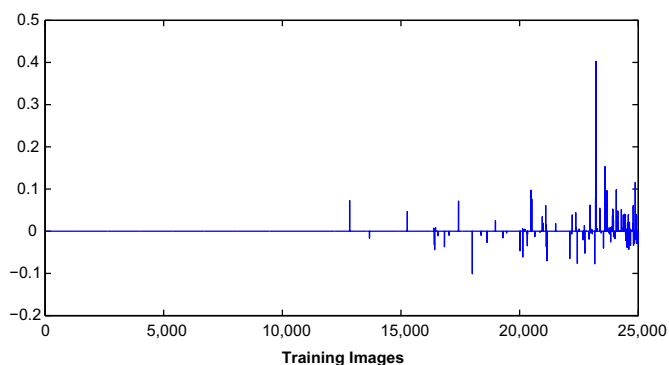


Fig. 1. Representation of a test image as sparse linear combination of the training set.

While the search for a sparse solution to a linear system is a difficult problem in general, foundational results in the theory of sparse representation show that the sparsest solution can be exactly recovered by solving a tractable optimization problem [14]. This is achieved by seeking α as the unique solution to the optimization problem:

$$\alpha^* = \operatorname{argmin} \|\alpha\|_1 \quad \text{s.t.} \quad \|\Psi\alpha - x_{\text{test}}\|_2 \leq \epsilon \quad (5)$$

where ϵ is the noise level in the observation. Eq. (5) can be rewritten as follows using a Lagrangian multiplier:

$$\alpha^* = \operatorname{argmin} \|\Psi\alpha - x_{\text{test}}\|_2^2 + \lambda \|\alpha\|_1 \quad (6)$$

Fig. 1 shows an example of vector α recovered by solving Eq. (6).

3.2. The role of feature extraction

Conventional visual representations for creating the dictionary Ψ , include local and global image descriptors. Local descriptors like SIFT [15] are found on the premise that images can be characterized by attributes computed on regions of the image; however, global (whole image) descriptors do not require any keypoint detection. They have the ability to characterize an entire image with a single vector. Also, they are fast to build and efficient to store. Moreover, the previous studies [16] show that global descriptors are highly effective in describing the crowd density when enough training samples are involved. Gist [17] represents an image in terms of its responses to a bank of Gabor filters with different frequencies and orientations. An image is divided into tiles and the final feature descriptor would be the mean response of tiles to steerable filters. HOG [3] counts the occurrences of gradient orientation in localized portions of an image. It is a window based descriptor densely sampled over all image points; the window is divided into a square grid and the distribution of edge orientations within each cell is computed.

Recently, in contrast to the hand-crafted features, learnt image features with deep network structures have shown great potential in various vision recognition tasks. Among these architectures, one of the greatest breakthroughs in image classification is the deep Convolutional Neural Network (CNN) [18], which has achieved the state-of-the-art performance in the large-scale object recognition task. Thanks to an optimized GPU implementation and new regularization techniques, Krizhevsky et al. [18] successfully applied CNNs to classification on the ImageNet dataset [19]. The feature representation learned by this network shows excellent performance not only on the ImageNet classification task it was trained for, but also on a variety of other recognition tasks [20,21]. In this paper, we would evaluate whether features extracted from a pre-trained CNN can be reused to the people counting task. To the best of our knowledge, this is the first attempt to employ deep

learning to solve the pedestrian counting problem. For training CNNs and subsequent feature extraction, we use the DeCAF code [20] and instead of training a network from scratch, we use a pre-trained model on ImageNet. We refer the reader to [18] for more details on the architecture and training algorithm, which we followed exactly. As the feature descriptor, we use $DeCAF_7$; the features taken from the final hidden layer, i.e. just before propagating through the final fully connected layer to produce the class predictions.

3.3. Dimensionality reduction

One major obstacle in real-world large-scale people counting is the large scale of the training data and the high dimensional feature vectors. A classical technique for addressing the problem of high dimensionality is to project the data into a much lower-dimensional feature space R^d ($d \ll m$), such that the projected data still retain the useful properties of the original images. From Matten et al.'s study [22] one may infer that, all non-linear dimensionality reduction techniques require the optimization of some parameters and some of them suffer from memory complexity issue. Also, when $m < n$, where m is the feature dimension and n is the number of training data, non-linear techniques have computational disadvantages compared to the classical linear algorithm, such as principal component analysis (PCA) [23]. Thus, we selected PCA as the baseline dimensionality reduction technique for our application; however, PCA is still expensive for high-dimensional and large-scale image data; the computational complexity of PCA is $O(m^2n + m^3)$ due to the matrix computation and SVD eigenvalue decomposition [24]. A statistically optimal way of dimensionality reduction is to project the data onto a random lower-dimensional subspace that captures as much of the variation of data as possible. This technique is called Random Projection (RP), which is computationally efficient and simple. In RP, the original m -dimensional data is projected to a d -dimensional subspace through the origin, using a random matrix $\Phi_{d \times m}$ whose columns have unit lengths. The key idea of random projection arises from the Johnson–Lindenstrauss lemma [25]: if points in a vector space are projected onto a randomly selected subspace of suitably high dimension, then the distances between the points are approximately preserved. In contrast, PCA does not guarantee (approximate) distance preservation between data point pairs. Random projection is computationally very simple; its complexity is bounded by $O(mdn)$ and this could be even less when the data matrix is sparse [24]. Applying such a linear projection on training and test images, gives us a new observation:

$$\begin{aligned} \tilde{x}_{\text{test}} &= \Phi x_{\text{test}} \in R^d \\ \tilde{\Psi} &= \Phi \Psi \in R^{d \times n} \end{aligned} \quad (7)$$

Then the sparsest solution to α , would be obtained via a lower complexity convex optimization:

$$\alpha^* = \operatorname{argmin} \|\tilde{\Psi}\alpha - \tilde{x}_{\text{test}}\|_2^2 + \lambda \|\alpha\|_1 \quad (8)$$

Since $d \ll m$, the complexity is significantly reduced. RP offers clear benefits over PCA; the choice of projection does not depend on the data, it is much faster and as the projected dimension is decreased and drops below a threshold, RP offers a gradual degradation in the performance; however, the degradation suffered by PCA is not necessarily so gradual. At smaller dimensions, PCA distorts the data and this is mainly because the performance of PCA is dependent on the sum of omitted eigenvalues [24]. The choice of random matrix Φ is one of the key points of interest. The elements of Φ are often Gaussian distributed, i.e. the entries are independently sampled from a zero mean normal distribution $N(0, 1)$ and each row is normalized to unit length (makes rows orthogonal),

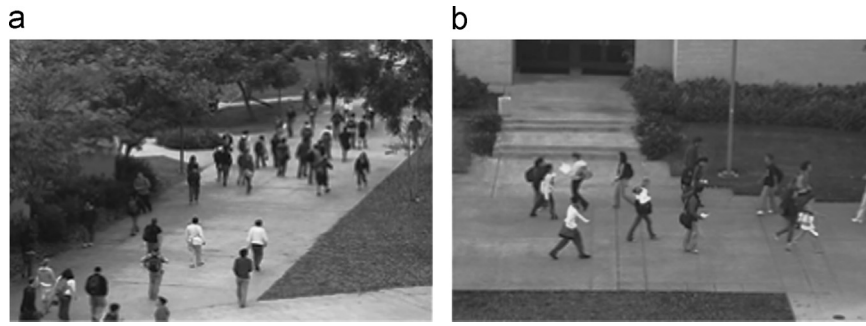


Fig. 2. UCSD pedestrian dataset. (a) Peds1 and (b) Peds2.

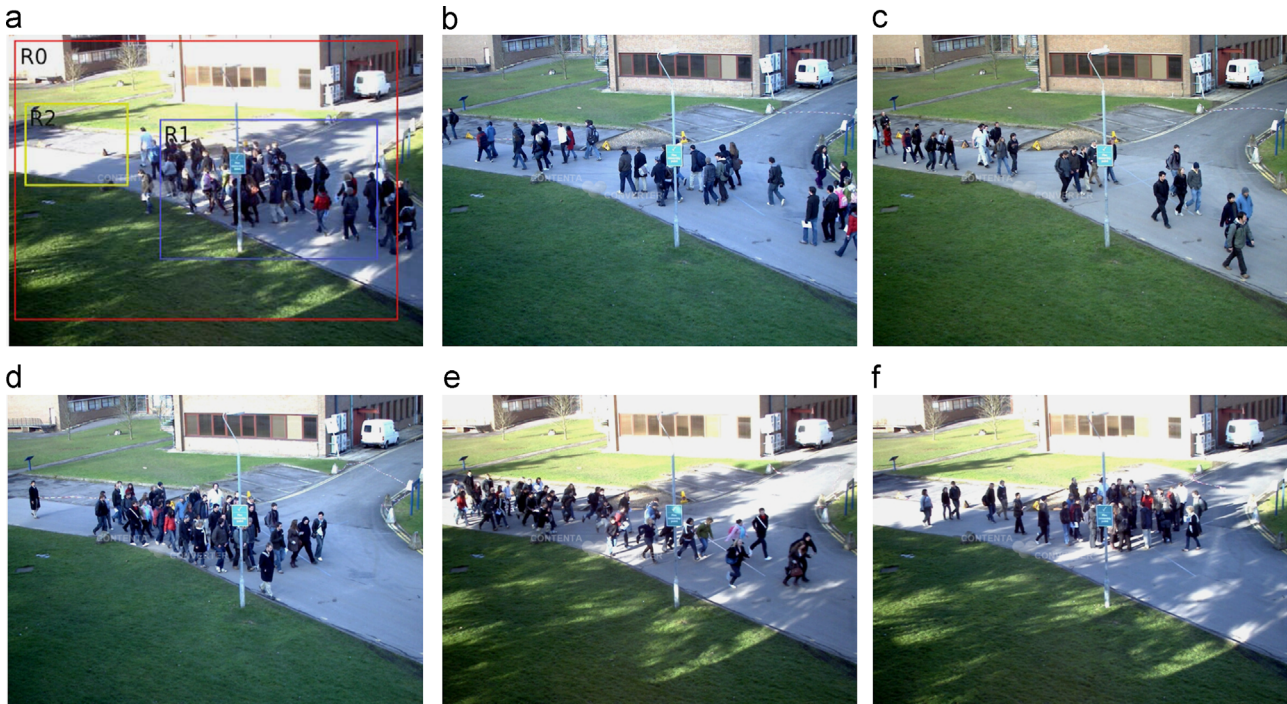


Fig. 3. PETS 2009 pedestrian dataset. (a) Regions of interest, (b) S1L1-1 Sequence, (c) S1L1-2 Sequence, (d) S1L2-1 Sequence, (e) S1L3-1 Sequence, (f) S1L3-2 Sequence.

competition and we would like to evaluate our proposed method on this dataset as well.

6. Results and discussion

6.1. Crowd counting results on UCSD dataset

We conduct extensive experiments to validate our proposed method in different scenarios. For all the experiments, training samples are randomly selected from the pedestrian dataset and the remaining ones are used for testing. Randomly choosing the training set ensures that our results and conclusions will not depend on any special choice of training data. In all the following graphs, the errors are shown in a logarithmic scale. Firstly, we examine the role of features for creating dictionary within our framework. We compare the performance on using different types of feature descriptors including Gist, HOG, CNN features obtained by DeCAF and raw sub-sampled binary image. Fig. 4 shows the MAE across various feature spaces in conjunction with five different training set sizes including 10k, 15k, 20k, 25k and 30k images on Peds1 and Peds2 datasets. According to this figure, Gist has the best overall performance over all training sizes in both datasets and it is closely followed by HOG and DeCAF. Interestingly

DeCAF works well; although not better than Gist and it confirms the previous studies which claim that the activations invoked by an image within the top layers of a large CNN provide a high-level descriptor of the visual content of the image. These results suggest that instead of learning a counting model from scratch in every new scene, the labelled data from other scenes or even from a complete different dataset could be exploited to compensate for the lack of labelled data in the new scene. To put it differently, our proposed people counting method is independent of both training dataset and feature representation which means the choice of an "optimal" feature transformation is no longer critical; even CNN features trained on a different dataset, should perform as well as any other carefully hand-engineered features. Furthermore, we observe that all the descriptors gain from involving more training samples. The reason is that, when we have more diversity in the training images, the new test image can be well represented as a sparse linear combination of the training set; which in turn leads to lower estimation error. Fig. 5 displays the crowd count estimates using 30k training images and Gist feature representation on Peds1 dataset; these estimates track the ground-truth well in most of the test set. As discussed earlier, in order to reduce the complexity of solving l_1 -minimization problem, random projection is employed as a fast and simple method to project the data into a much lower-dimensional space. Here, we demonstrate the

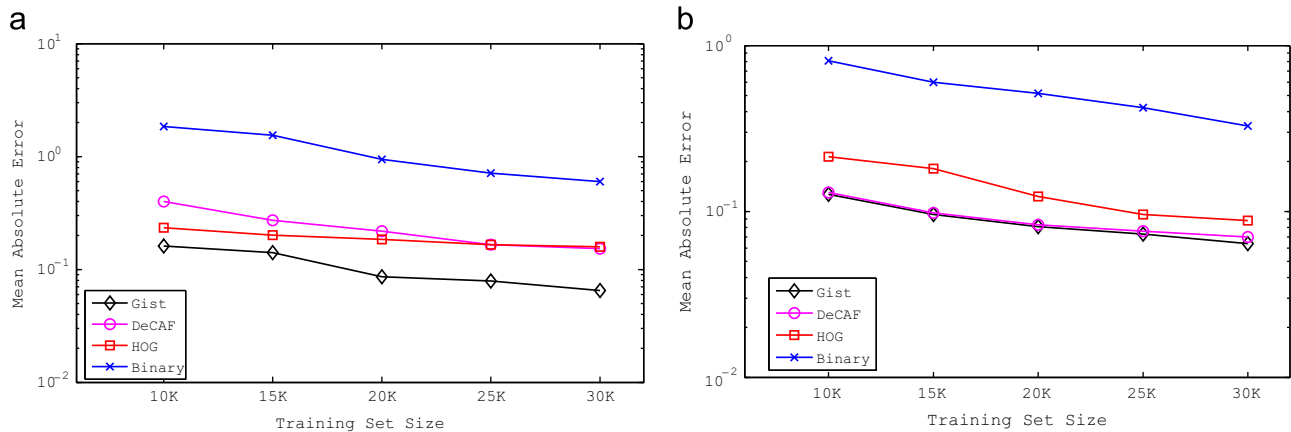


Fig. 4. MAE rate of SR-based method over various features and training sizes. (a) Peds1 and (b) Peds2.

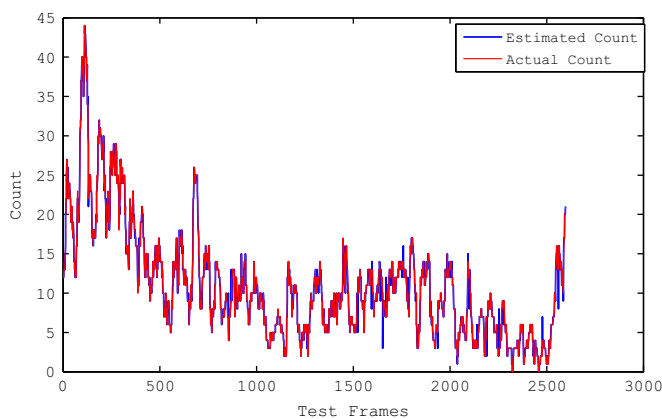


Fig. 5. Count estimation result on Peds1 dataset.

robustness of RP to preserve the similarities of the data points, while reducing the original dimension, across different feature representations in Fig. 6. We fix the training set size as 30k and use the Gaussian random matrix for the projections. We can see that Gist still performs the best in all the reduced dimensions in both datasets; also the projected vectors retain their original performance for up to even 90% reduction in the original dimension. Surprisingly, RP preserves the similarity of feature vectors well even when the data is projected to as low as 10% of original dimensions. Using such an efficient dimensionality reduction technique, the l_1 -minimization can be solved much faster without losing too much accuracy. Our next experiment concerns with the evaluation of RP vs. PCA as a baseline dimensionality reduction technique, in terms of accuracy. Since Gist is the most reliable and remarkable feature to represent the crowd density in our SR-based counting method, we decided to continue the upcoming experiments with it. Here, both RP and PCA are used to construct the lower-dimensional space on different dimension fractions of Gist using 30k training images; the accuracy (MAE) of the estimations has been summarized in Fig. 7. It is observed that the error rate is less affected using RP when dimension is reduced, in comparison with PCA. This is mainly due to the fact that PCA is dependent on the omitted eigenvalues, while RP continues to give accurate results in lower dimensions. It seems that Gist feature descriptors lie on a manifold that PCA is unable to handle it gracefully. We should also note that RP has much less time complexity than PCA which in total makes it the best candidate for dimensionality reduction in our people counting task. We then study the performance of two fast l_1 -solvers; Homotopy and DALM. We

evaluate the efficiency of these solvers on large training sets including 10k, 15k, 20k, 25k and 30k images. In this experiment, Gist feature descriptors of selected images are used to construct the dictionary. The MAE rate of estimations can be observed in Fig. 8. Homotopy produces smaller errors and outperforms DALM in both datasets and converges in fewer steps. The objective of next experiment is to compare different methods of count estimation described in Section 3.3. It can be observed in Fig. 9 that using minimum residual error in Eq. (10) slightly outperforms using maximization l_1 -ratio in Eq. (9); however, these two techniques are both much better than averaging largest entries (peaks). In this experiment, Gist feature descriptors of 30k images are used to construct the dictionary.

6.2. Undersampled SRC

In all of the above experiments, we assumed that we have sufficient labelled training samples, which is achieved via manual annotation of the UCSD dataset. So as to significantly reduce the amount of manual annotation and make our method much more applicable in practice, we use SSEN which enables us to annotate very few images. The goal of the following experiments is to evaluate the effectiveness of exploiting semi-supervised learning framework and especially SSEN for automatic labelling.

We start with a small dataset and choose l labelled samples from the ground-truth (g) by performing k -means clustering on the samples; then u sequential unlabelled frames are selected and the rest of samples ($g-l$) is used as the test set. We evaluate the inductive inference performance [37] of SSEN which is the error rate on the "unlabelled data in the test partition". We examine the effect of labelled and unlabelled data by measuring the MSE performance across labelled set $\{10, 50, 100, 200, 400, 600\}$ given unlabelled set $\{0, 200, 400, 800, 1000\}$ on Peds1 and Peds2 datasets. Images are represented by their Gist feature descriptor and all the λ parameters of SSEN are chosen by a 5-fold cross validation. Fig. 10 shows clearly that when the number of labelled samples is small, increasing the number of unlabelled samples remarkably improves the annotation performance; which means that manual labelling work can be greatly reduced without losing the performance. We follow an iterative procedure aiming to obtain an enlarged labelled dataset, in which we accept that the predictions tend to be correct. During each iteration, the unlabelled samples in the test partition are given predicted labels and then the most confident unlabelled samples, together with their predicted labels, are used to enlarge the training set. Ideally, these selected unlabelled instances can finally help to learn a better classifier. The learner is re-trained on the updated training set and the whole

model. We select Gaussian Process Regression (GPR) which was reported as one of the most accurate ones to handle the crowd density [6]. Although regression models could achieve promising results on small datasets, they suffer from serious weaknesses when they are generalized on large-scale datasets. Loy et al. [8] reveal that the actual performance of regression models can be quite different from what we anticipate, especially when they are applied to an unseen density, they tend to overestimate or underestimate subject to nature of data. Another key weakness is their poor tractability to large training datasets. Usually when more and more data is involved they are unable to adequately capture the non-linear trend in the feature space. They also need online training which takes a lot of time, in particular the time complexity of GPR is bounded by $O(n^3)$ which is a limitation in large datasets. Besides, selecting the optimal feature combination is an important step which is very dependant on the crowd structure and density.

We observed that gradient-based global descriptors (especially Gist) perform well in predicting count in large-scale datasets; this motivated us to examine the effectiveness of local gradient-based descriptors in the sparse representation framework as well. Initially, some local descriptors like SIFT [15] are extracted from an image and then each local descriptor is encoded into a sparse code using a dictionary learning technique. After obtaining the sparse codes, we pool them into a single vector using max pooling. Ge et al. [40] used this idea to derive a compact yet discriminative image representation from multiple types of features for large-scale image retrieval and achieved significant results on benchmark retrieval datasets; however, the aggregated local features cannot capture the crowd patterns properly.

We next compare SRC with three classical classifiers, namely, NN, linear SVM and non-linear SVM with RBF kernel within the context of people counting on large-scale datasets. We can observe that SRC outperforms NN in both datasets and this is essentially because NN independently evaluates the distance between the test sample and one training sample; in contrast, SRC uses a linear combination of all the training samples to represent the test sample and classify it into the class with the minimum deviation and this justifies SRC achievement. Essentially, the training samples are not uncorrelated and the distance between the test and training samples should not be independently calculated, rather; the relationship between different training instances should be taken into account. NN performs a linear search and adopts only one sample in $O(n)$; while SRC finds the sparsest representation in $O(n^2)$. The multi-class SVM classifier is implemented using the LIBSVM [41], which uses a one-against-one decomposition strategy. We select two kernels, linear and Gaussian radial basis function (RBF). Linear SVM is not appropriate for separating Gist features of different crowd densities; nevertheless, better performance is achieved using non-linear kernels like RBF; however our method is still better than best of SVM. When the classes in training data are non-separable by the SVM, SRC might have an advantage over the SVM classifier, depending on similarity between the test vector and the training exemplars from the same class. All the parameters and also the best model are estimated through 5-fold cross-validation over a large training set and this step is really time-consuming. Also, the time complexity of non-linear SVM is generally between $O(n^2)$ and $O(n^3)$ depending on the number of iterations [41]; however this depends a lot on the solving techniques. In brief, SRC outperforms traditional classifiers, meanwhile, it is faster and does not need any model/parameter selection. Finally, the proposed SR-based method is compared with state-of-the-art image retrieval methods including Bag of Features (BOF) [42] and aggregation-based representation, i.e. Fisher Vector [43] and VLAD [44]. In the former, a set of local image patches is sampled using a keypoint detector and a vector of visual descriptors is evaluated on each patch independently. The resulting

distribution of descriptors is then quantified to convert it to a histogram of votes for codebook centers and the resulting global descriptor vector is used as a characterization of the image. In Fisher and VLAD methods, a bag of local features in an image is converted into a global, fixed size vector representation through vector quantization. We observe that the evaluated retrieval methods, easily fail to estimate crowd density in all large-scale datasets and global image descriptors outperform the local ones in general. We also compare the average processing time spent on a test image to estimate its count using different methods on Peds1 dataset while using 30k training images in Table 5; notice that the training time has not been considered here.

On the whole, results suggest that SR-based people counting method is superior to all the other evaluated techniques and the errors of our method are significantly less than the others' including state-of-the-art counting by regression. In brief, the best performance can be achieved by using the Gist descriptor for creating the dictionary, Homotopy l_1 -minimizer and minimum-residual error as the count estimation method; meanwhile random projection is exploited as the dimensionality reduction technique.

6.4. SRC on small datasets

The proposed SR-based people counting method performs well in predicting counts in large datasets; however, we are also interested in evaluating it on smaller datasets, which are popular in counting by regression papers. We follow the experimental protocol in [6], where the training set consists of 1200 and 1000 frames for Peds1 and Peds2 datasets respectively, and the remaining 2800 and 3000 frames are held out as the test set. Table 6 presents the counting accuracy of our method versus [6] as the state-of-the-art regression-based method. Although the error rates are not as small as the larger datasets (e.g. 30k images); our proposed method still outperforms the regression-base method, even with CNN feature representation. Obviously, as the data keeps getting bigger, SRC is coming to play a key role in providing better estimations and the superiority of SRC would be more remarkable. Notice that these conditions are different from the "Undersampled SRC"; here the dictionary is rich enough to span variations of different classes under testing conditions because the

Table 5

Average processing time of a test image in seconds across different methods.

Method	Time	Details
SR-RP	0.55	Gist extraction: 0.1, Test time (l_1 -min): 0.45
Regression	2.10	SET extraction: 1, GPR test time: 1.10
SR-Pooling	0.80	Encoding: 0.6, SIFT extraction: 0.2, Query time: 0.005
NN	0.60	Gist extraction: 0.1, Linear search: 0.5
Lin-SVM	0.30	Gist extraction: 0.1, SMV test time: 0.2
RBF-SVM	0.30	Gist extraction: 0.1, SVM test time: 0.2
BOF	1.84	Encoding: 1.55, SIFT extraction: 0.2, Query time: 0.09
VLAD	0.37	Encoding: 0.15, SIFT extraction: 0.2, Query time: 0.02
Fisher	0.86	Encoding: 0.16, SIFT extraction: 0.2, Query time: 0.5

Table 6

Counting accuracy of SR-based and regression-based methods on UCSD dataset.

Method	Feature	Peds1		Peds2	
		MAE	MSE	MAE	MSE
SR-RP	Gist	0.55	4.13	0.74	2.03
SR-RP	DeCAF	0.87	6.69	0.97	3.62
GPR [6]	SET	3.65	7.41	1.58	2.16

Table 8
Crowd counting estimation error of different counting methods on PETS 2009.

Sequence	Region	Method							
		SR-RP (Our)	Chan [45]	Alahi [46]	Albiol [47]	Choduri [48]	Patzold [49]	Conte [50]	Subburaman [51]
S1L1-1	R0	1.30	2.46	-	1.42	1.29	2.75	1.38	5.95
	R1	1.23	2.28	-	-	2.23	2.58	2.14	1.90
	R2	0.71	0.99	-	-	0.70	1.38	7.60	2.50
S1L1-2	R0	1.01	1.41	4.20	1.77	3.26	2.35	1.14	2.08
	R1	0.70	0.69	2.30	-	3.18	1.58	0.80	1.86
	R2	0.99	1.23	1.87	-	1.04	1.58	0.87	0.86
S1L2-1	R1	1.95	5.89	6.50	1.94	3.70	6.37	2.18	2.40
	R2	1.90	4.48	4.00	-	4.17	6.08	3.25	1.40
S1L3-1	R1	0.86	0.98	0.90	1.36	0.67	4.70	2.95	7.00
S1L3-2	R1	0.01	-	8.83	-	16.50	24.38	-	16.00

proposed SR-RP method performs robustly throughout all the sequences and especially shows a promising performance on very dense crowd sequences such as S1L2-1 and S1L3-1.

6.6. Discussion and future work

Discriminative feature representation plays an important role for achieving the state-of-the-art performance in image classification. In particular, learning sparse representation has recently achieved impressive results and found many applications in machine learning and computer vision fields. Sparse coding represents an input data as a linear combination of a few items from a dictionary. The performance of sparse representation is closely related to the dictionary, which should faithfully and discriminatively represent the test image. The SRC algorithm naively uses all the training samples as the dictionary and achieves promising performance in many applications including people counting. Improvements in computational and memory efficiency of the SRC method make it a more practical solution to be implemented for real-time applications. To this aim, in this paper, we utilized fast l_1 -solvers and dimensionality reduction technique. For some applications, however, rather than using the entire set of training data as dictionary, it is computationally advantageous to learn a compact dictionary from training data. Different algorithms have been developed for learning such dictionaries in an unsupervised and supervised fashion; however recent research indicates that the dictionaries obtained in an unsupervised way, may not necessarily be the best for classification [54].

Many efforts have been dedicated to embedding the discriminative information into the representations via supervised learning [54]. Supervised learning approaches can be divided into different categories [55]: while the algorithms in the first group, learn multiple or class-specific dictionaries to promote discrimination, the second set of approaches incorporate discriminative terms into the objective function of dictionary learning and the last type of approaches learn a compact dictionary by merging or selecting dictionary items from an initially large dictionary. Taking a step further, some researchers designed dictionaries for the situations that the present training instances are different from the testing instances. For instance, Zhu et al. [56] utilized weakly labelled data from other visual domains as the auxiliary source data for enhancing the original learning system and their framework only requires a small set of labelled samples in the source domain.

Nevertheless, most dictionary learning algorithms are of high time complexity and converge slowly; moreover, they may get trapped in local minimum [55]; consequently, efficient dictionary learning on large-scale data still remains a challenging task. It is

worth noting that submodular optimization has become a sensible trend to solve large-scale problems in computer vision and as an example, Jiang et al. [55] exploited submodularity and monotonicity properties of object function to construct a dictionary from a set of dictionary item candidates. Possible future work includes exploring supervised dictionary learning methods to propose an efficient method to learn a compact and discriminative dictionary for large-scale people counting application. When the dictionary is large and the data dimension is high, learning sparse representations is a computationally challenging problem and this is the direction in which we would like to extend our work.

Furthermore, although the literature of image classification is predominated by local and global hand-crafted features; deep learning methodologies have been utilized recently to obtain machine-learned features for image classification and have shown great potential in various visual recognition tasks. In this paper, the features obtained by CNN yields encouraging results. Meanwhile, various architectures and techniques have been proposed to enhance the learning capacity. Recently, Shao et al. [57] proposed a multispectral neural network to learn features from multicolumn deep neural networks to obtain an effective low-dimensional embedding, which led to a more discriminative feature than that of CNN. However, deep architectures often require a large amount of labelled data for supervised training; their training would take very long time and a large number of hyper-parameters should be tuned. Alternatively, proposing an optimal solution can be considered as a generalized way to extract the most meaningful features for any user-defined application. As a successful example, in [58] the authors developed an evolutionary learning methodology to automatically generate domain-adaptive global feature descriptors for image classification using multiobjective genetic programming. An intriguing question for future work is whether this framework could be useful for people counting/detection applications.

7. Conclusion

In this paper, we proposed an extremely accurate and scalable people counting method based on sparse representation. Sparsity provides a powerful tool for inferring high-dimensional image data that have complex low-dimensional structure. Methods like l_1 -minimization offer computational tools to extract such structures and help harness the semantic of data. In order to reduce the computation complexity of l_1 -minimization solvers for finding such sparse representation, random projection is employed as a fast and simple dimensionality reduction method which preserves the similarities of the data vectors well. According to our extensive

- [48] S. Choudri, J.M. Ferryman, A. Badii, Robust background model for pixel based people counting using a single uncalibrated camera, in: 2009 Twelfth IEEE International Workshop on Performance Evaluation of Tracking and Surveillance (PETS-Winter), IEEE, 2009, pp. 1–8.
- [49] M. Patzold, R.H. Evangelio, T. Sikora, Counting people in crowded environments by fusion of shape and motion information, in: 2010 Seventh IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), IEEE, 2010, pp. 157–164.
- [50] D. Conte, P. Foggia, G. Percannella, F. Tufano, M. Vento, A method for counting people in crowded scenes, in: 2010 Seventh IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), IEEE, 2010, pp. 225–232.
- [51] V.B. Subburaman, A. Descamps, C. Carincotte, Counting people in the crowd using a generic head detector, in: 2012 IEEE Ninth International Conference on Advanced Video and Signal-Based Surveillance (AVSS), IEEE, 2012, pp. 470–475.
- [52] J. Ferryman, A.-L. Ellis, Performance evaluation of crowd image analysis using the PETS2009 dataset, *Pattern Recognit. Lett.* 44 (2014) 3–15.
- [53] H. Bay, T. Tuytelaars, L. Van Gool, Surf: speeded up robust features, in: *Computer Vision—ECCV 2006*, Springer, Berlin, Heidelberg, 2006, pp. 404–417.
- [54] J. Tang, L. Shao, X. Li, Efficient dictionary learning for visual categorization, *Comput. Vis. Image Underst.* 124 (2014) 91–98.
- [55] Z. Jiang, G. Zhang, L.S. Davis, Submodular dictionary learning for sparse coding, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, IEEE, 2012, pp. 3418–3425.
- [56] F. Zhu, L. Shao, Weakly-supervised cross-domain dictionary learning for visual recognition, *Int. J. Comput. Vis.* 109 (1–2) (2014) 42–59.
- [57] L. Shao, D. Wu, X. Li, Learning Deep and Wide: A Spectral Method for Learning Deep Networks.
- [58] L. Shao, L. Liu, X. Li, Feature Learning for Image Classification via Multi-objective Genetic Programming.

Homa Foroughi received her bachelor degree from BuAli Sina University, Iran, in 2004 and the M.Sc. degree from Ferdowsi University, Iran, in 2009, both in computer engineering. She is currently working toward the Ph.D. degree in computing science at University of Alberta, Canada. Her research interests include computer vision, image processing and machine learning. She is a student member of IEEE.

Nilanjan Ray received his bachelor degree in mechanical engineering from Jadavpur University, Calcutta, India, in 1995, master's degree in computer science from the Indian Statistical Institute, Calcutta, in 1997, and Ph.D. in electrical engineering from the University of Virginia, Charlottesville, in May, 2003. After having 2 years of postdoctoral research and a year of industrial work experience, he joined the department of Computing Science, University of Alberta in July 2006 and now, he is an associate professor at Computing Science, University of Alberta. His research area is image and video analysis, segmentation, object detection, image classification and object tracking.

Hong Zhang received the B.Sc. degree from North-eastern University, Boston, USA, in 1982, and the Ph.D. degree from Purdue University, West Lafayette, IN, USA, in 1986, both in electrical engineering. He conducted post-doctoral research at the University of Pennsylvania from 1986 to 1987 before joining the Department of Computing Science, University of Alberta, Canada, where he is currently a Professor and Director of the Centre for Intelligent Mining Systems. He is an NSERC Industrial Research Chair. His current research interests include robotics, computer vision, and image processing.