



Locality constraint distance metric learning for traffic congestion detection



Qi Wang*, Jia Wan, Yuan Yuan

School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an, China

ARTICLE INFO

Article history:

Received 14 November 2016
Revised 21 January 2017
Accepted 24 March 2017
Available online 29 March 2017

Keywords:

Distance metric learning
Locality constraint
Kernel regression
Traffic congestion analysis

ABSTRACT

In this paper, a locality constraint distance metric learning is proposed for traffic congestion detection. First of all, an accurate and unified definition of congestion is proposed and the congestion level analysis is treated as a regression problem in the paper. Based on that definition, a dataset consists of 20 different scenes is constructed for the first time since the existing dataset is not diverse for real applications. To characterize the congestion level in different scenes, the low-level texture feature and kernel regression is utilized to detect traffic congestion level. To reduce the influence among different scenes, a Locality Constraint Distance Metric Learning (LCML) which ensured the local smoothness and preserved the correlations between samples is proposed. The extensive experiments confirm the effectiveness of the proposed method.

© 2017 Elsevier Ltd. All rights reserved.

1. Introduction

As the development of the society, crowd and traffic has become more and more congested around the world and cause many problems [1,2]. The traffic jams not only waste our time and resources, but also create more pollutions and accidents. If the traffic congestion level can be automatically detected, it will be easier to relieve the traffic jams.

There are three popular devices that can be utilized to detect congestion level. The first one is Loop Detector [3] whose installation and maintenance is complicated and the detection range is very limited. The second one is GPS based smart cell phone and vehicle. This method is popular because of the wide detection range, but the precision of the congestion level detection is low. The last one is camera. The cameras are very popular nowadays, so, the detection range could be ensured. What's more, the detection could be very precise.

Many approaches are proposed to detect congestion level from videos [4,5]. Most of them detect congestion by analyzing the key points or moving areas. The number of key points or moving blobs and their speed can be used as the feature to estimate the congestion. These methods are effective to answer whether it is a traffic jam but can't determine the exact congestion level. Besides, these methods solve this problem in only one specific scene which is not useful for real applications.

Unfortunately, the research of traffic congestion analysis is limited because there exist many problems. The first one is that the researchers treated congestion detection as a classification task which simplifies the problem but limits its usage. The second is that there exists no large scale dataset containing different scenes for traffic congestion analysis. The most challenge is the different illumination, occlusion level, camera angles and road conditions in various scenes. These problems make the traffic congestion status hard to analyze.

To remedy these problems, an accurate and unified definition is first proposed. With this definition, the congestion analysis becomes a regression problem since the label of congestion level is a real number between 0 and 1. At the same time, a dataset consists of 20 different scenes is conducted to serve as the platform for congestion analysis.

To reduce the effect of different conditions which makes the congestion level hard to analyze, the metric learning [6,7] is utilized. However, traditional metric learning for regression predicts the value of a sample by measuring its distance to all the other training samples. But, the performance is bad since the generalization ability is limited and different scenes will affect the prediction. Thus, a locality constraint metric learning is proposed in which only several nearest neighbors are considered.

The contributions of the paper are summarized as follows:

- An accurate and unified definition of congestion is proposed. With the definition, the question of how congested the traffic status could be answered. Besides, the congestion level in different scenes can be compared in a unified framework.

* Corresponding author.

E-mail address: crabwq@gmail.com (Q. Wang).

- A dataset consisting of 20 different scenes is constructed for the first time to promote the research of traffic congestion analysis. The videos in this dataset contain different illumination, camera angles and road conditions. This dataset can serve as a platform for the research of congestion level regression.
- A locality constraint metric learning which ensured the local smoothness and preserved the correlations between samples is proposed for congestion level regression. Since the difference among various scenes affects the performance of congestion regression, this approach is proposed to reduce the defect across different scenes through constraining that only the neighbors of a testing sample can contribute to the prediction.

The remaining part of this paper is as follows. Relevant works are reviewed in Section 2 and the definition and dataset are presented in Section 3. Then, the details of the proposed method are elaborated in Section 4. After the experimental results are reported and discussed in Section 5, the conclusion and future works are presented in Section 6.

2. Related work

In this section, the relevant works of congestion detection and metric learning are briefly reviewed.

2.1. Congestion detection

The congestion detection algorithms can be divided into two categories. The first category is based on the analysis of key points and moving areas. Another is based on direct feature extraction and classification.

The assumption of the first category is that more congested traffic scenes contain more moving objects. Hu et al. [8] proposed an algorithm that classifies congestion videos based on the segmentation of moving vehicles. Firstly, the moving objects are segmented by background subtraction method [9]. Then, the speed of the moving blobs are calculated by Optical Flow [10]. Finally, the percentage of moving blobs and their speed can be served as the features. Fuzzy logic is utilized for the final decision. Sobral et al. [11] proposed an approach which combined the features of key points and moving blobs together. The crowd density is first evaluated by background subtraction algorithm. Then, the speed is estimated by Kanade–Lucas–Tomasi (KLT) algorithm [12]. These methods rely on the preprocessing like Background Subtraction and Tracking. Thus, the performance of these congestion detection methods is limited because of the uncertain preprocessing.

The purpose of the other category is to design congestion related features. Derpanis et al. [13] proposed the Spatialtemporal Orientation Analysis feature motivated by the visual dynamics of congested scenes. Riaz et al. [14] encoded motion information by analyzing the statistics of motion vectors. To combine the appearance and motion information together, Dallalzadeh et al. [15] proposed the symbolic representation. These methods don't rely on the preprocessing algorithms which make them work well in a specific scene. However, how to design the features that can cross congested scenes is still a challenging task.

2.2. Metric learning

Metric learning is the task to learn a good distance measurement. The aim of the metric learning is to minimize the distance of samples from the same class and maximize the distance from different classes [16].

Most of the distance learning algorithms are designed for classification task [17,18], such as image retrieval [19], person re-identification [20] and face recognition [21,22]. Large Margin Nearest Neighbor (LMNN) [23,24] is proposed to learn a Mahalanobis

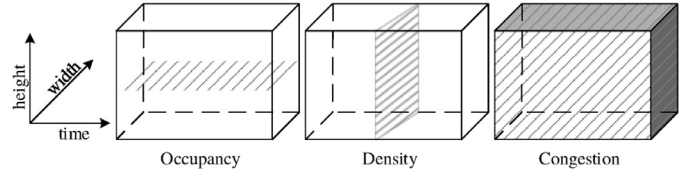


Fig. 1. The proposed definition is based on the time-space congestion.

distance for K Nearest Neighbor (KNN) classification. Information-Theoretic Metric Learning (ITML) [25] is also proposed for KNN classification. The difference is that the distribution parameterized by distance metric M is regularized to be closed to a prior distribution. Neighborhood Component Analysis (NCA) [26] is another linear metric learning algorithm which optimized the classification performance based on the Leave-one-out validation. To exploit the negative constraint lacking, Discriminative Component Analysis (DCA) [27] is proposed.

The metric learning can also be included into regression task. Metric Learning for Kernel Regression (MLKR) [28] is proposed to learn a distance metric for kernel regression. Sparse representation has been proved effective in many applications [29]. To combine the sparsity into the framework, Kernel Regression with Sparse Metric Learning (KRSML) [30] is proposed to regularized the distance metric with a mixed (2,1)-norm. Xiao et al. [31] proposed an application which utilized metric learning for human age estimation.

Recently, many deep metric learning methods are proposed with the development of deep learning. Hu et al. [32] proposed a deep metric learning method to compare the similarity of two faces by minimizing the intra-class variation and maximize the inter-class variation. Li et al. [33] proposed to learn a suitable metric with the help of community-contributed images. Song et al. [34] proposed a novel structural objective function on the lifted problem which is proved to be effective for image retrieval.

3. The definition and dataset

In this section, an unified and accurate definition is first presented. And then, the detail of the traffic congestion dataset and the corresponding labeling method are introduced.

3.1. Definition

The congestion level is defined as the occupancy of moving objects in the domain of space-time. In literature, the congestion can be measured by spatial congestion and temporal congestion which are called density and occupancy respectively [3]. As shown in Fig. 1, the density can only measure the congestion at a point of time, while the occupancy can only measure the congestion at a point of space. The proposed definition of congestion considers the spatial and temporal information simultaneously. This definition is accurate and unified, so the comparison of congestion level in different scenes becomes possible.

Formally, the congestion can be expressed as follows:

$$\text{congestion} = \frac{\sum_{x,y,t} f(x,y,t)}{\text{width} \times \text{length} \times \text{time}} \quad (1)$$

where $\text{congestion} \in (0, 1)$ is the congestion level, and time is the number of frames in a video clip. width and length are the width and length of a road. $f(x, y, t)$ is defined as:

$$f(x, y, t) = \begin{cases} 1, & \text{occupied} \\ 0, & \text{not occupied} \end{cases} \quad (2)$$

which indicates that whether a point is occupied by a moving object.



Fig. 2. Typical images of the traffic congestion dataset.

3.2. Traffic congestion dataset

Since the existing datasets are not diverse for real applications, a new dataset which consists of 20 different scenes is constructed. First, large amount of videos containing different streets and weather conditions are collected. The resolution of these videos are equal or larger than 1080×720 . The average length of these videos is 30 minutes. Different direction on same road is treated as different scenes. Typical images are shown in Fig. 2.

With the definition, the calculation of the real congestion level needs a pixel-wise labeling which is time-consuming. To remedy this, we suppose that vehicles and lanes have the same width. The assumption is based on a fact that there is only one row of vehicles on each lane even under extremely crowded circumstance at most of the time. With this assumption, the congestion level can be reduced as follows:

$$\text{congestion} = \frac{\sum_{y,t} f(y,t)}{\text{height} \times \text{time}} \quad (3)$$

where *height* can be seen as the length of the road.

With the simplification, the moving objects and road can be represented as a line along them and the congestion level can be calculated easily. The visualization of labeling is shown in Fig. 3. Specifically, the length of a lane on the road is represented by the length of a line along it. Similarly, the length of vehicles can be estimated in the same way. Then, the congestion level can be calculated by the fraction of total length of vehicles and the total length of lanes.

Since the perspective transformation can heavily affect the real congestion level, it is considered in the labeling procedure. The perspective transformation causes that the vehicles far from the camera will be smaller than the vehicles close to the camera. Motivated by this, the weight of pixels at the top of images (where the vehicles are far from the camera) should be larger. Specifically, the weight of pixels is decided by the width of road in images since the width of road is affected by perspective transformation as well. After the perspective is considered, the variation of the congestion level caused by perspective transformation is reduced. As shown

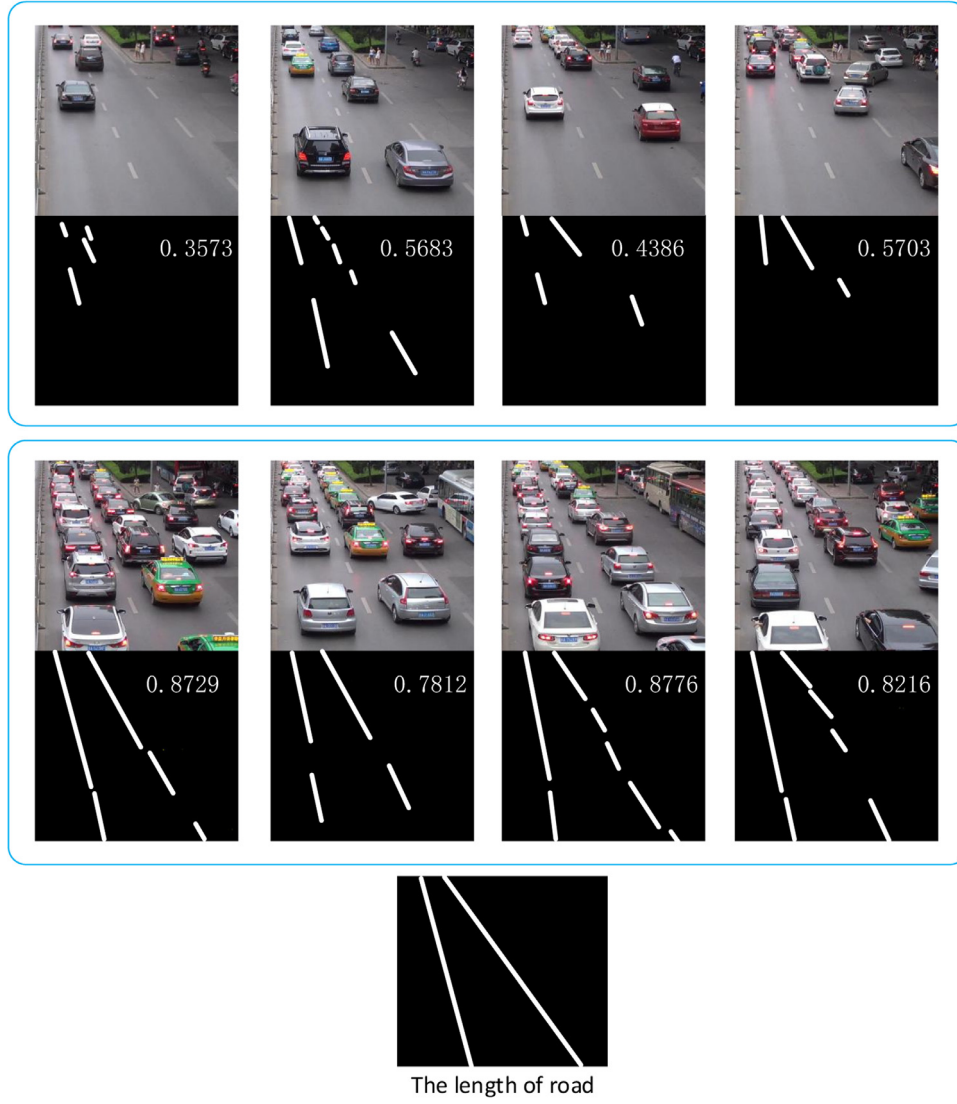


Fig. 3. The visualization of labeling method. In this figure, two different scenes is used for examples. The numbers on the binary image indicate the congestion level. In the proposed labeling method, each vehicle can be represented by a line along it. Since there are two lanes in the road, the length of road is represented by two lines as the bottom image shows.

in Fig. 4, the congestion level is more smooth after the perspective transformation is taken into consideration.

4. Our method

In this section, the details of the proposed method are presented. First of all, the texture feature of image is extracted as low-level features. Then, a distance metric is learned by the proposed locality constraint metric learning algorithm with the precomputed features. Based on the learned metric, the congestion level can be efficiently detected by kernel regression.

4.1. Locality constraint distance metric learning

The influence among different scenes makes the prediction of congestion level a hard task. The locality has been proved to be effective for face recognition [35] and subspace learning [36]. In this paper, the locality constraint is combined with distance metric learning to reduce the effect of different scenes. Different from traditional metric learning, this method only takes the neighbors of the testing sample into consideration, since the testing sample and its neighbors have high probability of belonging to the same scene.

That efficiently reduces the influences among different scenes as shown in Fig. 5.

The kernel regression can be treated as the weighted sum of training samples. The weight should be related to the similarity between samples. Formally, given a feature vector x_i of a sample, the corresponding congestion level \hat{y}_i can be calculated as:

$$\hat{y}_i = \frac{\sum_{j \neq i} y_j k_{ij}}{\sum_{j \neq i} k_{ij}}, \quad (4)$$

where k_{ij} refers to the kernel function which is defined as:

$$k_{ij} = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{d_{ij}}{\sigma}\right), \quad (5)$$

where d_{ij} is the distance of x_i and x_j . Note that, σ is fixed to 1 for simplification. Usually, the Euclidean distance is included as the distance measurement. However, many works have shown that the learned metric can generate better performance. Thus, the distance metric learning is utilized to learn a better distance measurement (Fig. 6).

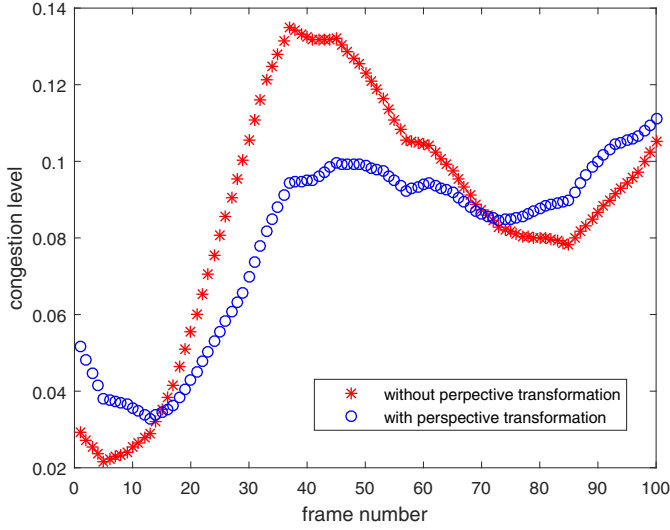


Fig. 4. The congestion level is more smooth after the perspective transformation is taken into consideration. In this figure, the horizontal axis is the frame number and the vertical axis is the congestion level. The red stars are congestion level without perspective transformation. The blue circles are congestion level with perspective transformation. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

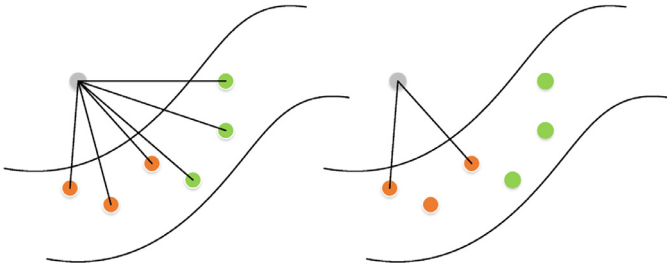


Fig. 5. The illustration of locality constraint metric learning. Suppose there are 6 samples in the training set which are divided into 2 classes (orange and green). The Metric Learning for Kernel Regression in the left image uses all the training samples to predict a new sample which will suffer from the influence among different scenes. The proposed method in the right image can eliminate such an influence.. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

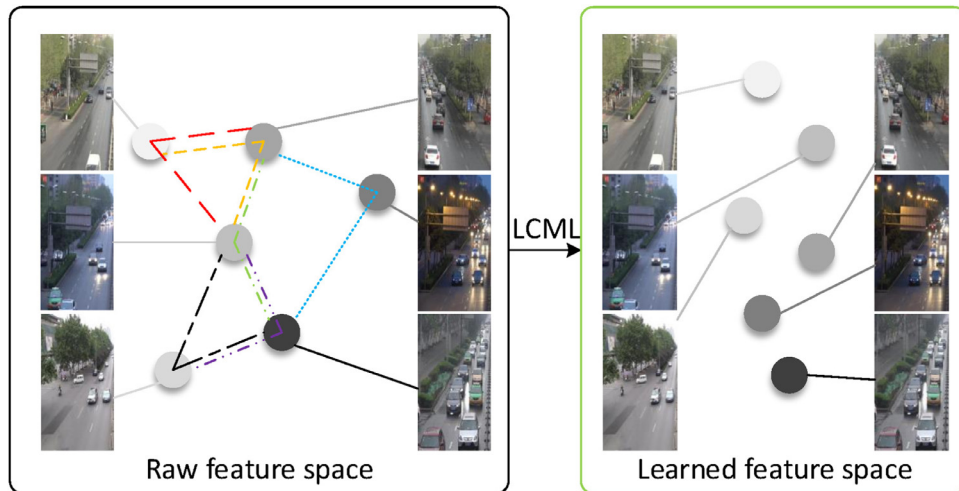


Fig. 6. The proposed metric learning algorithm constrains that the value of a sample can only be approximated by its neighbors, which minimizes the influence among different scenes. In this figure, each point represents the feature of the congestion level in the feature space. In LCML, only neighbors are considered for prediction as the dotted lines show.

In metric learning, the distance between two vector x_i and x_j is calculated as follows:

$$d_{ij} = (x_i - x_j)^T M (x_i - x_j), \quad (6)$$

where M is the learned distance metric that can transform the features to learned space and produce better performance. Note that M needs to be preserved to semi-definite, which is hard to satisfy. Motivated by [28], M can be decomposed to $A^T A$. Then, Eq. 6 can be reformulated as follows:

$$d_{ij} = \|A(x_i - x_j)\|^2. \quad (7)$$

To learn the distance metric A , the mean squared error between the ground truth and the prediction can be used as the loss function:

$$\mathcal{L}_{mse} = \sum_i \left((y_i - \hat{y}_i)^2 + \beta \sum_{j \neq i} (k_{ij} \times d_{ij}) \right), \quad (8)$$

where y_i is the ground truth and \hat{y}_i is the prediction. d_{ij} is the distance of two samples and k_{ij} can be treated as the weight in the kernel regression in Eq. 4. This regularization term is used to punish the weights of x_j that is far from x_i .

To ensure the divergence invariant under scaling of the feature space, the LogDet divergence [25] is utilized. Then, the final loss function is as below:

$$\mathcal{L} = \mathcal{L}_{mse} + \beta D_{ld}(A^T A, M_0), \quad (9)$$

where M_0 is the prior metric which is set as identity matrix in our experiments. The $D_{ld}(M, M_0)$ is defined as:

$$D_{ld}(M, M_0) = \text{tr}(MM^{-1}) - \log \det(MM_0^{-1}) - d. \quad (10)$$

Motivated by [37], $D_{ld}(M, M_0)$ can be replaced by:

$$D_{ld}(M, M_0) = \text{tr}(MM^{-1}) - \log \det(M). \quad (11)$$

4.2. Approximation and optimization

To minimize the problem in Eq. 8, an approximation is proposed. The aim of Eq. 8 is to find out some neighbors of sample x_i in the training set and give a correct prediction through the weighted average of these neighbors. That suggests we can use K neighbors for prediction instead of all samples. This greatly reduced the influence among different scenes. Thus, Eq. 1 can be reformulated as:

$$\hat{y}_i = \frac{\sum_{j \in N(i)} y_j k_{ij}}{\sum_{j \in N(i)} k_{ij}}. \quad (12)$$

where $N(i)$ is the set of neighbors of x_i . Consequently, \mathcal{L}_{mse} loss can be reduced as follows:

$$\mathcal{L}_{mse} = \sum_i (y_i - \hat{y}_i)^2. \quad (13)$$

This final problem can be efficiently solved by gradient decent algorithm and preserve the locality at the same time. To solve this optimization problem, the gradient of \mathcal{L} with respect to A is calculated as follows:

$$\frac{\partial \mathcal{L}}{\partial A} = \frac{\partial \mathcal{L}_{mse}}{\partial A} + \beta \frac{\partial \mathcal{D}_{ld}(A^\top A, M_0)}{\partial A}, \quad (14)$$

where $x_{ij} = x_i - x_j$ and

$$\frac{\partial \mathcal{L}_{mse}}{\partial A} = 4A \sum_i (\hat{y}_i - y_i) \sum_{j \in N(i)} (\hat{y}_j - y_j) k_{ij} x_{ij} x_{ij}^\top, \quad (15)$$

and

$$\frac{\partial \mathcal{D}_{ld}(AA^\top, M_0)}{\partial A} = 2A(M_0^{-1} + (A^\top A)^{-1}) \quad (16)$$

The details of the gradient decent procedure are shown in Algorithm 1. In this algorithm, X is the feature matrix containing n samples in \mathbb{R}^d and Y is the corresponding congestion level. The learned transformation A is obtained by minimizing the loss function \mathcal{L} defined in Eq. 9.

Algorithm 1 Locality constraint metric learning.

Input: X : feature matrix containing n samples in \mathbb{R}^d ;
 Y : n corresponding congestion level.

```

1: repeat
2:   Calculate  $\nabla A \leftarrow \frac{\partial \mathcal{L}}{\partial A}$  via Equation 14
3:   Initialize  $A_{best} \leftarrow 0$ ,  $\mathcal{L}_{best} \leftarrow 0$ 
4:   Calculate  $A' \leftarrow A - \delta \nabla A$ 
5:   Calculate  $\mathcal{L}' \leftarrow \mathcal{L}(A')$  via Equation 9
6:   if  $\mathcal{L}' < \mathcal{L}_{best}$  then
7:     Update  $A_{best}$ :  $A_{best} \leftarrow A'$ 
8:   end if
9: until  $A$  is convergence.

```

Output: The learned transformation A .

5. Experiments

In this section, extensive experiments are conducted to confirm the effectiveness of the proposed method. Firstly, the experimental parameters are selected through cross-validation. Then, the proposed method is compared to some traditional algorithms. After that, the effectiveness of the feature and classifier is evaluated. Lastly, the superiority of locality constraint metric learning is confirmed.

5.1. Parameter settings

The constructed dataset consists of 20 different scenes which has different lights, occlusions and road conditions. The average length of these videos is 30 minutes. The resolutions are 1280×720 and 1920×1080 . Typical images are shown in Fig. 2.

All experiments are performed for ten times and 5000 samples are used for training, 1500 for testing each time. Then, the average result is reported. The number of neighbors K is set as 100 and 10 for training and testing based on the result shown in Fig. 7. Note that the Mean Squared Error (MSE) is employed for the evaluation. The lower MSE indicates better performance. The parameter β in Eq. 9 is set as 0.1 through cross validation. The result is shown in Fig. 8.

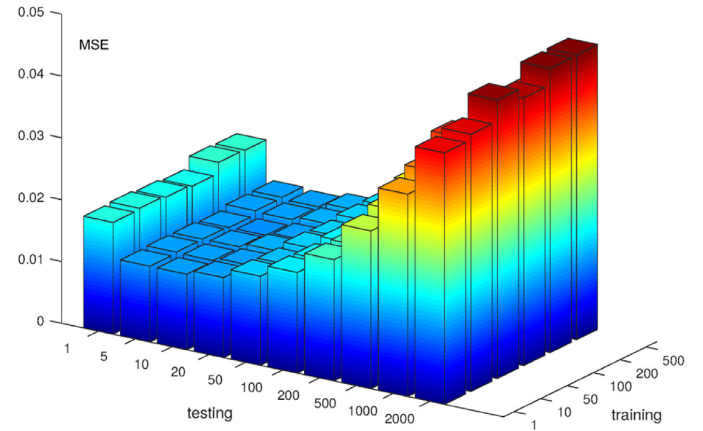


Fig. 7. The axes are K in the training, K in the testing, and mean squared error. For the best performance, K is set as 100 and 10 for training and testing. Note that the lower MSE indicates better performance.

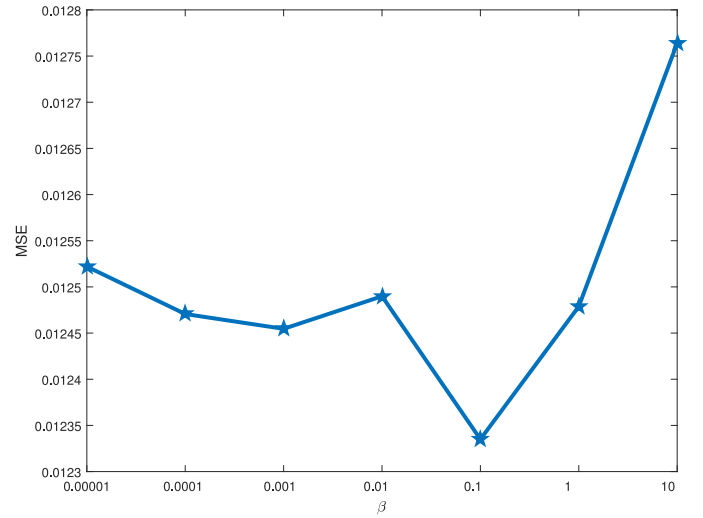


Fig. 8. The selection of parameter β . In this figure, the horizontal axis is the different β and the vertical axis is the mean squared error. Note that the lower MSE indicates better performance.

Table 1

Comparison of Linear Regression (LR), Kernel Regression (KR), Local Linear Embedding (LLE), Robust Principal Component Analysis (RPCA), Metric Learning for Kernel Regression (MLKR) and the proposed Locality Constraint Distance Metric Learning (LCML). MAE is mean absolute error and MSE is mean squared error.

Methods	LR	KR	LLE	RPCA	MLKR	LCML
MAE	0.152	0.082	0.136	0.082	0.087	0.076
MSE	0.036	0.013	0.032	0.012	0.014	0.010

5.2. The evaluation of the proposed method

In this section, the proposed method is compared to some traditional algorithms. The comparison methods are Linear Regression (LR), Kernel Regression (KR) [38] and Metric Learning for Kernel Regression (MLKR) [28]. In this experiment, the low-level texture feature is utilized as the representation of congestion level. The final results comparison is shown in Table 1.

As shown in Table 1, the performance of kernel regression is better than linear regression, which indicates that a non-linear regressor outperforms a linear regressor for congestion detection. It is unsurprising since the cross scenes congestion level detection is a non-linear regression problem. However, it is surprising that

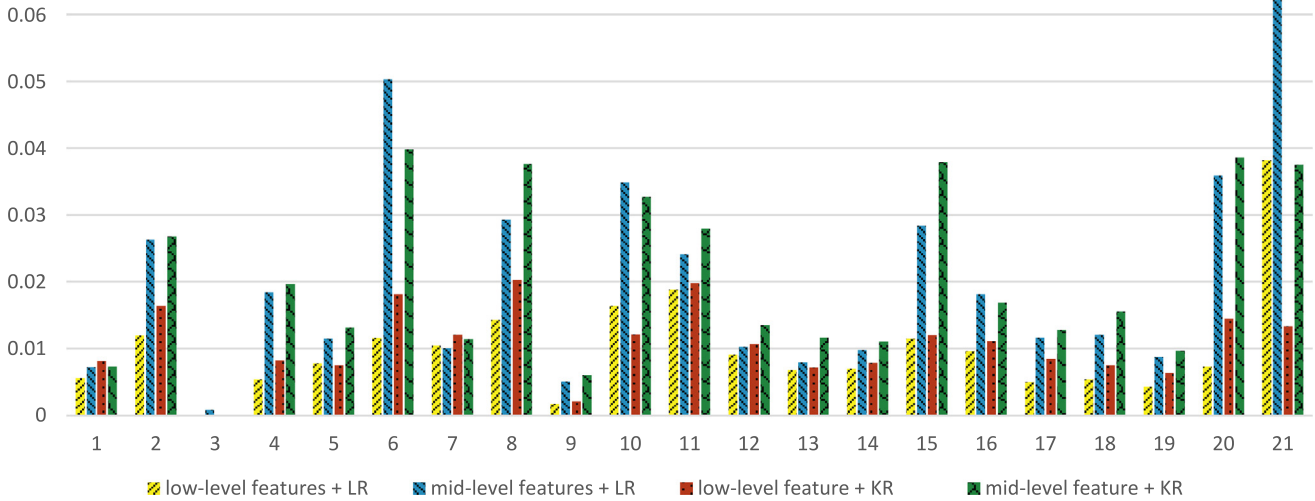


Fig. 9. The comparison of different features and classifiers. In this figure, 1 – 20 of the horizontal axis are the results of the experiments in which the training and testing are performed on only one scene. 21 of the horizontal axis is the result of the experiment in which the whole dataset (20 different scenes together) is included for training and testing.

Metric Learning for Kernel Regression is not superior to kernel regression. The result indicates that metric learning is ineffective under the influence of different scenes. After the locality constraint is added to the loss function, metric learning shows its superiority. Thus, the influence among different scenes will affect the prediction of Metric Learning for Kernel Regression, and locality constraint can reduce that influence efficiently.

The reason why locality constraint metric learning outperforms other methods is that the locality is preserved during the learning of distance metric. The locality is guaranteed by constraining that only neighbors of a sample can be included for prediction. With this restraint, the similar congestion scenes will have similar predictions since their neighbors should be similar as well.

5.3. The effectiveness of feature and classifier

To confirm that the feature is useful for congestion regression, the low-level feature [39] and the mid-level feature [40] are included as the comparison features. In this experiment, the Local Binary Patterns (LBP) is utilized as low-level feature. The extraction of mid-level feature consists of dictionary learning, feature encoding and feature pooling. Specifically, the Batch K -means [41,42] is included for dictionary learning, the sparse coding [43] is utilized to encode image descriptor and max-pooling is used to construct final representation. The results are shown in Fig. 9.

First of all, if the training and testing is performed on only one scene, the performance is better than the training and testing across different scenes. As shown in Fig. 9 that the most errors of 1 – 20 (single scene, and remind that the constructed dataset contains 20 different scenes) are lower than the error of 21 (multiple scenes). That confirms the influence of different scenes will drop the congestion detection performance.

Secondly, the low-level feature is superior to the mid-level features [44]. The mid-level feature concentrates on global information which is ineffective to detect congestion. For example, the mid-level features of different congestion images come from the same scene might be very similar which makes it harder to detect congestion level. In contrast, the concentration of low-level texture feature is local information which is more useful for congestion detection. For instance, the low-density images tend to present coarse texture, and the high-density images tend to present fine texture. [45]. Note that, the performance can be further boosted if

we can segment road [1] area to eliminate the influence of background.

Finally, the most popular regression method is Linear Regression (LR). However, Kernel Regression (KR) is shown to be superior than LR in the experiment. That indicates that the non-linear classifier is better than linear classifier for congestion level regression. It is unsurprising since the cross scene congestion level detection is a non-linear regression problem. The LLE is a typical manifold learning which can be used for dimensional reduction with the locality preserved. However, the manifold learning is based on an assumption that the high-dimensional data points lie on a low-dimension manifold which might not be satisfied in the training data. Thus, the performance of LLE is rather poor. The RPCA with low-rank and sparsity constraints achieves good performance which indicates that the sparsity is important. The proposed method achieves the best performance since the structural information among samples can be further exploited based on the sparsity. It is surprising that Metric Learning for Kernel Regression is not superior to kernel regression. The result indicates that metric learning is ineffective under the influence of different scenes. After the locality constraint is added to the loss function, metric learning shows its superiority. Thus, the influence among different scenes will affect the prediction of Metric Learning for Kernel Regression, and locality constraint can reduce that influence efficiently.

5.4. The effectiveness of locality constraint metric learning

To confirm the effectiveness of locality constraint metric learning, the Metric Learning for Kernel Regression (MLKR) [28], Local Linear Embedding (LLE) [46,47], and Robust Principal Component Analysis (RPCA) [48,49] are included as the comparison. A visualization of the regression results is shown in Fig. 11.

As shown in Fig. 11, the performance of linear regression is the worst since the congestion regression is a non-linear problem. The performance of Metric Learning for Kernel Regression is better. However, the influence among different scenes makes the prediction a hard work. The proposed locality constraint metric learning outperforms other methods since the effect among different scenes is reduced. Note that, the proposed method can be further improved by graph methods [50] for big data.

To further confirm that the locality constraint metric learning can efficiently reduce the influence of different scenes. An image retrieval experiment is conducted. In this experiment, 9 images



Fig. 10. Most of the neighbors and queries are from the same scene and similar congestion level with LCML. Only the neighbors which have same scene and similar congestion level are considered in LCML. Thus, the influence of different scenes is reduced and the performance of cross scene congestion level detection is improved. In this figure, the first column is queries and the rest are nearest neighbors. The images in red boxes are failures. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

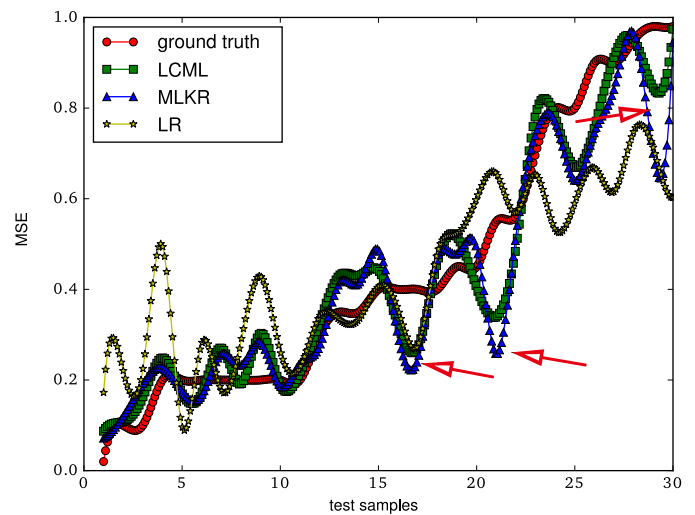


Fig. 11. The visualization of the regression results. The proposed method outperforms MLKR, as the red arrows point. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

from different scenes are randomly selected as the queries. 3 nearest neighbors of these queries are selected from the whole training set (containing 20 different scenes) with the proposed LCML. The result is shown in Fig. 10.

As shown in Fig. 10, most of the neighbors and the respectively queries are from the same scene and similar congestion level except for 3 failures which come from complex scenes that have no obviously road boundary. If all the training samples are included to predict the congestion level, the different scenes will influence each other, which decreases the detection performance. Only the neighbors which have same scene and similar congestion level are considered in LCML. Thus, the influence of different scenes is reduced and the performance of congestion level detection is increased with LCML.

6. Conclusion

To remedy the congestion detection problem, a locality constraint metric learning is proposed to reduce the influence among different scenes. An unified and accurate definition of congestion is first proposed to better describe the traffic congestion level. Based on the definition, a dataset consisting of 20 different scenes is constructed to serve as a platform for congestion detection. To solve this problem, the low-level texture feature and kernel regression which outperform mid-level feature and linear regression are included as the feature and classifier. Since the influence of different scenes makes the detection of traffic congestion a difficult task, a locality constraint metric learning which ensured the local smoothness and preserved the correlations between samples is proposed to reduce such an influence. The extensive experiments confirm the effectiveness of the proposed method.

In the future, some density based features will be exploited to better represent the congestion level. Besides, a hierarchical model will be exploited as well since the congestion level is a high level semantic conception.

Acknowledgment

This work was supported by the National Natural Science Foundation of China under Grant 61379094.

References

- [1] Y. Yuan, J. Fang, Q. Wang, Online anomaly detection in crowd scenes via structure analysis, *IEEE Trans. Cybern.* 45 (3) (2015) 548–561.
- [2] Y. Yuan, D. Wang, Q. Wang, Anomaly detection in traffic scenes via spatial-aware motion reconstruction, *IEEE Trans. Intell. Transp. Syst.* (2016), doi:10.1109/TITS.2016.2601655.
- [3] F.L. Hall, Traffic stream characteristics, *Traffic Flow Theory*, US Federal Highway Administration, 1996.
- [4] L. Li, L. Chen, X. Huang, J. Huang, A traffic congestion estimation approach from video using time-spatial imagery, in: *International Conference on Intelligent Networks and Intelligent Systems*, 2008, pp. 465–469.
- [5] F. Porikli, X. Li, Traffic congestion estimation using hmm models without vehicle tracking, in: *IEEE Intelligent Vehicles Symposium*, 2004, pp. 188–193.
- [6] J. Lu, G. Wang, P. Moulin, Localized multifeature metric learning for image-set-based face recognition, *IEEE Trans. Circuits Syst. Video Technol.* 26 (3) (2016) 529–540.
- [7] J. Lu, X. Zhou, Y.-P. Tan, Y. Shang, J. Zhou, Neighborhood repulsed metric learning for kinship verification, *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (2) (2014) 331–345.
- [8] S. Hu, J. Wu, L. Xu, Real-time traffic congestion detection based on video analysis, *J. Inf. Comput. Sci.* 9 (10) (2012) 2907–2914.
- [9] C. Zhan, X. Duan, S. Xu, Z. Song, M. Luo, An improved moving object detection algorithm based on frame difference and edge detection, in: *International Conference on Image and Graphics*, 2007, pp. 519–523.
- [10] B.K. Horn, B.G. Schunck, Determining optical flow, *Artif. Intell.* 17 (1–3) (1981) 185–203.
- [11] A. Sobral, L. Oliveira, L. Schnitman, F.D. Souza, Highway traffic congestion classification using holistic properties, in: *International Conference on Signal Processing, Pattern Recognition and Applications*, 2013.
- [12] B.D. Lucas, T. Kanade, An iterative image registration technique with an application to stereo vision, in: *Proceedings of the International Joint Conference on Artificial Intelligence*, 1981, pp. 674–679.
- [13] K.G. Derpanis, R.P. Wildes, Classification of traffic video based on a spatiotemporal orientation analysis, in: *IEEE Workshop on Applications of Computer Vision*, 2011, pp. 606–613.
- [14] A. Riaz, S.A. Khan, Traffic congestion classification using motion vector statistical features, in: *International Conference on Machine Vision*, 2013. 90671A–90671A.
- [15] E. Dallalzaheh, D. Guru, B. Harish, Symbolic classification of traffic video shots, in: *Advances in Computational Science, Engineering and Information Technology*, 2013, pp. 11–22.
- [16] Y. He, Y. Mao, W. Chen, Y. Chen, Nonlinear metric learning with kernel density estimation, *IEEE Trans. Knowl. Data Eng.* 27 (6) (2015) 1602–1614.
- [17] A. Globerson, S.T. Roweis, Metric learning by collapsing classes, in: *Advances in Neural Information Processing Systems*, 2005, pp. 451–458.
- [18] Z. Huang, R. Wang, S. Shan, X. Li, X. Chen, Log-euclidean metric learning on symmetric positive definite manifold with application to image set classification, in: *Proceedings of the International Conference on Machine Learning*, 2015, pp. 720–729.
- [19] H. Chang, D.-Y. Yeung, Kernel-based distance metric learning for content-based image retrieval, *Image Vision Comput.* 25 (5) (2007) 695–703.
- [20] W. Li, Y. Wu, J. Li, Re-identification by neighborhood structure metric learning, *Pattern Recognit.* 61 (2017) 327–338.
- [21] Z. Huang, R. Wang, S. Shan, X. Chen, Face recognition on large-scale video in the wild with hybrid euclidean-and-riemannian metric learning, *Pattern Recognit.* 48 (10) (2015) 3113–3124.
- [22] Z. Huang, R. Wang, S. Shan, X. Chen, Projection metric learning on grassmann manifold with application to video based face recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 140–149.
- [23] K.Q. Weinberger, J. Blitzer, L.K. Saul, Distance metric learning for large margin nearest neighbor classification, in: *Advances in Neural Information Processing Systems*, 2005, pp. 1473–1480.
- [24] K.Q. Weinberger, L.K. Saul, Distance metric learning for large margin nearest neighbor classification, *J. Mach. Learn. Res.* 10 (2009) 207–244.
- [25] J.V. Davis, B. Kulis, P. Jain, S. Sra, I.S. Dhillon, Information-theoretic metric learning, in: *Proceedings of the International Conference on Machine Learning*, 2007, pp. 209–216.
- [26] J. Goldberger, S.T. Roweis, G.E. Hinton, R. Salakhutdinov, Neighbourhood components analysis, in: *Advances in Neural Information Processing Systems*, 2004, pp. 513–520.
- [27] S.C.H. Hoi, W. Liu, M.R. Lyu, W. Ma, Learning distance metrics with contextual constraints for image retrieval, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2006, pp. 2072–2078.
- [28] K.Q. Weinberger, G. Tesauro, Metric learning for kernel regression, in: *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, 2007, pp. 612–619.
- [29] M. Wang, X. Liu, X. Wu, Visual classification by ℓ_1 -hypergraph modeling, *IEEE Trans. Knowl. Data Eng.* 27 (9) (2015) 2564–2574.
- [30] R. Huang, S. Sun, Kernel regression with sparse metric learning, *J. Intell. Fuzzy Syst.* 24 (4) (2013) 775–787.
- [31] B. Xiao, X. Yang, H. Zha, Y. Xu, T.S. Huang, Metric learning for regression problems and human age estimation, in: *Advances in Multimedia Information Processing*, 2009, pp. 88–99.
- [32] J. Hu, J. Lu, Y.-P. Tan, Discriminative deep metric learning for face verification in the wild, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1875–1882.
- [33] Z. Li, J. Tang, Weakly supervised deep metric learning for community-contributed image retrieval, *IEEE Trans. Multimedia* 17 (11) (2015) 1989–1999.
- [34] H.O. Song, Y. Xiang, S. Jegelka, S. Savarese, Deep metric learning via lifted structured feature embedding, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4004–4012.
- [35] X. Peng, L. Zhang, Z. Yi, K.K. Tan, Learning locality-constrained collaborative representation for robust face recognition, *Pattern Recognit.* 47 (9) (2014) 2794–2806.
- [36] X. Peng, Z. Yu, Z. Yi, H. Tang, Constructing the l2-graph for robust subspace learning and subspace clustering, *IEEE Trans. Cybern.* (2016).
- [37] E.Y. Liu, Z. Guo, X. Zhang, V. Jojic, W. Wang, Metric learning from relative comparisons by minimizing squared residual, in: *IEEE International Conference on Data Mining*, 2012, pp. 978–983.
- [38] H. Takeda, S. Farsiu, P. Milanfar, Kernel regression for image processing and reconstruction, *IEEE Trans. Image Process.* 16 (2) (2007) 349–366.
- [39] F. Lu, J. Huang, An improved local binary pattern operator for texture classification, in: *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2016, pp. 1308–1311.
- [40] Y. Yuan, J. Wan, Q. Wang, Congested scene classification via efficient unsupervised feature learning and density estimation, *Pattern Recognit.* 56 (2016) 159–169.
- [41] J.A. Hartigan, M.A. Wong, Algorithm as 136: a k-means clustering algorithm, *J. R. Stat. Soc. Ser. C (Applied Statistics)* 28 (1) (1979) 100–108.
- [42] A. Coates, A.Y. Ng, Learning feature representations with k-means, in: *Neural Networks: Tricks of the Trade - Second Edition*, 2012, pp. 561–580.
- [43] J. Yang, K. Yu, Y. Gong, T. Huang, Linear spatial pyramid matching using sparse coding for image classification, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 1794–1801.
- [44] F. Li, P. Perona, A bayesian hierarchical model for learning natural scene categories, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 524–531.
- [45] Automatic estimation of crowd density using texture, *Saf. Sci.* 28 (3) (1998) 165–175.
- [46] Y. Bengio, J. Paiement, P. Vincent, O. Delalleau, N.L. Roux, M. Ouimet, Out-of-sample extensions for lle, isomap, mds, eigenmaps, and spectral clustering, in: *Advances in Neural Information Processing Systems*, 2003, pp. 177–184.
- [47] X. Peng, J. Huang, Q. Hu, S. Zhang, A.M. Elgammal, D.N. Metaxas, From circle to 3-sphere: head pose estimation by instance parameterization, *Comput. Vision Image Understanding* 136 (2015) 92–102.
- [48] E.J. Candès, X. Li, Y. Ma, J. Wright, Robust principal component analysis? *J. ACM* 58 (3) (2011) 11:1–11:37.
- [49] X. Peng, S. Zhang, Y. Yang, D.N. Metaxas, PIEFA: personalized incremental and ensemble face alignment, in: *International Conference on Computer Vision*, 2015, pp. 3880–3888.
- [50] M. Wang, W. Fu, S. Hao, H. Liu, X. Wu, Learning on big graph: label inference and regularization with anchor hierarchy, *IEEE Trans. Knowl. Data Eng.* (2017), doi:10.1109/TKDE.2017.2654445.



Qi Wang received the B.E. degree in automation and Ph.D. degree in pattern recognition and intelligent system from the University of Science and Technology of China, Hefei, China, in 2005 and 2010 respectively. He is currently an associate professor with the School of Computer Science and the Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an, China. His research interests include computer vision and pattern recognition.



Jia Wan is currently working toward the M.E. degree in the School of Computer Science and the Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an, China. His current research interests include image classification and congestion analysis.

Yuan Yuan is currently a Full Professor with the School of Computer Science and the Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an, China. She has authored or coauthored over 150 papers, including about 100 in reputable journals such as IEEE Transactions and Pattern Recognition, as well as conference papers in CVPR, BMVC, ICIP, and ICASSP. Her current research interests include visual information processing and image/video content analysis.