

UNIVERSITY OF OTTAWA

MASTER DEGREE THESIS

**Person Re-identification Based on
Kernel Local Fisher Discriminant
Analysis and Mahalanobis Distance
Learning**

Author:

Qiangsen He

Supervisor:

Robert Laganriere

*A thesis submitted in fulfillment of the requirements
for the degree of Master of Applied Science*

in the

VIVA lab

School of Electrical Engineering and Computer Science

February 21, 2017

Declaration of Authorship

I, Qiangsen He, declare that this thesis titled, “Person Re-identification Based on Kernel Local Fisher Discriminant Analysis and Mahalanobis Distance Learning” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

University of Ottawa

Abstract

Faculty of Engineering

School of Electrical Engineering and Computer Science

Master of Applied Science

Person Re-identification Based on Kernel Local Fisher Discriminant Analysis and Mahalanobis Distance Learning

by Qiangsen He

Person re-identification (Re-ID) has become an intense research area in recent years. The main goal of this topic is to recognize and match individuals over time at the same or different locations. This task is challenging due to the variation of illumination, viewpoints, pedestrians' appearance and partial occlusion. Previous works mainly focus on finding robust features and metric learning. Many metric learning methods convert the Re-ID problem to a matrix decomposition problem by Fisher discriminant analysis (FDA). Mahalanobis based metric learning is a popular method to measure similarity; however, since directly extracted descriptors usually have high dimensionality, it's intractable to learn a high-dimensional semi-positive definite (SPD) matrix. Dimensionality reduction is used to project high-dimensional descriptors to a lower-dimensional space while preserving those discriminative information. In this paper, the kernel Fisher discriminant analysis (KLFDA) [38] is used to reduce dimensionality given that kernelization method can greatly improve Re-ID performance for nonlinearity. Then an SPD matrix is learned on lower-dimensional descriptors based on the limitation that the within-class distance is at least one unit smaller than the minimum interclass distance. This method is proved to have excellent performance compared with other advanced metric learning.

Acknowledgements

First, I would like to express my sincere gratitude to my supervisor, Professor Robert Laganieri, for the continuous support of my master study and research and for his immense knowledge, patience and enlightenment. His guidance helped me throughout my research and writing of this thesis. I would also like to thank Professor Jiying Zhao and Professor Andy Adler for being my thesis examiners.

My sincere thanks also goes to my friends, Di Pang, Dongfeng Gu, Xuelu Wang, Binhao Wang, Muye Jiang and other members of the VIVA lab, for all the fun we have had in the last two years, and for your kind help and encouragement. I would like to thank Lu Sun and Andres Solis Montero, the PhD candidates in the VIVA lab, for giving me suggestions and tips in image segmentation. I would also like to thank Yong Wang, the post doctoral fellow in our lab, for his kind revision suggestions and encouragement.

Last but not the least, I would like to thank my parents and sister for supporting me throughout the study in Canada and the writing of this thesis.

Contents

Declaration of Authorship	i
Abstract	ii
Acknowledgements	iii
1 Introduction	1
1.1 Basic concepts	3
1.2 Challenges	5
1.3 Proposed work	6
1.4 Contributions	8
1.5 Thesis organization	9
2 Related work	10
2.1 Appearance descriptors	10
2.2 Metric learning	14
2.3 Other methods for Re-ID	16
2.4 Some state-of-the-art works	17
2.5 Performance measuring	19
3 Extraction of Hierarchical Gaussian descriptors	22
3.1 Basic color and textural features	22
3.1.1 Color histogram descriptors on different color space	22
3.1.2 Local binary pattern (LBP)	24
3.1.3 Histogram of oriented gradients (HOG)	24
3.2 Influence of background segmentation on different basic descriptors	26
3.3 The hierarchical gaussian descriptor	27
3.3.1 Handling the background	30
3.3.2 Single pixel modelling	30

3.3.3	Riemannian manifold based SPD transformation	32
3.3.4	Integral image for fast region covariance computation	33
3.3.5	Some variants of Hierarchical Gaussian descriptor	35
3.3.6	Dimensionality and superiority analysis of Hierarchical Gaussian descriptor	37
4	Dimension reduction and Mahalanobis distance learning	39
4.1	Kernel local fisher discriminant analysis	39
4.1.1	Fisher discriminant analysis (FDA)	40
4.1.2	Locality preserving projection (LPP)	41
4.1.3	Local fisher discriminant analysis (LFDA)	42
4.1.4	Kernel local fisher discriminant analysis (KLFDA)	43
4.2	Mahalanobis distance	47
4.3	Gradient descent optimization	47
4.4	Metric learning based on sample pairs distance comparison	48
5	Experiment Settings	51
5.1	Datasets and evaluation settings	51
5.2	The influence of mean removal and L_2 normalization	53
5.3	Parameters setting	55
5.4	Results and experiment analysis	57
6	Conclusion	63
6.1	Contributions	63
6.2	Future work	64
6.2.1	Improve Hierarchical Gaussian descriptor	64
6.2.2	Influence of video-based foreground segmentation	64
6.2.3	Computational cost of gradient descent method	64
	Bibliography	65

List of Figures

1.1	Re-ID workflow	1
1.2	A typical single-shot Re-ID workflow	2
1.3	The VIPeR dataset	4
1.4	Samples from the prid_2011 dataset	4
1.5	VIPeR foreground	5
1.6	The workflow of proposed work; the left part is training and the right part is testing	7
1.7	Samples from the prid_2011 dataset	7
2.1	A CMC plot of gallery size is $M = 10$ and probe size is $N = 12$. .	21
3.2	A CMC comparison of color histogram on different color spaces . .	23
3.1	RGB and HSV visual comparison, the first row is RGB and second row is HSV for same views	23
3.3	A comparison of two patches with same entropy but different color distribution	24
3.4	LBP: by thresholding the neighbour pixels the pixels are transformed into a binary number	24
3.5	One LBP example	25
3.6	A demo of HOG feature with a cell size of six pixels	26
3.7	Foreground segmentation of individuals from VIPeR	27
3.8	A CMC comparison of foreground segmentation on HSV histogram descriptor tested on VIPeR	28
3.9	A CMC comparison of foreground segmentation on LBP feature tested on VIPeR	28
3.10	A CMC comparison of foreground segmentation on HOG feature tested on VIPeR	29
3.11	Integral image	34

3.12 A SLIC superpixel segmentation example	36
4.1 Example of five clusters belong to three classes	45
4.2 1-D distribution of dimension reduced data with gaussian kernel	46
4.3 1-D distribution of dimension reduced data with linear kernel	46
4.4 Steepest gradient descent	48
5.1 Pedestrians in prid_450 dataset	51
5.2 Pedestrians in prid_450s dataset	52
5.3 Pedestrians in GRID dataset	53
5.4 Rank 1 scores with respect to α on VIPeR	56
5.5 Rank 5 scores with respect to α on VIPeR	56
5.6 CMC curves on VIPeR comparing different metric learning	58
5.7 CMC curves on CUHK1 comparing different metric learning	59
5.8 CMC curves on prid_2011 comparing different metric learning	60
5.9 CMC curves on prid_450s comparing different metric learning	61
5.10 CMC curves on GRID comparing different metric learning	62

List of Tables

2.1	A CMC example of gallery size is $M = 10$ and probe size is $N = 12$	20
3.1	A rank score comparison of GOG_{fusion} Variants	35
3.2	A performance comparison between Gaussian of Gaussian descriptor and its variants on VIPeR dataset	37
4.1	Optimization algorithm of Mahanalobis distance matrix learning	50
5.1	Testing setting for different datasets	53
5.2	The influence of data preprocessing on VIPeR	54
5.3	The influence of data preprocessing on CUHK1	54
5.4	The influence of data preprocessing on prid_2011	54
5.5	The influence of data preprocessing on prid_450s	54
5.6	The influence of data preprocessing on GRID	55
5.7	Parameters setting	56
5.8	Performance of different metrics on VIPeR	57
5.9	Performance of different metrics on CUHK1	58
5.10	Performance of different metrics on prid_2011	59
5.11	Performance of different metrics on prid_450s	60
5.12	Performance of different metrics on GRID	61

Chapter 1

Introduction

People re-identification (Re-ID) has been an intense research topic in recent years, the main goal of which is to match an image of a given person with other images of the same person. Person Re-ID has great potential in video surveillance, target detection and tracking, and forensic search. However, it is quite challenging since the accuracy is greatly influenced by many factors like occlusion, illumination variation, camera settings and color response. In Re-ID, those images with known labels are called gallery images, and the image used to determine its label is called a probe image. The probe image and gallery images can be from the same or different camera views, so the viewpoint and illumination between the probe and gallery image can be quite different. Also, because of the different color response of different cameras, the images of the same person may look different in different cameras. Occlusions between the camera and target can also bring about much difficulty. In a word, images of the same person may look different while images of different persons may look quite the same.

Given a sequence or video of individuals, there are three steps to match a person. A simple workflow is shown in Figure 1. However, since most Re-ID datasets are already cropped manually or by an automatic detector, most Re-ID work will only focus on robust descriptors design and efficient matching algorithm aimed at those well-cropped images.

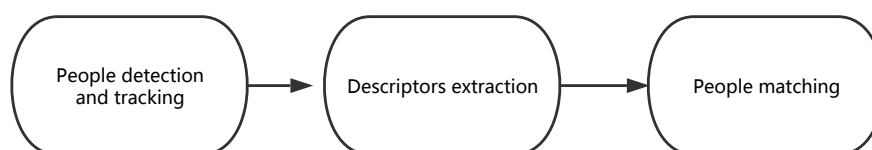


FIGURE 1.1: Re-ID workflow

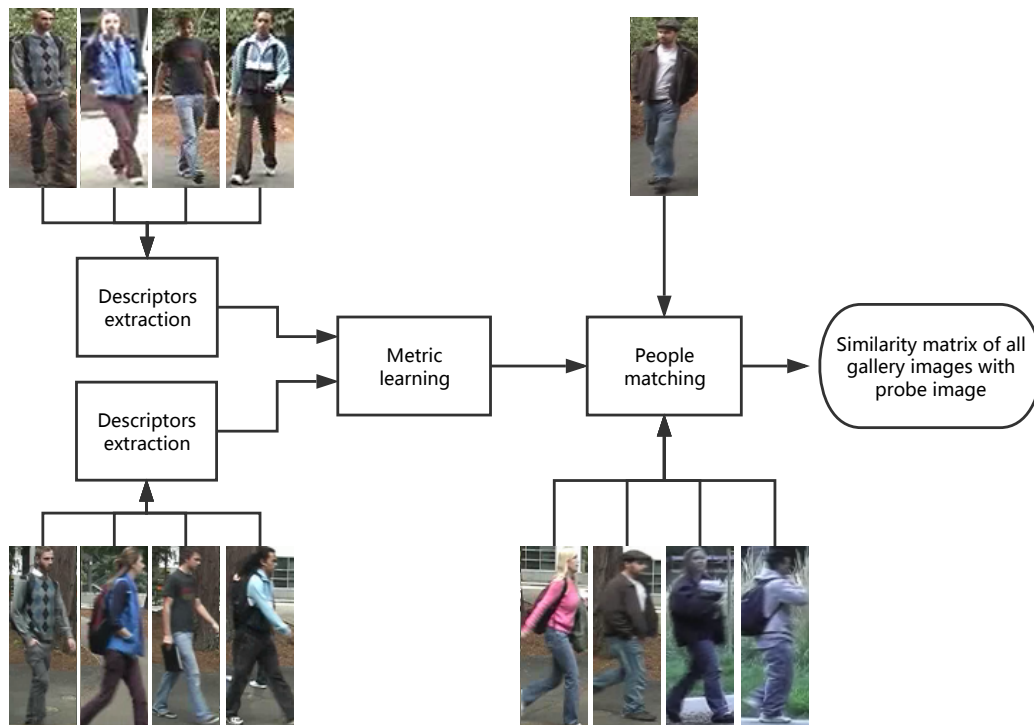


FIGURE 1.2: A typical single-shot Re-ID workflow

The first task in Re-ID is to design a robust descriptor to represent images. The descriptor is supposed to contain the key information for each captured person. Basically, the descriptors are supposed to be robust and discriminative. One straightforward way is to extract the color, textural information of images and then use descriptors to compute the similarity scores. But this method turns out to be not robust to the illumination variation, camera color response difference and camera angle settings. Therefore, many other advanced descriptors take into account the correlation of color, texture and position together to improve performance.

The second is to design the similarity computing methodology (i.e., compare the similarity of two descriptors). Previous methods use Euclidean distance, Bhattacharyya distance and Mahalanobis distance. The Euclidean distance is the easiest to match descriptors like color and texture descriptors, but not the most effective. Many creative metric learning methods have been proposed to compute descriptor similarity. Among them, the Mahalanobis distance based metric is the most popular. The goal in the Mahalanobis distance based metric is to learn a semi-positive definite

(SPD) matrix M , so that M satisfies predefined intraclass and interclass distance limitations.

1.1 Basic concepts

People re-identification can be divided into a few categories according to different conditions. Some general concepts are listed below.

open-set and close-set Re-ID [6] According to the gallery size and how the gallery size evolves, Re-ID can be divided into open-set Re-ID and close-set Re-ID. In close-set Re-ID, no new identities will be added to the gallery set, and gallery size remains the same as time goes by. Furthermore, the probe set will be a subset of the gallery set. That means, the number of unique identities in the gallery set will be equal or greater than the probe set. In open-set Re-ID, the gallery set will evolve as time goes by. Each time a probe image is entered into the recognition system, the system will judge if it has a corresponding match in the gallery set. If the probe image doesn't match any of the gallery images, it will be regarded as a new identity and will be added to the gallery set. The probe set is not necessarily the subset of the gallery set.

Long-term and short-term Re-ID According to the time interval between gallery and probe images, Re-ID can be divided into long-term and short-term Re-ID. In short-term Re-ID the time interval between gallery and probe images is small, say a few minutes or several hours. Conversely, long-term Re-ID refers to the case that the time interval between gallery and probe images is a few days or even longer. One difference to consider regarding the long time interval between gallery and probe images is the variation of individuals' clothing and appearance. If the gallery images were shot a few days prior, the same individual may have changed his suit or taken off his bag, causing the appearance to change. In this case, it will be much more difficult to recognize the same identity in long-term Re-ID. In most cases, we use short-term Re-ID, which guarantees that the appearance of the same person will remain the same and we will only need to consider the differences brought by other factors such as viewpoint variation and occlusions.

Single-shot and multi-shot Re-ID According to the size of the sample set for each person, Re-ID can be divided into single-shot and multi-shot approaches. In

single-shot Re-ID, only one image is provided for a person in a camera view. Single-shot Re-ID is challenging because only limited information can be extracted. One example is the VIPeR dataset shown in Figure 1.3. For each person in this dataset only one image is provided in each camera view, and the viewpoint of each view is different. In multi-shot Re-ID, a sequence of images is provided for a person in a camera view. Compared with single-shot Re-ID, more information, like temporal-spatial information, can be extracted from the sample set. One example of multi-shot dataset is the prid_2011 dataset 1.4, which provides a long sequence of images for each person in a single camera view.



FIGURE 1.3: The VIPeR dataset

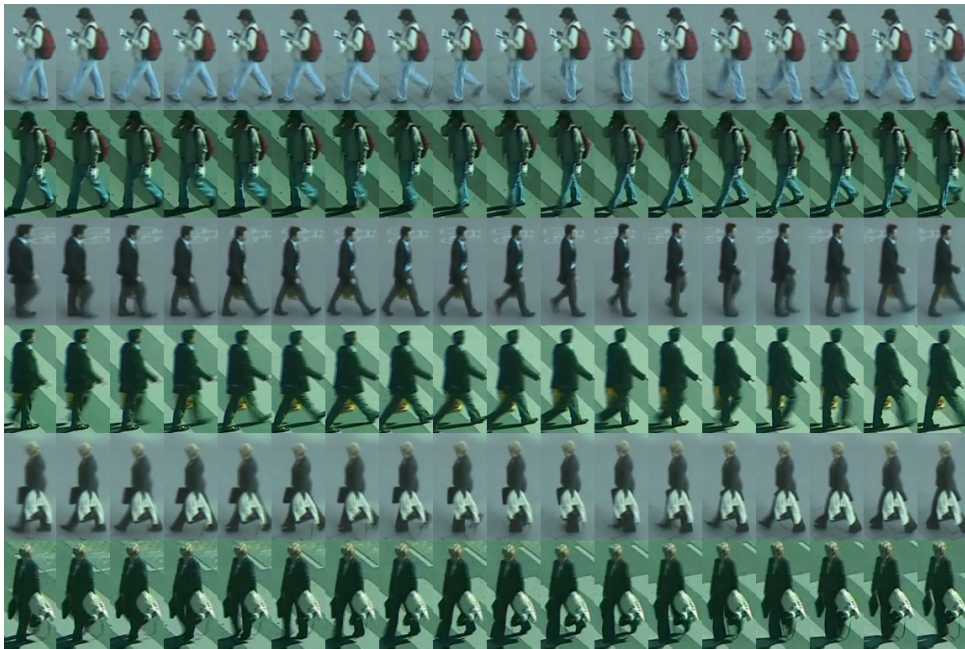


FIGURE 1.4: Samples from the prid_2011 dataset

1.2 Challenges

Detection, tracking and dataset labelling for supervised learning Though classical person re-identification focuses on robust descriptors design and matching algorithms, in real-time application, the detection and tracking has to be operated on video frames to get bounding boxes of individuals. A good detection and tracking algorithm is necessary for Re-ID. Furthermore, training the matching algorithm is a supervised process, thus we have to know the labels for those training data.

Descriptors design Good descriptors should be robust to the variation of people's postures, outer environment changes and camera settings. Though there have been many kinds of descriptors based on different property like color and texture, it is hard to judge which property is universally useful for different camera settings. In fact, the robustness, reliability and feasibility depends on different camera settings and viewing conditions. What's more, the pedestrian background may add many errors to descriptors, so it is important to quantify the impact of a noisy background. Many works have tried to use a segmented foreground of pedestrians, so it is important to design segmentation algorithms. The automatic foreground segmentation for a single frame is difficult because there isn't much information available compared with video background segmentation. Take the VIPeR dataset as an example. There is only one frame for each view of a certain person, thus the segmented foreground masks are imperfect, and chances are high that important body parts will be lost. A segmented foreground provided by [15] is shown in Figure 1.5.



FIGURE 1.5: VIPeR foreground

Efficient matching algorithm design When designing machine learning algorithms to match persons, there are many limitations. One of them is the small sample

size problem [47]. The extracted descriptors usually has a high dimension d but only a small number of samples $n(n \ll d)$, underfitting may appear because of insufficient data samples with high dimension. It is also necessary to take into consideration intra and inter distance of samples. The intraclass distance means the distance of two samples with the same class label, while interclass distance is the distance of samples with different class labels.

Feasibility, complexity and scalability When applying those descriptors and matching algorithms, we have to consider real-time performance. The Re-ID datasets usually have small sample sizes but in a surveillance network, much more pedestrians in different cameras can be presented simultaneously. A system like this has plenty of individuals to re-identify, which requires that the processing time for a single probe should be short for low latency. Because the gallery in this system evolves, it is crucial to design an algorithm that can determine if a person appearing in a current camera is new or has previously appeared in the gallery.

1.3 Proposed work

In many previous works, the kernel local Fisher discriminant analysis is used as a subspace learning method, and Euclidean distance is usually used in the subspace to measure similarity. In this thesis, the KLFDA [38] method is used as a dimension-reducing method to project high-dimensional descriptors to a lower-dimensional space. Compared with other dimension reduction methods, KLFDA is a supervised method that takes into consideration of intraclass and interclass information; therefore, less information is lost after dimension reduction. A Mahalanobis distance based matrix M is learned based on the limitation that the distance of people from the same class should be at least one unit smaller than the distance of people from different classes. A target function that penalizes large intraclass distance and small interclass distance is created. When the target function converges by iterative computation, the matrix M is thought to be optimal. It turns out that this metric learning exhibits good performance when compared with other metric learning methods. A workflow of proposed work is in Figure 1.6.

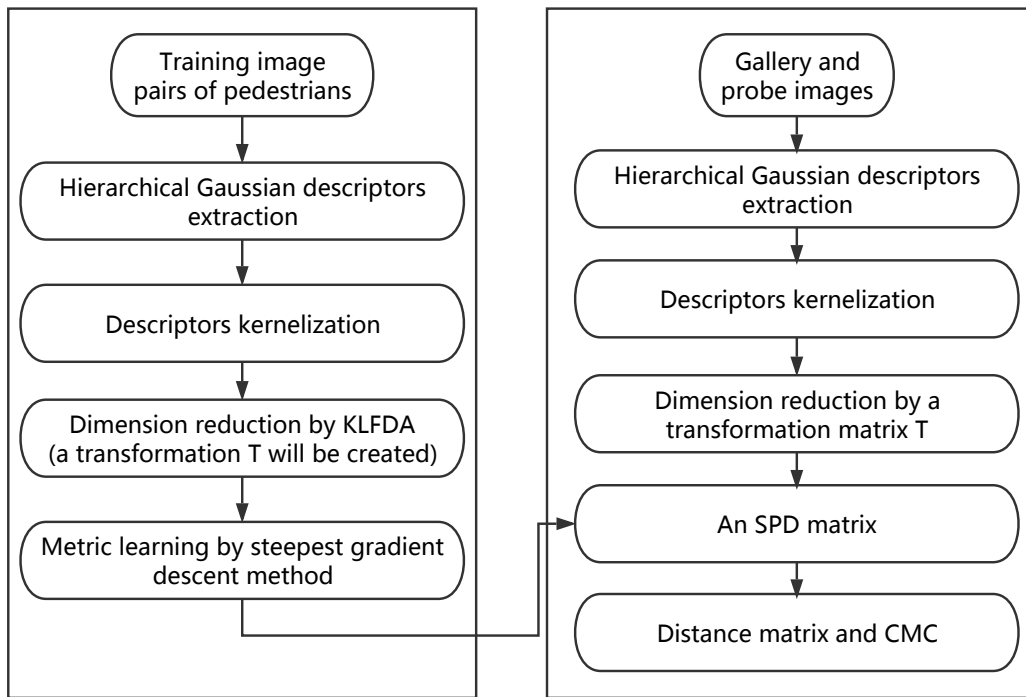


FIGURE 1.6: The workflow of proposed work; the left part is training and the right part is testing



FIGURE 1.7: Samples from the prid_2011 dataset

1.4 Contributions

In this paper, we have three contributions. The first is that we combined the KLFDA with distance comparison learning. Instead of learning the subspace with KLFDA and computing Euclidean distance in a lower-dimensional space, a Mahalanobis distance based matrix is learned under the limitation that the intraclass distance is at least one unit smaller than interclass distance. Compared with those advanced metrics including cross view quadratic analysis (XQDA) [22] and the null Foley-Sammon transfer (NFST), this proposed metric learning proves to have excellent performance on the VIPeR, CUHK1, prid_2011, prid_450s and GRID dataset.

Another contribution of this thesis is the study of influence of background subtraction on different descriptors. We found that background subtraction can improve the performance of some descriptors, like the histogram of HSV colorspace, but can also decrease the performance of other descriptors, like local binary pattern (LBP) and histogram of gradient (HOG). This comparison is shown in Chapter 3. The reason for this is that imperfect background segmentation brings in textural interference. If descriptors are color based and don't handle texture information, like HSV histogram descriptor, background segmentation can greatly improve the performance. However, if the descriptor extracts texture information, background segmentation will decrease its performance since the imperfect segmentation will mask out many parts of the foreground area, which will cause important textural information variation. Because segmentation algorithms will cause different influence on various features, in this thesis, a weighted map of images is used instead of using the background segmentation.

For the last contribution, some variants of hierarchical Gaussian descriptor have been tested. In one variant, LBP was used in the basic pixel feature. In another variant, superpixel segmentation was also applied to combine with hierarchical Gaussian descriptor, which implied that overlapping patch sampling is important in a hierarchical Gaussian descriptor. At last, the Gaussian mixture model (GMM) was also tested but it displayed the worst performance.

1.5 Thesis organization

In this thesis, Chapter 2 will give a brief introduction of previous work. Chapter 3 will explain the implementation of the hierarchical Gaussian descriptors used in this thesis. The performance of some variants of the hierarchical Gaussian descriptor is studied in this chapter. In Chapter 4, a detailed introduction of the kernel local Fisher discriminant analysis will be presented, and a detailed explanation of the metric learning on the lower-dimensional space based on relative distance limitation learning will also be provided. In Chapter 5, the used datasets and parameters and other experiment settings will be explained, and a detailed analysis of results will be presented there. Finally, the conclusion is given in Chapter 6.

Chapter 2

Related work

Previous works mainly focus on finding more discriminative descriptors and better metric learning. It is known that color and texture are the most important information in Re-ID. Researchers found that features based on a single attribute (like color or texture) are not robust to various datasets. Instead, combinations of different features are exploited to improve the performance. Most descriptors capture the local or global statistic color and texture information to characterize individuals. A brief introduction of those descriptors and metrics is given in this chapter.

2.1 Appearance descriptors

In most descriptors, the input image will be divided into a few subregions to model the complex human kinematics. Features of those subregions are extracted respectively and concatenated directly or characterized by their statistic properties. According to how those subregions are divided, there can be three kinds of models, fixed part-based models, adaptive models and learned part models [36].

In fixed part models, the size of body parts is predefined. One example is in [11, 50, 34], where a silhouette is divided into a fixed number of horizontal stripes, which mainly include the head, torso and legs. In [19], the input images are divided into three horizontal stripes, and the width of each stripe is respectively 16%, 29% and 55%. The fixed models predefine parameters including numbers of stripes and the stripe height.

In the adaptive part models, the size of each body parts may vary to fit predefined body part models. Take [15] for an instance; the silhouette of each person is divided into three parts horizontally, which include the head, torso and legs respectively. But the height of each stripe is different for various silhouettes, and it is computed

according to symmetry and asymmetry with two operators $C(i, \delta)$ and $S(i, \delta)$, where

$$\begin{aligned} C(i, \sigma) &= \sum_{B_{[i-\delta, i+\delta]}} d^2(p_i - \hat{p}_i) \\ S(i, \sigma) &= \sum \frac{1}{W\delta} |A(B_{[i, i-\delta]}) - A(B_{[i, i+\delta]})|. \end{aligned} \quad (2.1)$$

Here the $C(i, \delta)$ is called the *chromatic bilateral operator*, and it computes the Euclidean distance of two HSV blobs located symmetrically with respect to horizontal axis $y = i$, where $B_{[i-\delta, i+\delta]}$ is the blob with a height 2δ . $S(i, \delta)$ is called the *spatial covering operator* and it computes the difference of two foreground areas. Then the axis between the torso and legs is computed as follow:

$$i_{TL} = \arg \min_i (1 - C(i, \delta) + S(i, \delta)), \quad (2.2)$$

and the axis between the head and torso is computed with the following equation:

$$i_{HT} = \arg \min_i (-S(i, \delta)). \quad (2.3)$$

The axis dividing the left and right torso is

$$j_{LR} = \arg \min_j (C(j, \delta) + S(j, \delta)). \quad (2.4)$$

This method has good performance. But one shortcoming of this model is that an imperfect background segmentation causes noise and introduces errors regarding the position of axes.

Another part-based adaptive spatial-temporal model used in [7] characterizes a person's appearance using color and facial features. Few works exploit human face features. In this work, human face detection based on low resolution cues selects useful face images to build face models. Color features capture representative color as well as the color distribution to build a color model. This model handles multi-shot re-identification, and it also characterizes the color distribution variation of many consecutive frames. Besides, the facial features of this model is conditional. That is, in the absence of good face images, this model is only based on color features.

Some methods based on learned part models have been proposed. Part model detectors (statistic classifiers) are trained with manually labelled human body parts images, exploiting features related to edges contained in the images. A pictorial

structure (PS) is proposed in [16]. The PS model of a non-rigid body is a collection of part models with deformable configurations and connections with certain parts. The appearance of each part is separately modelled, and deformable configurations are implemented with spring-like connections. This model can quantitatively describe visual appearance and model the non-rigid body. In [2], the pictorial structure body model is made up of N parts and N corresponding part detectors.

Another example of the learned part model is in [8, 7]. The overall human body model consists of several part models; each model is made up of a spatial model and a part filter. For each part, the spatial model defines allowed arrangements of this part with respect to the bounding box. To train each model, the latent support vector machine (LSVM) is used, and four body parts are detected, namely the head, left and right torso and upper legs. Compared with other models, this model exploits a sequence of frames of an individual and thus captures appearance characteristics as well as the appearance variation over time.

According to the way to extract features for a model (a whole model or a part-based model), the feature can be implemented with different methods. The features can be divided into two categories: global and local features. Global features refer to features extracted from a whole image or region, and the size of the descriptor is usually fixed. In order to extract local feature of a specified image or region, we first divide the whole image into many equal blocks and compute the feature of each block. Both descriptors may deal with color, texture and shape. The color information is exploited most by extracting the color histogram within different color spaces. Descriptors based on texture, such as the scale-invariant feature transform (SIFT), speeded up robust features (SURF) and LBP, are also widely combined to improve performance.

Global color histogram is a frequently used global feature. For an three-channel image, like an RGB image, each channel is quantized into B bins separately. The final histogram could be a multi-dimensional or one-dimensional histogram. For instance, if $B = 8$, for a multi-dimensional histogram, there will be $8 \times 8 \times 8 = 512$ bins. But if we concatenate the three-dimensional bins together, the dimension can be reduced to $8 + 8 + 8 = 24$ bins while the performance of this reduced descriptor doesn't decrease. This method can be applied on other color spaces like HSV and Lab, etc.

Local color histogram usually splits the specified model or region into many

equal-sized blocks and computes the global feature of each block. The feature can be based on color, texture and interest points. SIFT [24] is a kind of local feature based on the interest points. The salient interest points (identifiable over rotating and scaling) are selected by the interest operator. This algorithm detects key points by computing the difference of Gaussian (*DoG*) images of different scales σ with the equation

$$D(x, y, \sigma) = (G(x, y, k_1\sigma) - G(x, y, k_2\sigma)) * I(x, y). \quad (2.5)$$

Here $G(x, y, k_1\sigma)$ is the Gaussian function with deviation $k_1\sigma$, $I(x, y)$ is the image. The *DoG* images are compared to find their extrema as key points. With key points localization and other processing, descriptors describing key points are created as SIFT descriptors.

The maximally stable color region (MSCR) is used in [15]. The MSCR derives from the maximally stable extreme region (MSER) and detects the region with a stable color cluster. It uses an agglomerative clustering algorithm to compute color clusters, and by looking at the successive time steps of the algorithm, the extension of color is implemented. The detected color region is described with a nine-dimensional vector containing the area, averaging color, centroid and second moment matrix. With this vector, the color region detected makes it easy to do scale and affine transforms.

Recurrent highly structured patches (RHSP) are also used in [15]. This feature captures patches with highly recurrent color and texture characteristics from extracted silhouette pixels. This feature is extracted from the following steps. First, random and probably overlapping small image patches are extracted from silhouette pixels. Then, to capture those patches with informative texture, the entropy of each patch (the sum of the three channels' entropy) is computed, and we discard those patches with entropy smaller than a specified threshold. In the next step, some transforms are performed on the remaining patches to select those remaining invariant to the transforms. Subsequently, the recurrence of each patch is evaluated with the local normalized cross correlation (LNCC) function. This evaluation is only performed on a small region containing the patch instead of the whole image. Then the patches with high recurrence are clustered to avoid patches with similar content. Finally, the Gaussian cluster is applied to maintain the patch nearest to the centroid of each cluster for each cluster.

In [3], a new 3D model called SARC3D is used to characterize the individual. Compared with those 2D models, this model combines the texture and color information with their location information together to get a 3D model. This model starts with an approximate body model with a single shape parameter. By precise 3D mapping, this parameter can be learned and trained with even quite few images (even one image is feasible). This model's construction is driven by the frontal, top and side views extracted from various videos, and for each view, the silhouette of a person is extracted to construct the 3D graphical model. The final body model is sampled to get a set of vertices from the previously learned graphic body model. Compared with other models, this model has a robust performance when dealing with partial occlusion, people posture variation and viewpoint variations since the model is based on silhouettes from three viewpoints.

Descriptors combining color and texture are most often used in re-identification. In [5], a signature called asymmetry-based histogram plus epitome (AHPE) was proposed. This work starts with a selection of images to reduce image redundancy, which is caused by correlated consecutive sequences. This descriptor combines global and local statistical descriptors of human appearance, focusing on overall chromatic content via histogram and on the recurrent local patches via epitome analysis. Similar to the SDALF descriptor [15], the HPE descriptor consists of three components, the chromatic color histogram, the generic epitome and local epitome. The chromatic color histogram is extracted in the HSV color space, which turns to be robust to illumination changes. Here, color histogram is encoded into a 36-dimensional feature space $[H = 16, S = 16, V = 4]$. The authors define epitome by extracting generic and local epitome.

2.2 Metric learning

The second step of Re-ID is to design the metric learning to match descriptors, i.e., the way to compare how similar two descriptors are. Many different metric learning methods have been proposed [20, 33, 27, 46, 50, 41, 38, 42, 44, 47, 13] to get smaller intraclass distance and larger interclass distance. Generally, for two $d \times 1$ -dimensional input vectors $\mathbf{x}_1, \mathbf{x}_2$, any symmetric positive semi-definite matrix \mathbf{M} defines a pseudo-metric with the form of $D = (\mathbf{x}_1 - \mathbf{x}_2)^T \mathbf{M} (\mathbf{x}_1 - \mathbf{x}_2)$. Many widely

used distance metrics exploit this rule. Previous methods includes the Euclidean distance, Bhattacharyya distance and Mahalanobis distance methods. The Euclidean distance, which is the most common used distance, is a special case of Mahalanobis distance when the M is an identity matrix. One example of metric learning is the probabilistic relative distance comparison model proposed in [50]. This model exploits all the possible positive person pairs and negative pairs so that for each person, the between-class distance is larger than the within-class distance. Compared with the other distance learning models proposed, this model solves the matrix M by an iterative optimization algorithm. Suppose z is an image descriptor of a person; the task is to identify another image descriptor z' of the same person from z'' of a different person by using a distance model $f(\cdot)$, so that $f(z, z') < f(z, z'')$. The authors convert the distance learning problem to a probability comparison problem by measuring the probability of the distance between a relevant pair of images being smaller than that of a related irrelevant pair as

$$P(f(z, z') < f(z, z'')) = (1 + \exp^{(f(z-z')-f(z-z''))})^{-1}. \quad (2.6)$$

Here the author assumes that probability of $f(z, z')$ and $f(z, z'')$ is independent, therefore, using the maximal likelihood principal, the optimal function can be learned as

$$\begin{aligned} f &= \arg \min_f r(f, O) \\ r(f, O) &= -\log\left(\prod_{O_i} P(f(\mathbf{x}_i^p) < f(\mathbf{x}_i^n))\right). \end{aligned} \quad (2.7)$$

$O = \{O_i = (\mathbf{x}_i^p, \mathbf{x}_i^n)\}$, \mathbf{x}_i^p , \mathbf{x}_i^n are intraclass and interclass vector differences respectively, and \mathbf{x}_i^p , \mathbf{x}_i^n are defined as

$$\begin{aligned} \mathbf{x}_i^p &= |z - z'| \\ \mathbf{x}_i^n &= |z - z''|. \end{aligned} \quad (2.8)$$

The distance function $f(\cdot)$ here is parameterized as the Mahalanobis distance function $f = \mathbf{x}^T M \mathbf{x}$, $M \geq 0$. Here M is a semi-positive definite (SPD) matrix. In this way, the distance function learning problem is transformed to a matrix optimization problem. The author used an iteration algorithm to compute matrix M . One shortcoming for this algorithm is that it is computationally expensive because for each person it compares all the possible negative pair distances with corresponding

positive pair distance.

Single-shot image-based person representation suffers from the small sample size problem. This is why multi-shot Re-ID has been proposed. Since there are a sequence of images for each individual, there are much more cues to exploit. In [46], the author simplified computing of Mahalanobis matrix by applying the new limitations on datasets. The author finds that when using video based person representation the difference of inter-class may be more obscure than that of still image based representation. Therefore, the author proposed the top-push distance learning. For a person video sequence, the maximal intraclass distance should be 1 unit smaller than the minimal distance of interclass distance. Another constraint introduced is that sum of all intra-class distance should be as small as possible, so the final target function is summarized as

$$f(D) = (1 - \alpha) \sum_{\mathbf{x}_i, \mathbf{x}_j, y_i=y_j} D(\mathbf{x}_i, \mathbf{x}_j) + \alpha \sum_{\mathbf{x}_i, \mathbf{x}_j, y_i=y_j} \max\{D(\mathbf{x}_i, \mathbf{x}_j) - \min_{y_i \neq y_k} D(\mathbf{x}_i, \mathbf{x}_k) + \rho, 0\} \quad (2.9)$$

here $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k$ are feature descriptors and y_i, y_j, y_k are class labels.

Some previous works exploit Fisher discriminant analysis and convert the problem to matrix eigenvalue decomposition problem. In [33] the local fisher discriminant analysis (LFDA) is used, and in [44] the kernel version several linear metrics are proposed and it proves that kernelization improves Re-ID performance. In [22] the log ratio of Bayesian based estimation is proposed to have advanced performance. In [47] the authors propose to make sample points of the same class collapse to the same point in the null space while different points are project to different points. In [45] the semantic representation of color names are exploited. In this thesis a SPD matrix M is learned to meet certain intraclass and interclass distance limitations.

2.3 Other methods for Re-ID

Beside descriptors and metrics mentioned above, there are some other methods for Re-ID. Convolutional neural network (CNN) have been exploited in Re-ID. One advantage of neural network Re-ID is the preprocessing of images can be skipped (we can also say the preprocessing is included in convolutional layers). The input of this

structure can be straight-forward grey images or color images. Traditional neural network has too many weights to train. Convolutional neural network can avoid this problem while retaining high performance. Compared with classical neural network architecture, the convolutional neural network exploits receptive field, weights sharing and pooling to reduce weights number and thus decreases computational cost. When dealing with multi-shots and video based re-identification neural network is proven to have better performance. In [26] the author proposes a recurrent neural network layer and temporal pooling to combine all time-steps data to generate a feature vector of the video sequence. In [10] the author proposes a multi-channel layers based neural network to jointly learn both local body parts and whole body information from input person images. In [43] a convolutional neural network learning deep feature representations from multiple domains is proposed, and this work also proposes a domain guided dropout algorithm to dropout CNN weights when learning from different datasets.

There are many other works based on convolutional neural networks. However, person re-identification may be one of the areas where CNN's performance may be poorer than regular machine learning methods for the small sample size problem (SSS). In most datasets, the sample size of each pedestrian is quite small. Especially in single shot Re-ID only one frame is provided in each view for each person, this is why Re-ID more often rely on classical machine learning.

2.4 Some state-of-the-art works

Recently many works have been proposed and improved Re-ID performance by much margin. In this section those advanced descriptors and metrics are introduced.

Cross view quadratic discriminant analysis (XQDA) is proposed in [22]. Let define the sample difference $\Delta = \mathbf{x}_i - \mathbf{x}_j$, where \mathbf{x}_i and \mathbf{x}_j are two feature vectors. Δ is called intrapersonal difference when their labels satisfy $y_i = y_j$ and extrapersonal difference when $y_i \neq y_j$. The intrapersonal and interpersonal variation can be defined as Ω_I and Ω_E . The authors convert Re-ID problem to distinguishing Ω_I from Ω_E . In [28] each one of intrapersonal and interpersonal class is modelled with a multivariate Gaussian distribution, and in [28] it has been proved that both Ω_I and Ω_E have zero mean. Under the zero-mean distribution, the probability of observing Δ in Ω_I and

the probability of observing Δ in Ω_E can be denoted as

$$\begin{aligned} P(\Delta|\Omega_I) &= \frac{1}{(2\pi)^{d/2}|\Sigma_I|^{1/2}} \exp^{-\frac{1}{2}\Delta^T\Sigma_I^{-1}\Delta} \\ P(\Delta|\Omega_E) &= \frac{1}{(2\pi)^{d/2}|\Sigma_E|^{1/2}} \exp^{-\frac{1}{2}\Delta^T\Sigma_E^{-1}\Delta} \end{aligned} \quad (2.10)$$

where Σ_I and Σ_E are the covariance matrix of Ω_I and Ω_E , then the probability ratio between the interpersonal pairs and intrapersonal pairs can be denoted as

$$\begin{aligned} r(\Delta) &= \frac{P(\Delta|\Omega_E)}{P(\Delta|\Omega_I)} = \frac{\frac{1}{(2\pi)^{d/2}|\Sigma_E|^{1/2}} \exp^{-\frac{1}{2}\Delta^T\Sigma_E^{-1}\Delta}}{\frac{1}{(2\pi)^{d/2}|\Sigma_I|^{1/2}} \exp^{-\frac{1}{2}\Delta^T\Sigma_I^{-1}\Delta}} \\ r(\Delta) &= C \exp^{-\frac{1}{2}\Delta^T(\Sigma_E^{-1}-\Sigma_I^{-1})\Delta} \end{aligned} \quad (2.11)$$

C is the constant term, by taking log and deserting the constant term, we have

$$r(\Delta) = \Delta^T(\Sigma_I^{-1} - \Sigma_E^{-1})\Delta = (\mathbf{x}_i - \mathbf{x}_j)^T(\Sigma_I^{-1} - \Sigma_E^{-1})(\mathbf{x}_i - \mathbf{x}_j) \quad (2.12)$$

In [22] a subspace W is learned so that

$$r(\Delta) = (\mathbf{x}_i - \mathbf{x}_j)^T W(\Sigma_I'^{-1} - \Sigma_E'^{-1})W^T(\mathbf{x}_i - \mathbf{x}_j) \quad (2.13)$$

and $\Sigma_I' = W^T\Sigma_IW$, $\Sigma_E' = W^T\Sigma_EW$. Therefore, a subspace $M(W) = W(\Sigma_I'^{-1} - \Sigma_E'^{-1})W^T$ is learned in this work.

Null Foley-Sammon transform In [47] a null space is proposed so that with this space the intraclass points collapse to a same point in the null space while inter-class points are projected to different points. Given the within class scatter \mathbf{S}^w and between class scatter \mathbf{S}^b , an optimal projection matrix \mathbf{W} is computed so that

$$\begin{aligned} \mathbf{w}_i^T \mathbf{S}^w \mathbf{w}_i &= 0 \\ \mathbf{w}_i^T \mathbf{S}^b \mathbf{w}_i &> 0 \end{aligned} \quad (2.14)$$

here \mathbf{w}_i is the i_{th} column in \mathbf{W} .

It can be noticed that like many other metric learnings, XQDA and NFST evolve to matrix decomposition and eigenvalue selection problem. In this paper, those two metrics are used to compare with proposed metric. In [25] it has been shown GOG +

XQDA outperforms many other combinations including MetricEnsemble [32], SC-NCD [45], SemanticMethod [37], etc. In [47] it has been shown LOMO + NFST outperforms metrics including LMNN [41], KCCA [42], ITML [13], KLFDA[38], MFA [44], KISSME [20], SimilarityLearning [9], SCNCD [45], Mid-level Filters [49] and Improved Deep [14]. Based on the result that XQDA and NFST outperform other metrics, XQDA and NFST are used in this thesis to compare with our proposed metric learning.

2.5 Performance measuring

There are a few measures of Re-ID such as cumulative matching curve(CMC), receiver operating characteristic curve (ROC) and synthetic recognition rate (SRR). Specifically, CMC is used as a 1:m re-identification system and ROC is used for 1:1 re-identification system. SRR curve indicates the probability that any of given fixed number of matches is correct.

For the open-set Re-ID problem, the ROC [7, 8] curve is adopted. ROC represents the relationship between Re-ID accuracy vs false accept rate (FAR). In ROC the x-axis is FAR and the y-axis is the accuracy. Re-ID accuracy and FAR are defined by equations

$$\begin{aligned} Accuracy &= \frac{TP_s}{N_p} \\ FAR &= \frac{MM_s + FP_s}{N_p} \end{aligned} \quad (2.15)$$

the true positives (TPs) is the number of probe IDs that are correctly accepted. N_p is the number of all positive probes. The FAR is expressed by the mismatches (MM), false positives (FPs) and N_p . MMs is those probe IDs that are incorrectly matched to the galley when in fact those probe IDs exist in the gallery. The FPs is number of those probe IDs incorrectly matched to the gallery when they don't exist in the gallery actually.

For the closed set Re-ID problem, the mostly used is the cumulative matching curve (CMC). Closed set Re-ID assumes that the probes are contained in gallery set and the performance is measured by ranking the gallery with respect to each probe. The CMC curve describes the probability of right match given a list of computed similarity score, and the first ranked ID is regarded as the matched individual.

The proportion of uncertainty removed (PUR) is proposed in [33]. PUR measure the entropy difference between before ranking and after ranking. For a certain probe before ranking all those samples in gallery set has equal probability and the entropy is $\log(S)$, S is the number of observations in gallery set. After ranking the entropy is $\sum M(r) \log(M(r))$, $M(r)$ is the probability that the right match is within the top r ranked observations. With normalization the PUR value is computed by equation

$$PUR = \frac{\log(S) - \sum M(r) \log(M(r))}{\log(S)} \quad (2.16)$$

In this thesis, the CMC curve is used to measure Re-ID performance because most datasets are tested under closed set Re-ID settings. In this metric the $CMC(k)$ stands for the probability that the right match is within the top k matches. Suppose a set of gallery $G = \{G_1, G_2, \dots, G_m\}$ and a set of probe $P = \{P_1, P_2, \dots, P_n\}$, for each identity P_i there should be a right match in the gallery set. However, there could be identities that appear in gallery set but not in probe set. A $m \times n$ similarity matrix can be computed. Then for each probe identity P_i , a sorted list of gallery identities can be list as $S(P_i) = \{G_{(1)}, G_{(2)}, \dots, G_{(m)}\}$ so that their similarity with P_i descends. Suppose the right match of P_i is at the position k of $S(P_i)$, $k \leq m$, then G_i has a rank value of k . Therefore, the CMC can be calculated as

$$CMC(k) = \frac{1}{n}(\#k_l \leq k) \quad (2.17)$$

where k_l is the list of rank values of $P = \{P_1, P_2, \dots, P_n\}$, and $\#k_l \leq k$ means the number of rank values that is smaller than k . Therefore, CMC curve always ascends and converges at 1. A perfect CMC curve is supposed to have a high rank 1 score and approaches 1 as fast as possible.

An example of CMC computing is as follow. Suppose the gallery size is $M = 10$ and probe size is $N = 15$. By computing the similarity and ranking them we have a rank score frequency table as

TABLE 2.1: A CMC example of gallery size is $M = 10$ and probe size is $N = 12$

RankScore	1	2	3	4	5	6	7	8	9	10	11	12
Frequency	5	3	2	1	1	0	0	0	0	0	0	0
CumulativeSum	5	8	10	11	12	12	12	12	12	12	12	12

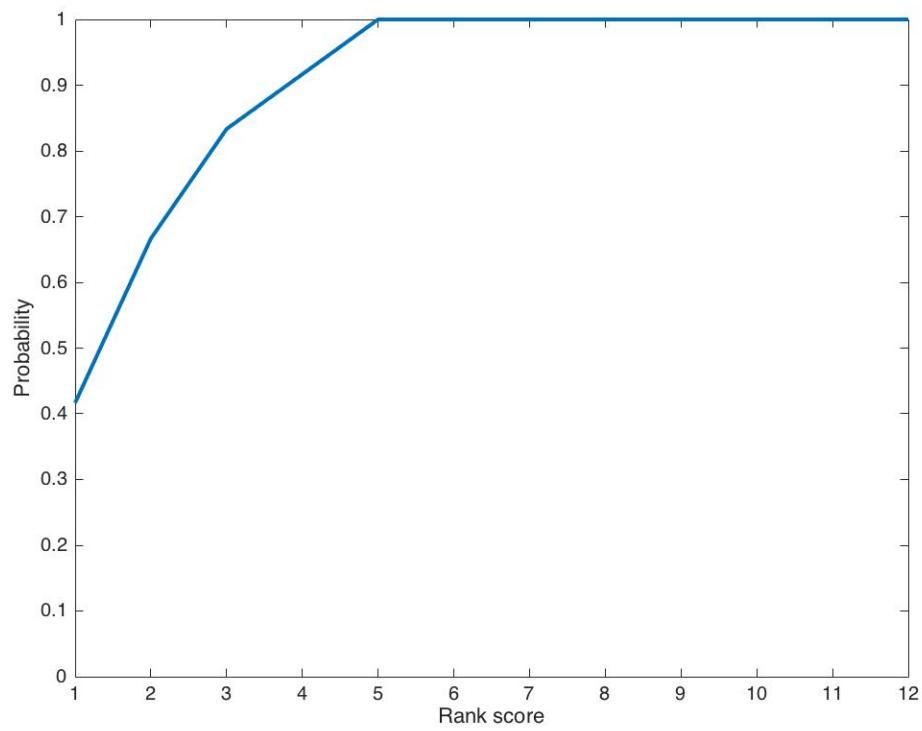


FIGURE 2.1: A CMC plot of gallery size is $M = 10$ and probe size is $N = 12$

Chapter 3

Extraction of Hierarchical Gaussian descriptors

In person re-identification, it is very important to choose robust descriptor to represent person. A good descriptor should be robust to variations of illumination, viewpoint, and camera color response. Most descriptors try to capture the color and texture information. In this chapter, we will first introduce some basic descriptors and compare their performance on VIPeR dataset, then a detailed introduction of a hierarchical descriptor will be presented in section 3.3 of this Chapter.

3.1 Basic color and textural features

3.1.1 Color histogram descriptors on different color space

Histogram descriptor extracts color statistics information of input images. A popular histogram extracting method is to divide input image into a few horizontal stripes and extract color histogram of each stripe, then they are concatenated to produce a histogram descriptor of the whole image. Color space selection has much influence on descriptor performance. HSV color space is commonly used in computer vision and image processing for target detection and tracking. The HSV descriptor has better performance than RGB histogram descriptor since HSV color separates image intensity from color information. Thus HSV color space is more robust to illumination variation. An unsupervised CMC performance comparison among different color spaces on VIPeR dataset is given in Figure 3.2. In this comparison camera A views are used as probe set and camera B views are used for gallery set. We can find

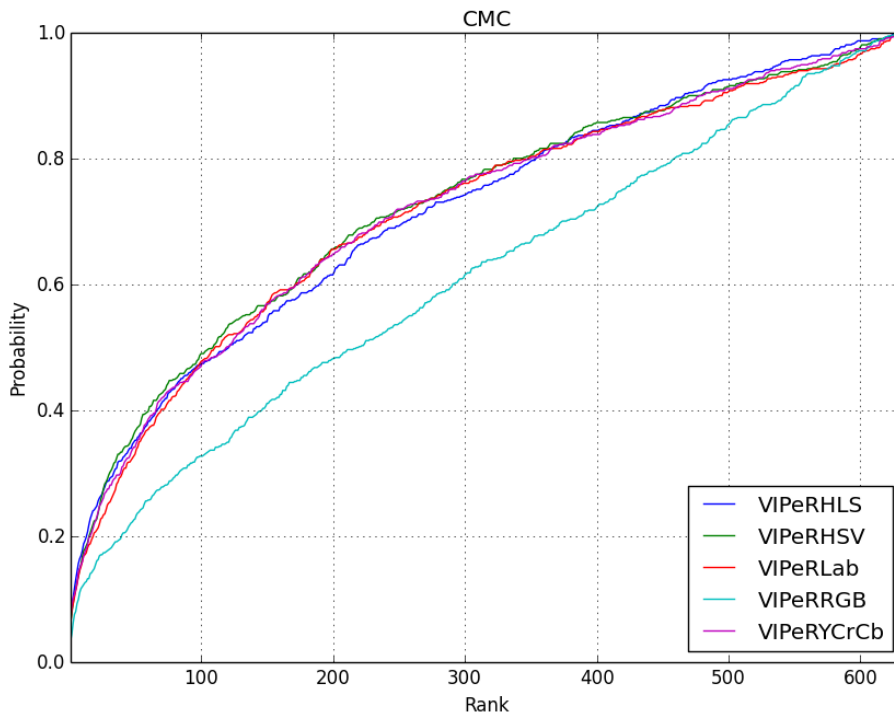


FIGURE 3.2: A CMC comparison of color histogram on different color spaces

that those color spaces separating intensity information outperform RGB color space by a large margin.



FIGURE 3.1: RGB and HSV visual comparison, the first row is RGB and second row is HSV for same views

Shortcoming of histogram based descriptor The performance of histogram descriptors suffers from ignoring the spatial information. Since it doesn't consider the

relative distribution of color patches. Images with same kind color patches but different distribution may have the same histogram descriptor. One example is shown in Figure 3.3.

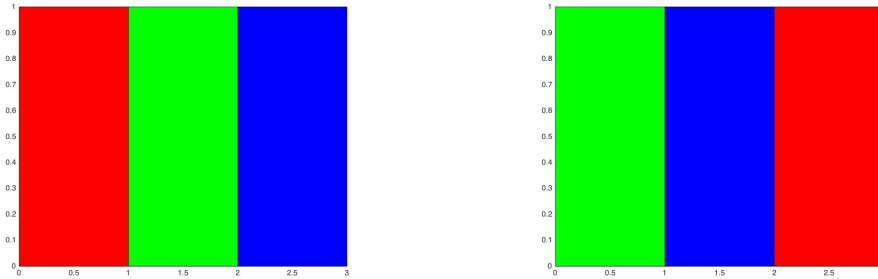


FIGURE 3.3: A comparison of two patches with same entropy but different color distribution

3.1.2 Local binary pattern (LBP)

Local binary pattern [30, 29] extracts the texture information with efficient computing and has been used on people detection and recognitions. Figure 3.5 is an example of LBP. by thresholding neighbour pixel of center pixel (shown in Figure 3.4), the pixels are transformed into a binary integer. There are many extended LBPs like tLBP [39], VLBP [48], OCLBP [4]. Besides, LBP is well known for its robustness to monotonic illumination variation.

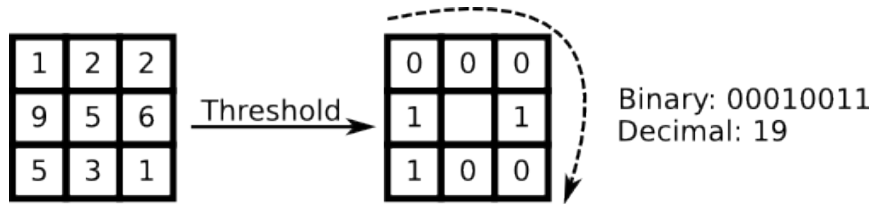


FIGURE 3.4: LBP: by thresholding the neighbour pixels the pixels are transformed into a binary number

3.1.3 Histogram of oriented gradients (HOG)

The HOG [12] descriptor also extracts textural information of images by gradient computing. A brief introduction about its gradient computation is presented here, more details can be found in [12]. HOG feature computes the gradient of input intensity image $I(x, y)$ by equations

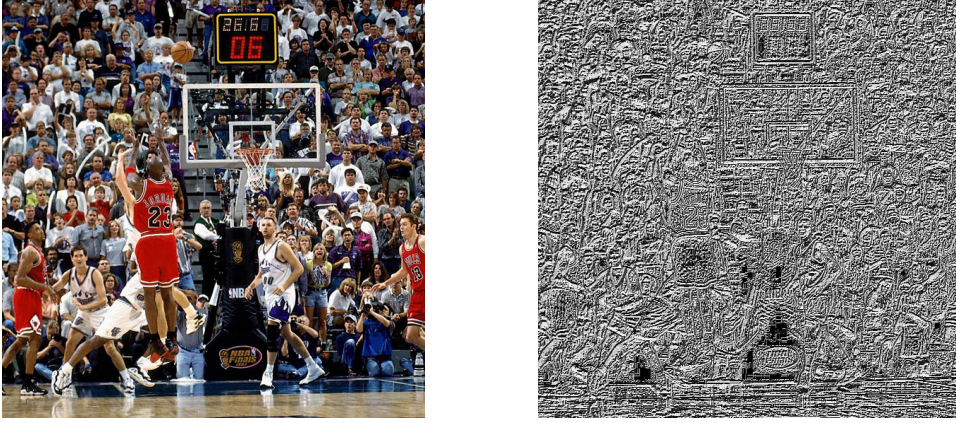


FIGURE 3.5: One LBP example

$$\begin{aligned} I_x &= \frac{\partial I}{\partial x}, \\ I_y &= \frac{\partial I}{\partial y}, \end{aligned} \quad (3.1)$$

the gradient can be computed fast by some discrete derivative masks below, like 1-D Sobel masks:

$$\begin{aligned} \text{Centered} : M_c &= [-1, 0, 1] \\ \text{Uncentered} : M_{uc} &= [-1, 1] \end{aligned} \quad (3.2)$$

or 2-D Sobel masks:

$$\begin{aligned} D_x &= \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \\ D_y &= \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \end{aligned} \quad (3.3)$$

or 3×3 Sobel masks:

$$\begin{aligned} S_x &= \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \\ S_y &= \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} \end{aligned} \quad (3.4)$$

Using different masks will lead to different performance. Besides, gaussian smoothing is often performed before gradient computing. It has been shown that using 1-D Sobel without gaussian smoothing has the best performance. A HOG

feature demo is shown in Figure 3.6.

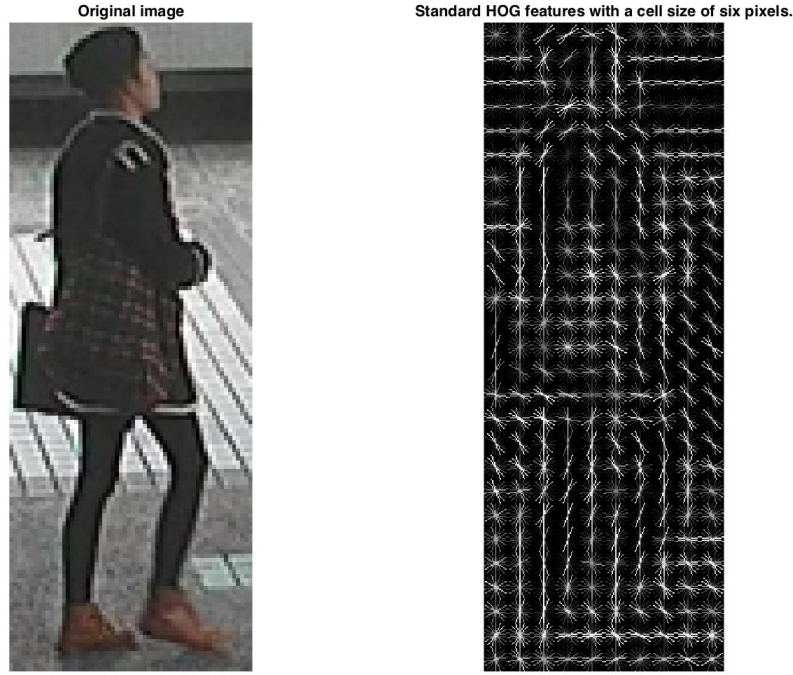


FIGURE 3.6: A demo of HOG feature with a cell size of six pixels

3.2 Influence of background segmentation on different basic descriptors

Many works try to minimize impact of background noise of pedestrians' image. It is easier to automatically segment foreground from a sequential frames or video than a single frame. In [15] the author provides foreground masks for all images following the algorithm in [18]. Authors in [18] propose stel (structure element) component analysis to model spatial correlations in image class structure by using probabilistic index maps (PIMs). A structure element is an area of an image with the same assigned index s ($s = 1, 2, \dots, n$), n is the total number of stels in an image. In PIMs the indices assignment can be denoted probabilistically as a map $q(s_i^t = s)$ over an image, in this equation the location is $i = (i, j)$, q is the probability of pixel at i of the t -th image (suppose there are many images for each class) belonging to the s -th stel. The authors propose that pixels in a structure element of an image class follow a

shared distribution which can be modelled with local measurements like pixel intensity value. The number of stels n is set to 2 to achieve background and foreground segmentation. That is, pixels of foreground stel share a distribution while pixels of background stel sharing another distribution. Both the foreground and background stels are modelled by a mixture of C Gaussians. Some of those segmented foregrounds are shown in Figure 3.7 and it is obvious that certain body parts like head and feet are lost. To compare those loss impact on color and textural descriptors, a comparison of foreground segmentation on HSV color histogram descriptors, LBP and HOG descriptor is given in Figure 3.8, Figure 3.9 and Figure 3.10.



FIGURE 3.7: Foreground segmentation of individuals from VIPeR

We can find that foreground segmentation decreases LBP and HOG's performance but increases HSV color histogram's performance on VIPeR dataset greatly. The reason for this is imperfect foreground segmentation causes body parts (like head and feet) loss and masks out many parts in torso and legs. Besides, in images of some individuals, a part of background scene is regarded as foreground. Since HSV color histogram doesn't handle spatial distribution but only color entropy, foreground segmentation improves its performance greatly. But since LBP and HOG handle texture for each sample patch, their performance suffer from those body parts loss and little black patches from background. What's more, we can infer that imperfect foreground segmentation will also decrease other textural feature's performance.

3.3 The hierarchical gaussian descriptor

The hierarchical gaussian descriptor is proposed by in [25], this descriptor uses a two-level gaussian distribution to model an individual. This descriptor densely samples the image and models each hierarchical structure with gaussian distribution and to outperform many other works. In this thesis, all input images are sized to 128×64 .

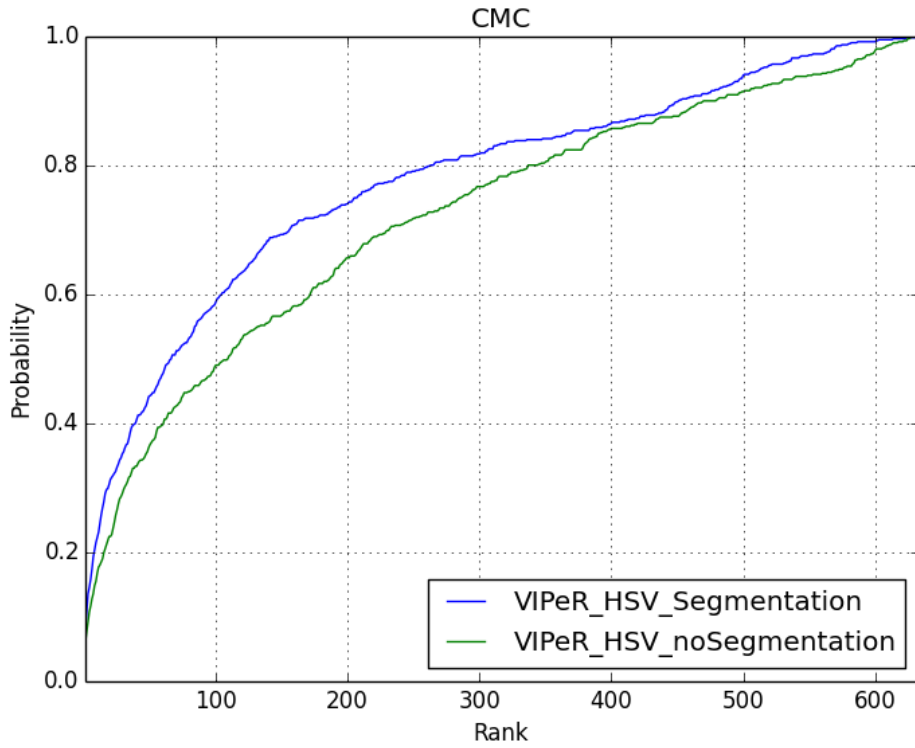


FIGURE 3.8: A CMC comparison of foreground segmentation on HSV histogram descriptor tested on VIPeR

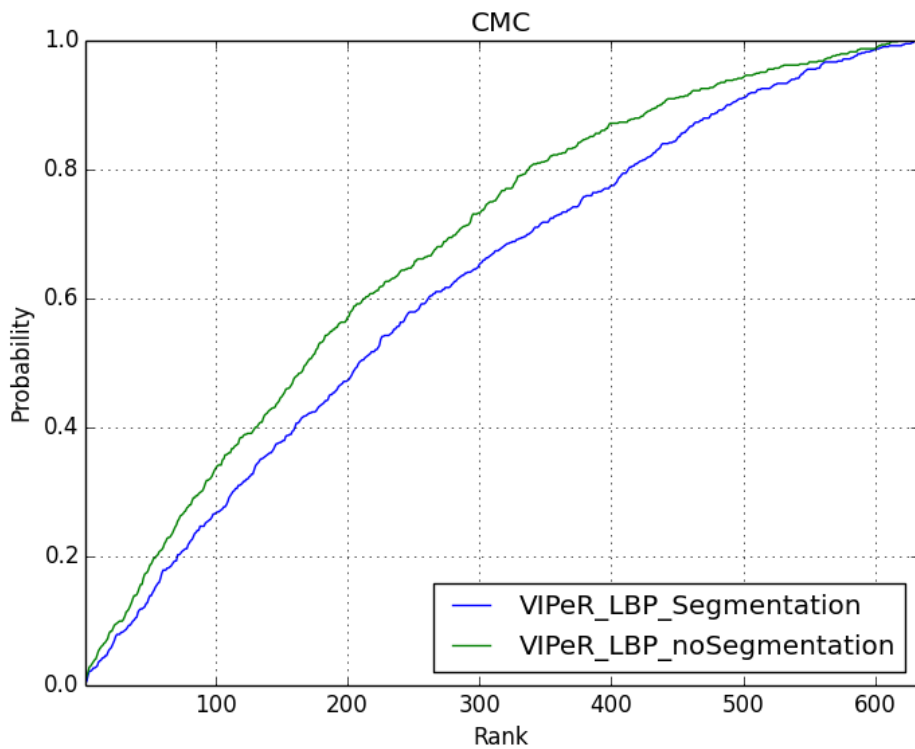


FIGURE 3.9: A CMC comparison of foreground segmentation on LBP feature tested on VIPeR

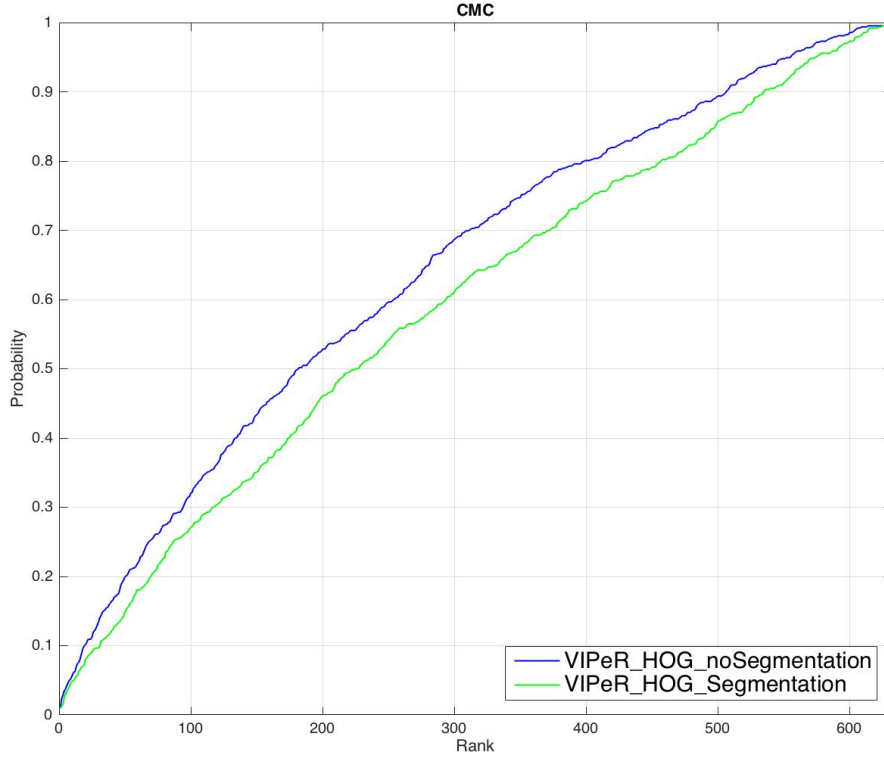


FIGURE 3.10: A CMC comparison of foreground segmentation on HOG feature tested on VIPeR

Firstly it divides the image into 7 overlapping horizontal slices, each slide has size 32×64 and slides are overlapping by 16 pixels vertically. In each slide, densely sampled square patches have a size of $s \times s$ pixels ($s = 5$ in this thesis), and small patches overlaps with each other by 2 pixels. So there is a two-level structure in this image, small patches and slides. The small patches are first modelled with a multivariate gaussian distribution, then with those small patch gaussian distributions, the slide containing those patches is modelled with another multivariate gaussian distribution.

In a certain horizontal slice for each small patch, it is modelled by a multivariate gaussian distribution $G_p(\mathbf{f}_i; \boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$, and G_p can be transformed to a vector \mathbf{p} by SPD mapping. Again, when all small patches are modelled and vectorized, the same process is repeated. Each slide is characterized by a multivariate gaussian distribution $G_r(\mathbf{p}; \boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r)$.

After the slide is modelled with $G_r(\mathbf{p}; \boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r)$, the same transformation will be operated on G_r so that it is vectorized as a vector \mathbf{v} . At last all computed \mathbf{v} are concatenated to consist of the descriptor of current image.

3.3.1 Handling the background

In the previous section the impact of background subtraction on different features' performance have been studied. We've concluded that imperfect foreground segmentation decreases textural feature's performance. Besides, in hierarchical gaussian descriptor there is no histogram-based feature computing. So in this thesis, when computing the pixel basic feature f_i , the foreground segmentation is not adopted. But when modelling the region gaussian $G_r(p; \mu_r, \Sigma_r)$ a weighted map is computed for each patch with equation

$$N(x; \mu_0, \sigma_0) = \frac{1}{\sigma_0 \sqrt{2\pi}} \exp \frac{(x - \mu_0)^2}{\sigma_0^2} \quad (3.5)$$

and here $\mu_0 = \frac{W_0}{2}$, $\sigma_0 = \frac{W_0}{4}$, W_0 is the number of patches in horizontal direction.

3.3.2 Single pixel modelling

In this hierarchical model, it is very important to have a full representation for every single pixel. To fully characterize single pixel, a d -dimensional vector is used to represent it. In this vector, there could be any predefined properties like coordinates, color values, texture and filter response. Suppose the original image is in RGB color space, the gaussian of gaussian descriptor uses a 8-dimensional vector f_i , and $f_i = (y, M_0, M_{90}, M_{180}, M_{270}, R, G, B)$. The y component is the y coordinate of pixel, and $M_{\{\theta \in 0^\circ, 90^\circ, 180^\circ, 270^\circ\}}$ is the quantized gradient information in 4 directions. The last three components are the color values of the specified color space.

In all the benchmark dataset, all the images are cropped with a bounding box, and the pedestrian in an image can be at left or right of center, while in the vertical direction the head and feet of pedestrian is very close the image edge. For each pixel, the y coordinate is more correlated than x coordinate, so only y coordinate is chosen for pixel modelling.

Then the M is to characterize the texture with the gradient histogram. Different M values is the magnitude of gradient in every direction. Firstly the gradient in x and y direction are computed by two gradient filters h_x and h_y , and we have

$$\begin{aligned} h_x &= [-1, 0, 1] \\ h_y &= -h'_x \end{aligned} \quad (3.6)$$

Then by convolving those two filters with the intensity image I , the horizontal and vertical gradient I_x, I_y can be computed, so the orientation and magnitude can be computed by following equations:

$$\begin{aligned} O(i, j) &= (\arctan(\frac{I_y(i, j)}{I_x(i, j)} + \pi) * 180/\pi \\ M(i, j) &= \sqrt{(I_x(i, j)^2 + I_y(i, j))^2} \end{aligned} \quad (3.7)$$

The orientation are quantized into four bins by a soft voting algorithm [31]. For each pixel its corresponding gradient orientation is decided by its nearest bin's direction. To make the descriptor to focus on the gradient components with high values, the gradient and orientation are multiplied as follow,

$$M_\theta = MO_\theta, \quad (3.8)$$

To model the patch with a multi-variate gaussian distribution, we have to estimate its mean value and the covariance matrix. A multi-variate gaussian model has the form

$$G_p(\mathbf{f}_i; \boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p) = \frac{\exp(\frac{1}{2}(\mathbf{f}_i - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1} (\mathbf{f}_i - \boldsymbol{\mu}_p))}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_p|} \quad (3.9)$$

where $\boldsymbol{\mu}_p$ is the estimated mean value, and $\boldsymbol{\Sigma}_p$ is the estimated covariance matrix of current small patch.

To estimate the parameters for this gaussian model based on sampled patches pixel features, the maximal likelihood estimate(MLE) is used. According MLE algorithm, we have the following estimated parameters

$$\boldsymbol{\mu}_p = \frac{1}{n} \sum_{i=1}^n \mathbf{f}_i, \quad (3.10)$$

$$\boldsymbol{\Sigma}_p = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{f}_i - \boldsymbol{\mu})(\mathbf{f}_i - \boldsymbol{\mu})^T, \quad (3.11)$$

where n is the number of pixels in current patch. When the gaussian model is computed, the next step is to model all the patch gaussians. But it is a complex problem to directly model those multivariate gaussian functions. So some transformation will be operated on the estimated parameters $\boldsymbol{\mu}_p$ and $\boldsymbol{\Sigma}_p$.

3.3.3 Riemannian manifold based SPD transformation

As described before this hierarchical gaussian descriptor is a stochastic feature, so operations like computing mean and covariance need to be operated on previous summarized gaussian distributions. Mean and covariance operation in Euclidean space can not be directly finished on previous estimated gaussian functions. A transformation is needed to make stochastic summarization feasible on patch gaussian function. In fact, the multivariate gaussian model is a Riemannian manifold and can be embedded into a semi-positive definite matrix (SPD) space. The gaussian function is mapped into a vector space by Equation 3.12. A d -dimensional multivariate gaussian function can be mapped into a $d + 1$ -dimensional SPD_+ space. According to [23], the mapping can be denoted as

$$G(\mathbf{x}_i; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \sim \mathbf{P}_i = |\boldsymbol{\Sigma}_i|^{1/(d+1)} \begin{bmatrix} \boldsymbol{\Sigma}_i + \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T & \boldsymbol{\mu}_i \\ \boldsymbol{\mu}_i^T & 1 \end{bmatrix} \quad (3.12)$$

The covariance matrix $\boldsymbol{\Sigma}_i$ can be singular for small number of pixels within the patch, to avoid this problem a regular factor λ is added to $\boldsymbol{\Sigma}_i$ so that $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_i + \lambda \mathbf{I}$.

After this mapping, the $n + 1$ -dimensional SPD matrix needs to be transformed into a vector. The matrix logarithm is used to transform it to tangent space. A $d + 1$ -dimensional SPD matrix can be mapped as a $d * (d + 3)/2 + 1$ vector, which can be denoted as $SPD_i^+ \sim \mathbf{p}_i = \text{vec}(\log(\mathbf{P}_i))$. Since \mathbf{P}_i is a positive symmetric matrix, it can be compressed by half, i.e only the upper triangular elements are preserved. To ensure its norm-1 remains the same after compression, the magnitude of off-diagonal elements in \mathbf{P}_i are timed by $\sqrt{2}$. Let $\mathbf{Q} = \log \mathbf{P}_i$, we have

$$\mathbf{p}_i = [\mathbf{Q}_{1,1}, \sqrt{2}\mathbf{Q}_{1,2}, \sqrt{2}\mathbf{Q}_{1,3}, \dots, \sqrt{2}\mathbf{Q}_{1,d+1}, \quad (3.13)$$

$$\mathbf{Q}_{2,2}, \sqrt{2}\mathbf{Q}_{2,3}, \dots, \sqrt{2}\mathbf{Q}_{2,d+1}, \dots, \mathbf{Q}_{d+1,d+1}] \quad (3.14)$$

Again we model the slide with a multivariate gaussian distribution \mathbf{G}_r by equation

$$G_p(\mathbf{p}_j; \boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r) = \frac{\exp(\frac{1}{2}(\mathbf{p}_j - \boldsymbol{\mu}_r)^T \boldsymbol{\Sigma}_r^{-1} (\mathbf{p}_j - \boldsymbol{\mu}_r))}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_r|} \quad (3.15)$$

we have known that all patches inside this slide been represented as vector \mathbf{p}_j , thus

the mean μ_r and covariance matrix Σ_r of G_r can be computed by MLE with equations

$$\mu_r = \frac{1}{m} \sum_{j=1}^m p_j, \quad (3.16)$$

$$\Sigma_r = \frac{1}{m-1} \sum_{j=1}^m (p_j - \mu_r)(p_j - \mu_r)^T, \quad (3.17)$$

m is the total number of patches inside current slide. Again μ_r and Σ_r are mapped to a SPD space by equation

$$G(p; \mu_r, \Sigma_r) \sim P_r = |\Sigma_r|^{1/(d+1)} \begin{bmatrix} \Sigma_r + \mu_r \mu_r^T & \mu_r \\ \mu_r^T & 1 \end{bmatrix} \quad (3.18)$$

and P_r is vectorized by Equation 3.13.

When all horizontal slides' descriptor computed they are concatenated to form the descriptor for the whole image.

3.3.4 Integral image for fast region covariance computation

To compute estimated parameters for all those overlapping small patches, the time complexity of computing one by one each patch is high because there are many repeating computations. To compute the estimated covariance matrix Σ for every small patch with size of $W \times H$, the integral image is used to reduce time complexity. The integral image [40] is an intermediate representation to fast compute rectangle area sum in an image. Each pixel value in integral image is the sum of all the pixels inside the rectangle bounded by current pixel and the upper left pixel. That is, the integral image $S(x, y)$ for image $I(x, y)$ is

$$S(x', y') = \sum_{x < x', y < y'} I(x, y), \quad (3.19)$$

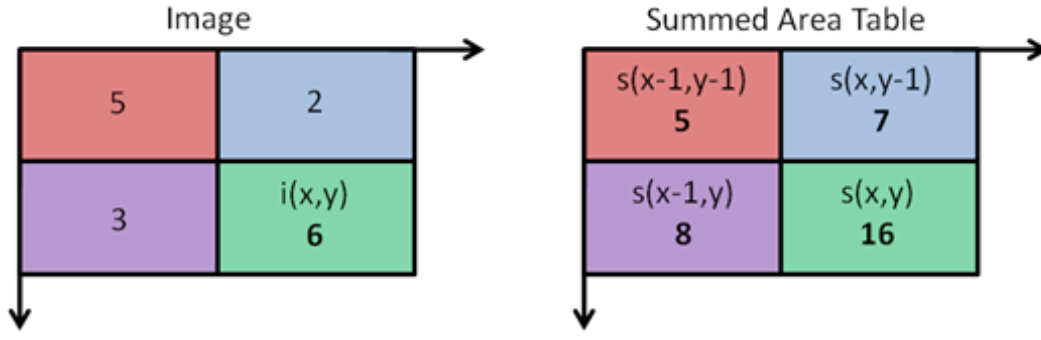


FIGURE 3.11: Integral image

By using integral image any rectangular region sum can be computed in constant time.

To compute the covariance matrix of a certain rectangle area in a $W \times H \times d$ -dimensional feature tensor F , suppose \mathbf{I}_F is the $W \times H \times d$ tensor of integral images of F , we have

$$\mathbf{I}_F(x', y', i) = \sum_{x < x', y < y'} F(x, y, i), i = 1 \dots d \quad (3.20)$$

and suppose the $\mathbf{C}(x', y', i, j)$ is the $W \times H \times d \times d$ tensor of second order integral images, we have

$$\mathbf{C}(x', y', i, j) = \sum_{x < x', y < y'} F(x, y, i) F(x, y, j), i, j = 1 \dots d. \quad (3.21)$$

let $\mathbf{I}_{x,y}$ be the d -dimensional vector in \mathbf{I}_F , $\mathbf{C}(x, y)$ be the $d \times d$ -dimensional matrix in \mathbf{C} ,

$$\begin{aligned} \mathbf{I}_{x,y} &= [\mathbf{I}_F(x', y', 1) \dots \mathbf{I}_F(x', y', d)]^T \\ \mathbf{C}_{x,y} &= \begin{bmatrix} \mathbf{C}(x, y, i, 1) & \dots & \mathbf{C}(x, y, i, d) \\ & \ddots & \\ \mathbf{C}(x, y, d, 1) & \dots & \mathbf{C}(x, y, d, d) \end{bmatrix} \end{aligned} \quad (3.22)$$

Then for any rectangle regions $R(x', y'; x'', y'')$, where (x', y') is the upper left coordinate and (x'', y'') is the lower right coordinate, the covariance matrix can be

compute as

$$\begin{aligned} \mathbf{C}_R(x', y'; x'', y'') = & \frac{1}{n-1} [\mathbf{C}_{x'', y''} + \mathbf{C}_{x', y'} - \mathbf{C}_{x'', y'} - \mathbf{C}_{x', y''} \\ & - \frac{1}{n} (\mathbf{I}_{x'', y''} + \mathbf{I}_{x'', y''} - \mathbf{I}_{x', y''} - \mathbf{I}_{x'', y'}) (\mathbf{I}_{x'', y''} + \mathbf{I}_{x'', y''} - \mathbf{I}_{x', y''} - \mathbf{I}_{x'', y'})^T] \end{aligned} \quad (3.23)$$

where n is the number of feature vector in F , and $n = (x'' - x')(y'' - y')$. By creating the integral image the covariance of any rectangular area in F can be computed in $O(d^2)$ time.

When all patches in a region are computed, the same process is repeated to compute the region gaussian.

3.3.5 Some variants of Hierarchical Gaussian descriptor

In hierarchical gaussian descriptor the basic pixel feature \mathbf{f}_i characterizes texture, color and y coordinate of an image. The importance of those three characteristics is ordered as color values > gradient components > y coordinate. A rank-1, rank-5, rank-10, and rank-20 comparison of three $\text{GOG}_{\text{fusion}}$ variants on VIPeR dataset are listed below. We can conclude that color values have the most influence on Re-ID accuracy. Based on this conclusion, in $\text{GOG}_{\text{fusion}}$ variants studied the color values will remain unchanged and y coordinate and M components will be replaced or removed. In one $\text{GOG}_{\text{fusion}}$ variant we replace the M components in \mathbf{f}_i with LBP.

TABLE 3.1: A rank score comparison of $\text{GOG}_{\text{fusion}}$ Variants

GOG _{fusion} Variants	Rank(%)			
	1	5	10	20
GOG _{fusion}	47.97	77.44	86.80	93.70
GOG _{fusion} without y	44.80	76.20	86.80	93.20
GOG _{fusion} without M	41.10	70.50	81.10	90.30
GOG _{fusion} without color values	8.40	21.00	31.70	45.50

Another variant is to combine superpixel segmentation with $\text{GOG}_{\text{fusion}}$. Superpixel algorithms cluster image pixels into perceptually atomic regions, which can be used to reduce redundancy and computation complexity. One successful superpixel algorithm is SLIC [1]. To combine superpixel segmentation and $\text{GOG}_{\text{fusion}}$, each horizontal slice of an image is first segmented by SLIC algorithm into many nonrigid

groups (An example of SLIC segmentation is shown in Figure 3.12). Each superpixel group is modelled with a multivariate gaussian function by summarizing basic feature f_i of pixels inside this group. With similar SPD mapping and vectorization every group can be summarized by a vector. At last the same process is repeated on those superpixels to get region gaussians. Compared with original hierarchical gaussian descriptor, the difference is overlapping square patch sampling in hierarchical gaussian are changed to nonrigid and non-overlapping superpixel sampling in the variant. Hopefully, superpixel segmentation can group pixels with similar property, which may decrease error caused by single gaussian modelling.

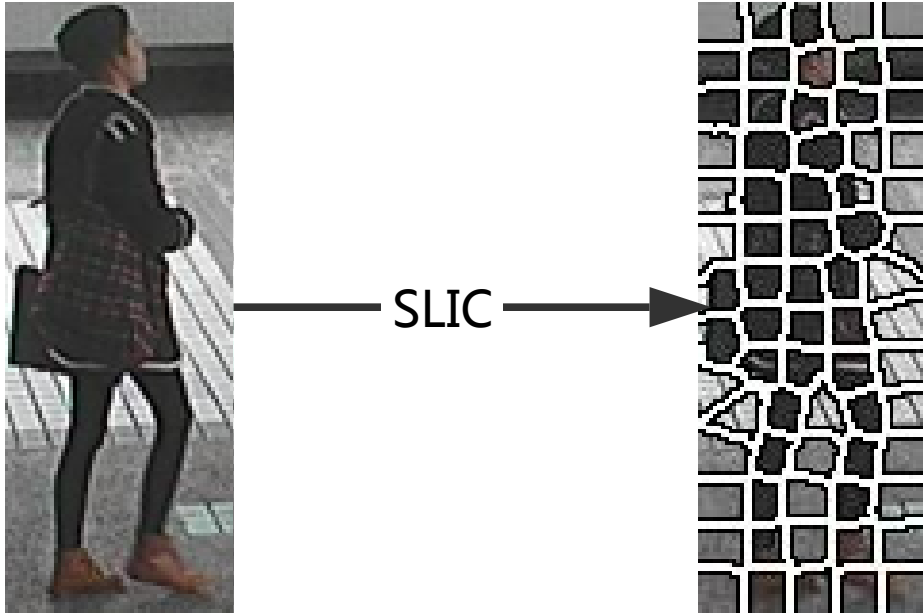


FIGURE 3.12: A SLIC superpixel segmentation example

Besides, gaussian mixture model (GMM) based descriptor is also tested. GMM models an image or a region with a mixture of gaussian distributions. But it is computationally expensive to compute similarity of GMMs. Also, it is extremely difficult to create a two-level hierarchical GMM model. Sparse represented Earth Mover's Distance (SR-EMD) [21] is used here to measure similarity. All those variants' performance have been listed in Table 3.2. It can be concluded that GOG_{fusion} has the best performance. $GOG_{fusion}^{MreplacedByLBP}$'s performance decreases heavily because LBP components in f_i is less robust. $GOG_{fusion}^{Superpixel}$ might suffer from its non-overlapping superpixel patch sampling. GMM + SR-EMD model has worst

performance because of two reasons: (1) This model is not trained; (2) A better similarity measurement is needed to for GMMs.

TABLE 3.2: A performance comparison between Gaussian of Gaussian descriptor and its variants on VIPeR dataset

Variant descriptors	Rank(%)				
	1	5	10	15	20
GOG _{fusion} +XQDA	47.97	77.44	86.80	91.27	93.70
GOG _{fusion} MreplacedByLBP+XQDA	40.70	72.53	83.16	88.45	91.90
GOG _{fusion} Supapixel+XQDA	42.72	74.84	85.22	89.34	92.25
GMM + SR-EMD(EMD-theta)	11.40	21.50	29.70	37.00	42.40

3.3.6 Dimensionality and superiority analysis of Hierarchical Gaussian descriptor

It has been known that combination of descriptors of different color space can greatly improve re-ID performance. In this project, the hierarchical gaussian descriptor in RGB color space is the base descriptor. Descriptors in three more color space {HSV, Lab, nRGB} are extracted. The nRGB color space is calculated as

$$\begin{aligned}
 nR &= \frac{R}{R + G + B}, \\
 nG &= \frac{G}{R + G + B}, \\
 nB &= \frac{B}{R + G + B},
 \end{aligned} \tag{3.24}$$

since nB can be calculated with nR and nG , in this color space only the first two channel values are used to reduce redundancy. Therefore, for color spaces {RGB, HSV, Lab, nRnG}, the corresponding dimension of pixel feature is {8, 8, 8, 7}. After the matrix to vector transformation, the dimension of patch gaussian vector of each channel is {45, 45, 45, 36}. Again after the patch gaussian to region gaussian transformation, the dimension of each channel is {1081, 1081, 1081, 703}. Suppose there are 7 horizontal slides in each image, the dimension of concatenated descriptor of each channel is {7567, 7567, 7567, 4921}. If four color space are all used, the dimension is the sum of each channel as 27622.

Hierarchical gaussian descriptor has a few advantages compared with other descriptors. Firstly, it has a full consideration of color, texture and y coordinate information. The color information is adopted by adding color values of different color space to \mathbf{f}_i . Textural information is also adopted by those four gradient components in \mathbf{f}_i . Secondly, by Riemannian manifold based SPD mapping, it provides a solution to summarize many multi-variate gaussian functions. The parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ of a multi-variate gaussian function can be fused into a SPD matrix. With later vectorization this SPD matrix is transformed to a vector. In this mapping process correlation of different components in \mathbf{f}_i are fully taken into account, which leads to a dimension increasing from basic pixel feature vector \mathbf{f}_i to patch gaussian vector \mathbf{p}_j . Thirdly, by a two-level gaussian model, this model guarantees to extract the overall statistical information of an image while being robust to local details variation caused by factors like viewpoint changes and partial occlusion.

It is intractable to directly learn a Mahalanobis distance matrix from concatenated Hierarchical Gaussian descriptor. In next Chapter KLFDA is applied to reduce the dimensionality of extracted descriptor from 27622 to $n - 1$, where n is the number of different classes. With this supervised nonlinear dimensionality reduction, discriminative information among different classes are preserved. At last a Mahalanobis matrix is learned on the dimension reduced descriptors by gradient descent method.

Chapter 4

Dimension reduction and Mahanalobis distance learning

To deal with high-dimensional descriptors, dimension reduction are firstly performed by kernel local fisher discriminant analysis (KLFDA). Then Mahanalobis distance metric learning based on limitations between interclass and intraclass distance is applied on dimension reduced data.

4.1 Kernel local fisher discriminant analysis

The extracted hierarchical gaussian descriptors have high dimension, it is intractable to learn a SPD matrix with such high dimension. Dimension reduction is required to learn a subspace. Among those methods to reduce dimension, principal component analysis (PCA) is often used. However, PCA is an unsupervised dimension reduction and may have low performance for two reasons. (1), PCA is to maximize the variance of dimension-reduced data, and as an unsupervised method it doesn't take into consideration the between and within classes' relations. It is very likely that the descriptors of different classes can be mixed up after dimension reduction; (2) PCA may suffer from the small sample size problem. In some Re-ID datasets, there may be one or two images for each pedestrian in each viewpoint (like VIPeR), if the dimension of the descriptor is much bigger than sample size, much information can be lost with PCA. In this thesis, the kernel local fisher discriminant analysis (KLFDA) is used to reduce dimension.

KLFDA is the kernel version of LFDA, and LFDA is a combination of Fisher discriminant analysis [33], the locality preserving projection [17] and kernel method. A brief introduction of FDA, LPP and kernel method is given below.

4.1.1 Fisher discriminant analysis (FDA)

FDA is a supervised dimension reduction and its input contains the class labels.

Given a set of d -dimensional observations $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$, where $i \in \{1, 2, \dots, n\}$, the label $y_i \in \{1, 2, \dots, l\}$, for each sample descriptor \mathbf{x}_i , a linear transformation with transformation matrix \mathbf{T} can be defined by equation

$$\mathbf{z}_i = \mathbf{T}^T \mathbf{x}_i \quad (4.1)$$

\mathbf{T} has dimension $d \times m$, \mathbf{z}_i is the m ($m < d$) dimensional vector. In FDA two matrices are defined as the intraclass scatter matrix $\mathbf{S}^{(w)}$ and between class scatter matrix $\mathbf{S}^{(b)}$,

$$\begin{aligned} \mathbf{S}^{(w)} &= \sum_{i=1}^l \sum_{j:y_j=i} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \\ \mathbf{S}^{(b)} &= \sum_{i=1}^l n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T \end{aligned} \quad (4.2)$$

where n_i is the number of classes with class label i , $\boldsymbol{\mu}_i$ is the mean of samples whose label is i , and $\boldsymbol{\mu}$ is the mean of all samples,

$$\begin{aligned} \boldsymbol{\mu}_i &= \frac{1}{n_i} \sum \mathbf{x}_i, \\ \boldsymbol{\mu} &= \frac{1}{n} \sum \mathbf{x}_i \end{aligned} \quad (4.3)$$

The Fisher Discriminant Analysis transform matrix \mathbf{T} can be represented as

$$\mathbf{T} = \arg \max \frac{\mathbf{T}^T \mathbf{S}^{(b)} \mathbf{T}}{\mathbf{T}^T \mathbf{S}^{(w)} \mathbf{T}} \quad (4.4)$$

This equation can be solved by Lagrange multiplier method, we define a Lagrange function

$$L(\mathbf{t}) = \mathbf{t}^T \mathbf{S}^{(b)} \mathbf{t} - \lambda (\mathbf{t}^T \mathbf{S}^{(w)} \mathbf{t} - 1) \quad (4.5)$$

Then the differential respect to \mathbf{t} is

$$\frac{\partial L(\mathbf{t})}{\partial \mathbf{t}} = 2\mathbf{S}^{(b)} \mathbf{t} - 2\lambda \mathbf{S}^{(w)} \mathbf{t} \quad (4.6)$$

let

$$\frac{\partial L(\mathbf{t})}{\partial \mathbf{t}} = 0 \quad (4.7)$$

we can get

$$\mathbf{S}^{(b)} \mathbf{t}_i = \lambda \mathbf{S}^{(w)} \mathbf{t}_i \quad (4.8)$$

here \mathbf{t}_i is the i_{th} column of \mathbf{T} , and the optimization problem is converted to a eigenvalue decomposition problem. \mathbf{T} is the set of eigenvectors of $\frac{\mathbf{S}^{(b)}}{\mathbf{S}^{(w)}}$.

Fisher discriminant analysis tries to minimize the intraclass scatter matrix while maximizing the interclass scatter matrix. \mathbf{T} is computed by the eigenvalue decomposition. \mathbf{T} can be represented as the set of all the corresponding eigenvectors, as $\mathbf{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_k)$.

FDA has a form similar to signal and noise ratio, however, the FDA dimension reduction may have poor performance because it doesn't consider the locality of data. An example of this is the multimodality [38]. Multimodality is when many clusters are formed in the same class.

4.1.2 Locality preserving projection (LPP)

In [17] locality preserving projection (LPP) is proposed to exploit data locality. An affinity matrix is created to record the affinity of sample \mathbf{x}_i and \mathbf{x}_j , typically the range of elements in $\mathbf{A}_{i,j}$ is $[0, 1]$. There are many manners to define a $n \times n$ affinity matrix \mathbf{A} , usually two sample points with a smaller distance has a higher affinity value than those with bigger distance value. One of them is if \mathbf{x}_i is within k-nearest neighbours of \mathbf{x}_j then $\mathbf{A}_{i,j} = 1$ otherwise $\mathbf{A}_{i,j} = 0$.

Another diagonal matrix \mathbf{D} can be defined that each diagonal element is the sum of corresponding column in \mathbf{A} ,

$$\mathbf{D}_{i,i} = \sum_{j=1}^n \mathbf{A}_{i,j} \quad (4.9)$$

then the LPP transform matrix is defined as follow,

$$\mathbf{T}_{LPP} = \arg \min_{\mathbf{T} \in \mathbf{R}^{d \times m}} \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j} ||\mathbf{T}^T \mathbf{x}_i - \mathbf{T}^T \mathbf{x}_j|| \quad (4.10)$$

so that $\mathbf{T}^T \mathbf{X} \mathbf{D} \mathbf{X}^T \mathbf{T} = \mathbf{I}$. Suppose the subspace has a dimension of m , then LPP transform matrix \mathbf{T} can be represented as

$$\mathbf{T}_{LPP} = \{\phi_{d-m+1} | \phi_{d-m+2} | \dots | \phi_d\}$$

and each ϕ in T is the eigenvector of following fomula,

$$\mathbf{X} \mathbf{L} \mathbf{X}^T \phi = \gamma \mathbf{X} \mathbf{D} \mathbf{X}^T \phi \quad (4.11)$$

where γ is corresponding eigenvalue of ϕ , and $\mathbf{L} = \mathbf{D} - \mathbf{A}$.

4.1.3 Local fisher discriminant analysis (LFDA)

LFDA [33] combines FDA and LPP and has better performance. The key in LFDA is it assigns weights to elements in $\mathbf{A}^{(w)}$ and $\mathbf{A}^{(b)}$, so that,

$$\begin{aligned} \mathbf{S}^{(w)} &= \frac{1}{2} \sum_{i=1}^l \sum_{j:y_j=i} \mathbf{A}_{i,j}^w (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \\ \mathbf{S}^{(b)} &= \frac{1}{2} \sum_{i=1}^l \mathbf{A}_{i,j}^b (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T \end{aligned} \quad (4.12)$$

where

$$\begin{aligned} \mathbf{A}_{i,j}^{(w)} &= \begin{cases} \mathbf{A}_{i,j}/n_c & y_i = y_j \\ 0 & y_i \neq y_j \end{cases} \\ \mathbf{A}_{i,j}^{(b)} &= \begin{cases} (\frac{1}{n} - \frac{1}{n_c}) \mathbf{A}_{i,j} & y_i = y_j \\ \frac{1}{n} & y_i \neq y_j \end{cases} \end{aligned} \quad (4.13)$$

where y_i is the class label of sample point \mathbf{x}_i . So the transformation matrix \mathbf{T}_{LFDA} can be computed by equation

$$\mathbf{T}_{LFDA} = \arg \min_{\mathbf{T}} \left(\frac{\mathbf{T}^T \mathbf{S}^{(b)} \mathbf{T}}{\mathbf{T}^T \mathbf{S}^{(w)} \mathbf{T}} \right) \quad (4.14)$$

Again this problem can be solved by eigenvalue decomposition by Equation 4.8.

When applying the LFDA to original high-dimensional descriptors, one problem is the computation cost. Suppose the vector data has a dimension of d , LFDA has to solve the eigenvalue a matrix with dimension $d \times d$. In some descriptors d could be more than 20000 and the computation cost is intractable.

4.1.4 Kernel local fisher discriminant analysis (KLFDA)

KLFDA [38] is the nonlinear version of LFDA. Most dimensionality reduction methods including PCA, LDA and LFDA are linear dimensionality reduction methods. However, when descriptors data are non-linear in feature space, its hard to capture its between-class discriminant information with linear reduction methods. One alternative method is to nonlinearly map input descriptors \mathbf{x}_i to higher-dimensional feature space Φ by a function $\phi(\mathbf{x}_i)$, again the LFDA is performed in feature space Φ . Thus the transformation matrix \mathbf{T} can be computed by equation

$$\mathbf{T} = \arg \min \frac{\mathbf{T}^T \mathbf{S}_{\phi}^{(b)} \mathbf{T}}{\mathbf{T}^T \mathbf{S}_{\phi}^{(w)} \mathbf{T}} \quad (4.15)$$

where $\mathbf{S}_{\phi}^{(b)}$ and $\mathbf{S}_{\phi}^{(w)}$ is the between class scatter and within class scatter in mapped feature space Φ .

Note that the transformation matrix $\mathbf{T} \in \Phi$, it is computationally expensive to explicitly compute the mapping function ϕ and perform LFDA in feature space Φ because the dimension of Φ may be infinite. Rather than explicitly computing, the mapping function ϕ can be implicit and the feature space Φ can be defined by the inner product of features in Φ . Kernel trick is used here and a kernel function can be defined as the inner product of mapped vectors $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$ by equation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle, \quad (4.16)$$

the $\langle \cdot \rangle$ is the inner product. There are many kinds of kernel like linear kernel, polynomial kernel and radial basis function (RBF) kernel. In this paper the RBF kernel is adopted. A RBF kernel is defined as

$$k_{RBF}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2). \quad (4.17)$$

Suppose \mathbf{X} is the sample descriptors matrix, and we have

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n), \quad (4.18)$$

and the label vector is $\mathbf{l} = (l_1, l_2, \dots, l_n)$. Then the kernel matrix of \mathbf{X} can be computed as following equation:

$$\mathbf{K} = \phi(\mathbf{X})^T \phi(\mathbf{X}) \quad (4.19)$$

and we have

$$\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2) \quad (4.20)$$

In [38] the authors proposed fast computation of LFDA by replacing $\mathbf{S}^{(b)}$ with the local scatter mixture matrix $\mathbf{S}^{(m)}$ defined by

$$\begin{aligned} \mathbf{S}^{(m)} &= \mathbf{S}^{(b)} + \mathbf{S}^{(w)} \\ \mathbf{S}^{(m)} &= \frac{1}{2} \sum_{i,j=1} \mathbf{A}_{i,j}^{(m)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T \end{aligned} \quad (4.21)$$

and

$$\mathbf{A}_{i,j}^{(m)} = \mathbf{A}_{i,j}^{(w)} + \mathbf{A}_{i,j}^{(b)} \quad (4.22)$$

according to indentify (Fukunaga, 1990)

$$\text{trace}((\mathbf{T}^T \mathbf{S}^{(w)} \mathbf{T})^{(-1)} (\mathbf{T}^T \mathbf{S}^{(m)} \mathbf{T})) = \text{trace}((\mathbf{T}^T \mathbf{S}^{(w)} \mathbf{T})^{(-1)} (\mathbf{T}^T \mathbf{S}^{(b)} \mathbf{T})) + m \quad (4.23)$$

Equation 4.14 is equal to

$$\mathbf{T}_{LFDA} = \arg \min_{\mathbf{T}} \left(\frac{\mathbf{T}^T \mathbf{S}^{(m)} \mathbf{T}}{\mathbf{T}^T \mathbf{S}^{(w)} \mathbf{T}} \right) \quad (4.24)$$

and it can be transformed into a eigenvalue decomposition problem

$$\mathbf{S}^{(m)} \mathbf{t}_i = \lambda \mathbf{S}^{(w)} \mathbf{t}_i \quad (4.25)$$

Also with the replacement of $\mathbf{S}^{(m)}$, in [38] the author summarized that

$$\mathbf{S}^{(m)} = \mathbf{X} \mathbf{L}^{(m)} \mathbf{X}^T \quad (4.26)$$

where $L^{(m)} = D^{(m)} - A^{(m)}$, and $D^{i,i} = \sum_{j=1}^n A_{i,j}^{(m)}$. Also $S^{(w)}$ can be represented as

$$S^{(w)} = XL^{(w)}X^T \quad (4.27)$$

where $L^{(w)} = D^{(w)} - A^{(w)}$, and $D^{i,i} = \sum_{j=1}^n A_{i,j}^{(w)}$. Therefore, Equation 4.25 can be represented as

$$XL^{(m)}X^T t_i = \lambda XL^{(w)}X^T t_i \quad (4.28)$$

the eigen vector t_i can be represented as $t_i = X\gamma$, vector $\gamma_i \in R^n$, with this replacement, we left multiply X^T to Equation 4.28 to get

$$X^T XL^{(m)}X^T X\gamma_i = \lambda X^T XL^{(w)}X^T X\gamma_i \quad (4.29)$$

and by the kernel trick, its represented as

$$KL^{(m)}K\gamma_i = \lambda KL^{(w)}K\gamma_i \quad (4.30)$$

One example of using KLFDA to reduce dimension and classify the nonlinear data clusters can be shown in Figures 4.1, 4.2 and 4.3. Three classes with five clusters are distributed on a 2-D plane, by KLFDA dimension reduction its 1-D dimension reduced data distribution are shown in Figure 4.2 and Figure 4.3. It shows that for those clusters the gaussian kernel are better than linear kernel because the dimensionality reduced data are more separate when using gaussian kernel function.

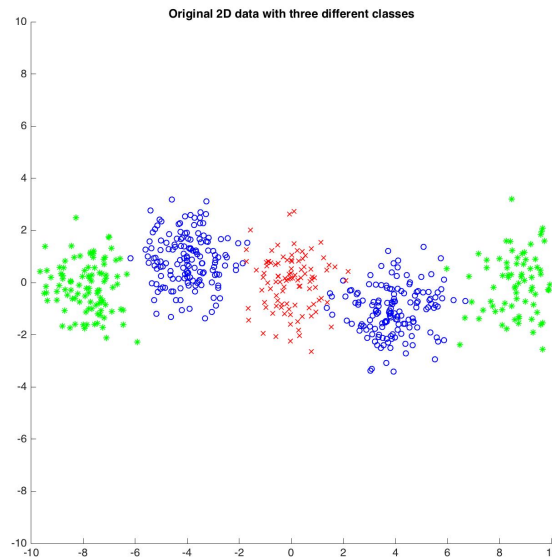


FIGURE 4.1: Example of five clusters belong to three classes

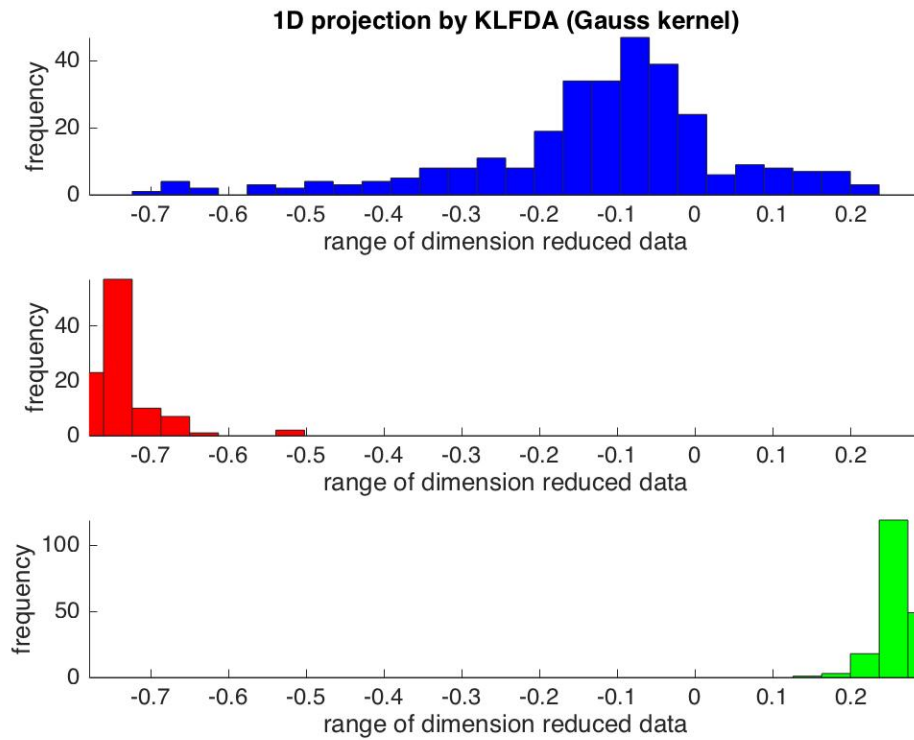


FIGURE 4.2: 1-D distribution of dimension reduced data with gaussian kernel

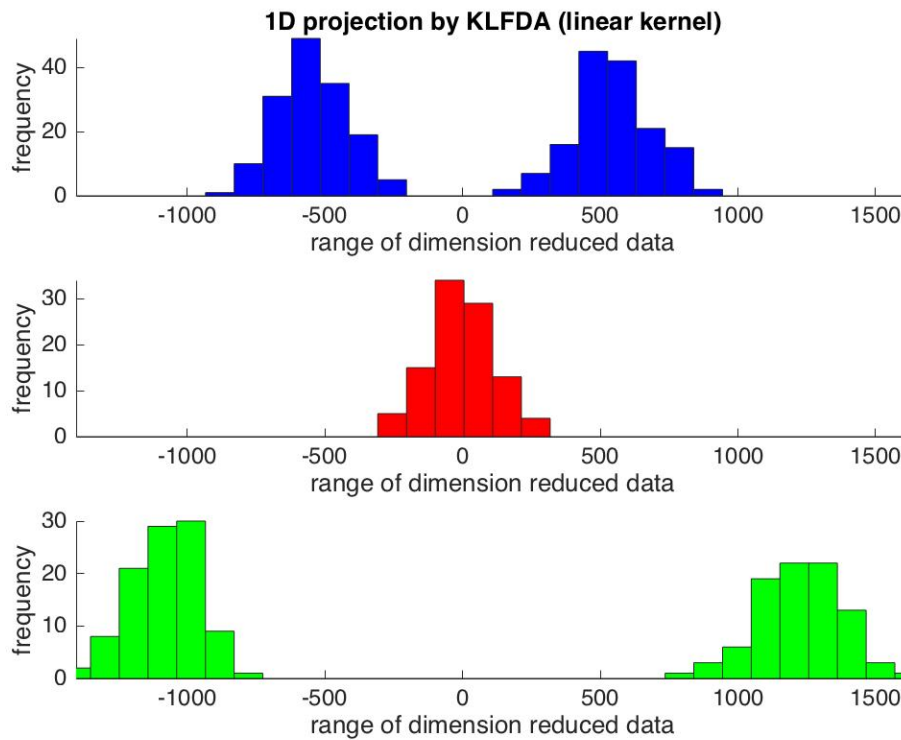


FIGURE 4.3: 1-D distribution of dimension reduced data with linear kernel

4.2 Mahalanobis distance

The Mahalanobis distance [35] based metric learning has received much attention in similarity computing. The Mahalanobis distance of two observations \mathbf{x} and \mathbf{y} is defined as

$$D(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{M} (\mathbf{x} - \mathbf{y}), \quad (4.31)$$

where \mathbf{x} and \mathbf{y} are $d \times 1$ observation vectors, \mathbf{M} is a semi-positive definite matrix. Since \mathbf{M} is a SPD matrix, \mathbf{M} can be decomposed as $\mathbf{M} = \mathbf{W}^T \mathbf{W}$, and Mahalanobis distance can also be written as

$$D(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{W}^T \mathbf{W} (\mathbf{x} - \mathbf{y}) = \|\mathbf{W}(\mathbf{x} - \mathbf{y})\| \quad (4.32)$$

Therefore, Mahalanobis distance can be regarded as a variant of Euclidean distance.

4.3 Gradient descent optimization

Given a multivariate function $F(\mathbf{x})$, \mathbf{x} is a d -dimensional vector, if $f(\mathbf{x})$ is continuous and differentiable in the neighbour of point \mathbf{x} for all \mathbf{x} , then $f(\mathbf{x})$ decreases fastest in the direction of negative gradient of F at \mathbf{x} . To compute the minimum of $F(\mathbf{x})$, an iterative method can be used by updating F with respect to \mathbf{x} . If the updating step λ is small enough, by updating \mathbf{x} with

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \lambda \mathbf{G} \quad (4.33)$$

we have

$$F(\mathbf{x}_{t+1}) \geq F(\mathbf{x}_t). \quad (4.34)$$

This process is repeated until certain condition is met, generally when gradient $\|\mathbf{G}\| \leq \eta$, η is a very small positive integer.

Analysis of steepest gradient descent method The advantages of gradient descent are that it is always downhill and it can avoid the saddle points. Besides, it is very efficient when initial value of $F(\mathbf{x})$ is far from minimum. However, there are a few shortcomings of gradient descent method. The first one is the convergence value of gradient descent might be the local minima of $F(\mathbf{x})$ if $F(\mathbf{x})$ is not monotonic. In

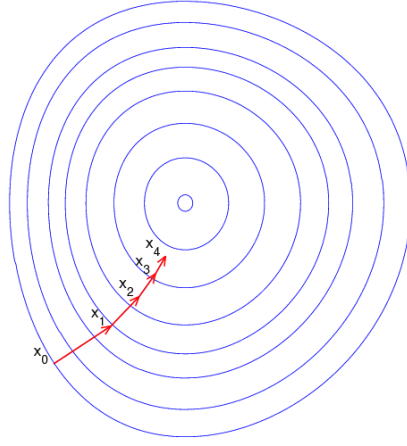


FIGURE 4.4: Steepest gradient descent

this case the convergence value will depend on the initial value of \mathbf{x} . Another shortcoming is the converging speed goes very slow when approaching the minimum, one example is zigzag downhill case. The third shortcoming is linear search in gradient descent might cause some problem.

4.4 Metric learning based on sample pairs distance comparison

Inspired by [46], a similar metric learning based on iteration computation is used. For a sample descriptor \mathbf{x}_i , its positive pairwise set is defined as $\{\mathbf{x}_i, \mathbf{x}_j\}$, where class ID $y_i = y_j$. Also the negative pairwise set can be defined as $\{\mathbf{x}_i, \mathbf{x}_j\}$, where $y_i \neq y_j$. Similar with [50], this method is also based on similarity comparison. In [50], for all possible positive and negative pairs, the distance between positive pairs must be smaller than the distance between negative pairs. Since it has to compare all possible positive and negative pairs, its computation complexity is quite expensive. To reduce the complexity, a simplified version is proposed as the top-push distance metric learning [46]. Since re-identification is a problem of ranking, it is desired that the rank-1 descriptor should be the right match. Given a Mahalanobis matrix M , for samples $\mathbf{x}_i, i = 1, 2, 3, \dots, n$, n is the number of all samples, the requirement is distance between positive pair should be at least 1 unit smaller than the minimum distance of all negative pair. This can be denoted as

$$D(\mathbf{x}_i, \mathbf{x}_j) + \rho < \min D(\mathbf{x}_i, \mathbf{x}_k), y_i = y_j, y_i \neq y_k. \quad (4.35)$$

ρ is a slack variable and $\rho \in [0, 1]$. This equation can be transformed into a optimization problem as

$$\min \sum_{y_i=y_j} \max\{D(\mathbf{x}_i, \mathbf{x}_j) - \min_{y_i \neq y_k} D(\mathbf{x}_i, \mathbf{x}_k) + \rho, 0\}. \quad (4.36)$$

However, the equation above only penalizes the minimum interclass distance. Another term is needed to penalize large intraclass distance. That is, the sum of intraclass distance should be as small as possible. This term is denoted as

$$\min \sum_{y_i=y_j} D(\mathbf{x}_i, \mathbf{x}_j). \quad (4.37)$$

To combine equations above, a ratio factor α is assigned to Equation (4.36) and (4.37) so that the target function can be denoted as

$$\begin{aligned} f(\mathbf{M}) = & (1 - \alpha) \sum_{\mathbf{x}_i, \mathbf{x}_j, y_i=y_j} D(\mathbf{x}_i, \mathbf{x}_j) + \\ & \alpha \sum_{\mathbf{x}_i, \mathbf{x}_j, y_i=y_j} \max\{D(\mathbf{x}_i, \mathbf{x}_j) - \min_{y_i \neq y_k} D(\mathbf{x}_i, \mathbf{x}_k) + \rho, 0\} \end{aligned} \quad (4.38)$$

In this way the problem is transformed to an optimization problem. Notice that Equation 4.31 can be denoted as

$$D(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \mathbf{M} (\mathbf{x} - \mathbf{y}) = \text{trace}(\mathbf{M} \mathbf{X}_{i,j}) \quad (4.39)$$

where $\mathbf{X}_{i,j} = \mathbf{x}_i * \mathbf{x}_j^T$, and trace is to compute the matrix trace. Therefore, Equation 4.38 can be transformed as follow,

$$\begin{aligned} f(\mathbf{M}) = & (1 - \alpha) \sum_{y_i=y_j} \text{trace}(\mathbf{M} \mathbf{X}_{i,j}) \\ & + \alpha \sum_{y_i=y_j, y_i \neq y_k} \max\{\text{trace}(\mathbf{M} \mathbf{X}_{i,j}) - \text{trace}(\mathbf{M} \mathbf{X}_{i,k}) + \rho, 0\} \end{aligned} \quad (4.40)$$

To minimize Equation 4.40, the gradient descent method is used. The gradient respect to \mathbf{M} is computed as

$$\mathbf{G} = \frac{\partial f}{\partial \mathbf{M}} = (1 - \alpha) \sum_{y_i=y_j} \mathbf{X}_{i,j} + \alpha \sum_{y_i=y_j, y_i \neq y_k} (\mathbf{X}_{i,j} - \mathbf{X}_{i,k}) \quad (4.41)$$

The iteration process can be summarized as in Table 4.1. In each iteration, to make sure the updated M is a SPD matrix, first a eigenvalue decomposition is performed on M , and we have

$$M = V\Lambda V^T \quad (4.42)$$

here Λ is a diagonal matrix and its diagonal elements are eigenvalues. Then the negative eigenvalues in V are removed and the corresponding eigenvectors in V are also removed. Then M is restored by Equation (4.42).

TABLE 4.1: Optimization algorithm of Mahalanobis distance matrix learning

Gradient optimization algorithm for target function
Input Descriptors of training person pairs
Output A SPD matrix
Initialization
Initialize M_0 with eye matrix I ;
Compute the initial target function value f_0 with M_0 ;
Iteration count $t = 0$;
while (not converge)
Update $t = t + 1$;
Find x_k for all sample points x_i , where $y_i \neq y_k$;
Update gradient G_{t+1} with Equation 12;
Update M with equation : $M_{t+1} = M_t - \lambda G_t$;
Project M_{t+1} to the positive semi-positive definite space;
Update the target value $f _{M=M_{t+1}}$;
end while
return M

Chapter 5

Experiment Settings

5.1 Datasets and evaluation settings

VIPeR dataset is the most used dataset in person re-identification. In this dataset there are 632 different individuals and for each person there are two outdoor images from different viewpoints. All the images are scaled into 48×128 . In this experiment the we randomly select 316 individuals from cam A and cam B as the training set, the rest images in cam A are used as probe images and those in cam B as gallery images. This process is repeated 10 times to reduce error.

CUHK1 dataset contains 971 identities from two disjoint camera views. The cameras are static in each pair of view and images are listed in the same order. For each individual, there are two images in each view. All images are scaled into 60×160 . In this paper, we randomly select 485 image pairs as training data and the rest person pairs are used for test data.



FIGURE 5.1: Pedestrians in prid_450 dataset

Prid_2011 dataset consists of images extracted from multiple person trajectories recorded from two different, static surveillance cameras. Images from these cameras contain a viewpoint change and a stark difference in illumination, background and camera characteristics. Camera view A shows 385 persons, camera view B shows 749 persons. The first 200 persons appear in both camera views, The remaining persons in each camera view complete the gallery set of the corresponding view. Hence,

a typical evaluation consists of searching the 200 first persons of one camera view in all persons of the other view. This means that there are two possible evaluation procedures, either the probe set is drawn from view A and the gallery set is drawn from view B. In this paper, we randomly select 100 persons that appeared in both camera views as training pairs, and the remaining 100 persons of the 200 person pairs from camera A is used as probe set while the 649 remaining persons from camera B are used for gallery images.

Prid_450s dataset contains 450 image pairs recorded from two different, static surveillance cameras. Additionally, the dataset also provides an automatically generated, motion based foreground/background segmentation as well as a manual segmentation of parts of a person. The images are stored in two folders that represent the two camera views. Besides the original images, the folders also contain binary masks obtained from motion segmentation, and manually segmented masks. In this test, we randomly select 225 persons from each of two camera views as the training set, and the remaining persons are left as gallery and probe images.



FIGURE 5.2: Pedestrians in prid_450s dataset

GRID There are two camera views in this dataset. Folder probe contains 250 probe images captured in one view (file names starts from 0001 to 0250). Folder gallery contains 250 true match images of the probes (file names starts from 0001 to 0250). Besides, in gallery folder there are a total of 775 additional images that do not belong to any of the probes (file name starts with 0000). These extra images should be treated as a fixed portion in the testing set during cross validation. In this paper, we randomly select 125 persons from those 250 persons appeared in both camera views as training pairs, and the remaining persons in probe folder is used as probe images while the remaining 125 persons and those 775 additional persons from gallery folder are used as gallery images. A brief summarization of test settings is in Table 5.1.



FIGURE 5.3: Pedestrians in GRID dataset

TABLE 5.1: Testing setting for different datasets

Dataset	training	probe	gallery	cam_a	cam_b
VIPeR	316	316	316	632	632
CUHK1	485	486	486	971	971
Prid_2011	100	100	649	385	749
Prid_450s	225	225	225	450	450
GRID	125	125	900	250	1025

5.2 The influence of mean removal and L_2 normalization

In [25], mean removal and L_2 normalization is found to improve performance by 5.1%. The reason for this is mean removal and normalization can reduce the impact of extremas of descriptors. When testing proposed metric learning, we find the mean removal can slightly improve performance. A comparison between performance of original descriptors and preprocessed descriptors is shown in Tables 5.2, 5.3, 5.4, 5.5, 5.6, all those datasets are tested by proposed metric. The original GOG means no mean removal and normalization. It shows that the mean removal and normalization has a slight improvement around 0.5% on the performance on all five datasets. Since preprocessing are required to test XQDA, the mean removal and normalization are applied on descriptors in this experiment.

TABLE 5.2: The influence of data preprocessing on VIPeR

	Rank(%)				
Terms	1	5	10	15	20
Original GOG	43.01	74.91	84.87	89.81	93.32
Preprocessed GOG _{rgb}	43.77	74.84	85.25	90.32	93.89
Original GOG _{fusion}	48.77	77.47	87.41	91.52	94.27
Preprocessed GOG _{fusion}	48.32	76.90	87.78	91.93	94.49

TABLE 5.3: The influence of data preprocessing on CUHK1

	Rank(%)				
Terms	1	5	10	15	20
Original GOG _{rgb}	56.11	83.77	90.10	92.65	94.28
Preprocessed GOG _{rgb}	55.91	84.24	90.41	93.15	94.67
Original GOG _{fusion}	57.10	84.65	90.35	92.88	94.65
Preprocessed GOG _{fusion}	56.67	84.49	90.51	93.31	94.84

TABLE 5.4: The influence of data preprocessing on prid_2011

	Rank(%)				
Terms	1	5	10	15	20
Original GOG _{rgb}	24.80	52.10	63.20	69.90	72.90
Preprocessed GOG _{rgb}	23.80	52.20	63.50	70.20	73.50
Original GOG _{fusion}	32.20	56.60	67.00	73.10	77.70
Preprocessed GOG _{fusion}	32.30	57.40	66.30	73.40	78.00

TABLE 5.5: The influence of data preprocessing on prid_450s

	Rank(%)				
Terms	1	5	10	15	20
Original GOG _{rgb}	60.93	84.31	91.29	94.00	96.18
Preprocessed GOG _{rgb}	60.71	84.53	91.29	94.13	96.27
Original GOG _{fusion}	63.07	86.67	92.53	95.20	96.98
Preprocessed GOG _{fusion}	62.80	86.58	92.36	95.29	96.89

TABLE 5.6: The influence of data preprocessing on GRID

Terms	Rank(%)				
	1	5	10	15	20
Original GOG _{rgb}	22.96	41.92	51.68	58.72	64.64
Preprocessed GOG _{rgb}	22.64	43.68	52.00	59.04	65.04
Original GOG _{fusion}	24.32	44.56	54.80	62.40	66.64
Preprocessed GOG _{fusion}	23.92	44.64	54.88	62.32	66.40

5.3 Parameters setting

In this experiment, there are a few parameters for the iteration computing including slack variable ρ , maximal iteration T , gradient step λ , the interclass and intraclass limitation factor α and the updating ratio β . Firstly the slack variable ρ is initialized as 1 to ensure the minimum interclass distance is 1 larger than intraclass distance at least. The step size of gradient updating λ is initialized as 0.01. When target value f increases, λ is scaled by a factor 0.5, and λ is scaled by 1.01 when target value f decreases. To judge if target value converges, the threshold β is defined as the ratio target value change versus previous target value, that is, $\beta = \frac{(f_{t+1}-f_t)}{f_t}$. According many experiment trials, when it satisfies $\beta = 10^{-5}$, the target value converges and the iteration is stopped. The maximal iteration times t is set to 100 since the target value f will converge in around 15 iterations. The last parameter for the iteration is α , to know the best value for α , we tried 11 different values ranges from 0 to 1 with a step of 0.1, and find that the rank 1 and rank 5 scores reach maxima at interval $[0.7, 0.8]$, as shown in 5.4 and 5.5. Then another ten trials are performed with alpha ranging from $[0.7, 0.8]$ with a step of 0.01. The best α value should have as large top rank scores as possible and at last we find that the optimal value for α is 0.76. A form of all parameters are shown in Table 5.7.

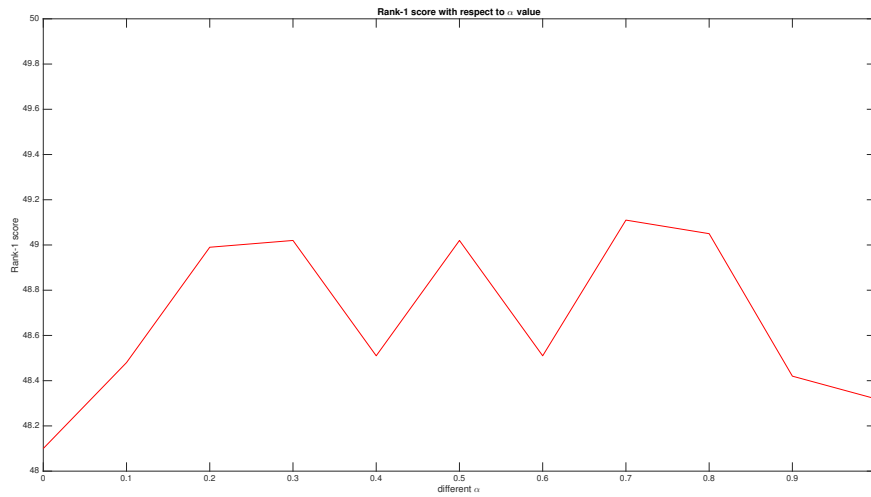
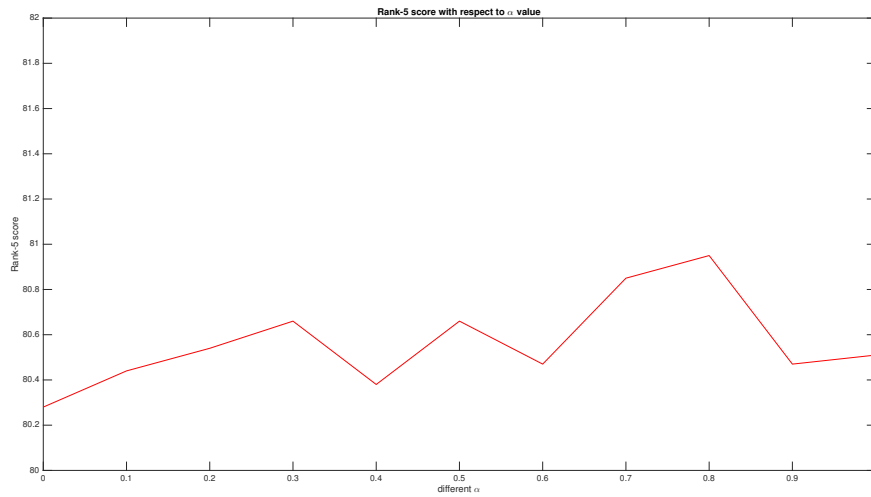
FIGURE 5.4: Rank 1 scores with respect to α on VIPeRFIGURE 5.5: Rank 5 scores with respect to α on VIPeR

TABLE 5.7: Parameters setting

Paramters	α	β	λ	t	ρ
Values	0.76	10^{-5}	0.01	100	1

Performance measuring The cumulative matching curve is used to measure the descriptor performance. The score means the probability that the right match is within the top n samples. A better CMC curve is expected to have a high rank-1 value and reaches 1 as fast as possible.

5.4 Results and experiment analysis

In this paper, we compare proposed metric with other state-of-the-art metrics including NFST [47], XQDA [22]. NFST is a metric which learn a null space for descriptors so that the the same class descriptors will be projected to a single point to minimize within class scatter matrix while different classes are projected to different points. This metric is a good solution to small sample problems in person re-identification. XQDA is quite similar with many other metrics, which learns a projection matrix W and then a Mahalanobis SPD matrix M is learned in the subspace. Those two metric are proved to have state-of-the-art performance with many other methods. The GOG_{rgb} in all forms stands for the hierarchical gaussian descriptor in RGB color space while GOG_{fusion} stands for the descriptor concatenated by four different color spaces {RGB, Lab, HSV, nRnG}.

VIPeR A comparison form is given in Table 3. Some of recent results are also included in this form. We can find that the rank scores are better than those of NFST and XQDA in terms of both GOG_{rgb} and GOG_{fusion} . More specifically, the rank 1, rank 5, rank 10, rank 15 and rank 20 scores of proposed metric learning are 0.76%, 0.92%, 1.39%, 1.08%, 1.52% higher than those of $GOG_{rgb} + XQDA$, and the rank 1, rank 5, rank 10, rank 15 and rank 20 GOG_{fusion} scores of proposed metric learning are 0.35%, -0.54%, 0.98%, 0.66%, 0.79% higher than $GOG_{fusion} + XQDA$ respectively. Also we can see that the proposed metric learning has a better performance than NFST.

TABLE 5.8: Performance of different metrics on VIPeR

	Rank(%)				
Methods	1	5	10	15	20
$GOG_{rgb}+NFST$	43.23	73.16	83.64	89.59	92.88
$GOG_{rgb}+XQDA$	43.01	73.92	83.86	89.24	92.37
$GOG_{rgb}+Proposed$	43.77	74.84	85.25	90.32	93.89
$GOG_{fusion}+NFST$	47.15	76.39	87.31	91.74	94.49
$GOG_{fusion}+XQDA$	47.97	77.44	86.80	91.27	93.70
$GOG_{fusion}+Proposed$	48.32	76.90	87.78	91.93	94.49

CUHK1 We can find that the rank 1, rank5, rank 10, rank 15, rank 20 score

of GOG_{rgb} combined with proposed metric are 5.4%, 4.18%, 3.31%, 2.16%, 1.46% higher than XQDA, and 0.31%, 1.22%, 1.34%, 1.17%, 1.11% than NFST. Also the rank 1, rank 5, rank 10, rank 15, rank 20 score of GOG_{fusion} combined with proposed metric are 4.57%, 2.64%, 0.70%, 1.33%, 0.83% higher than GOG_{fusion} combined with XQDA, and 0.41%, 0.83%, 0.88%, 1.09%, 1.14% than GOG_{fusion} combined with NFST.

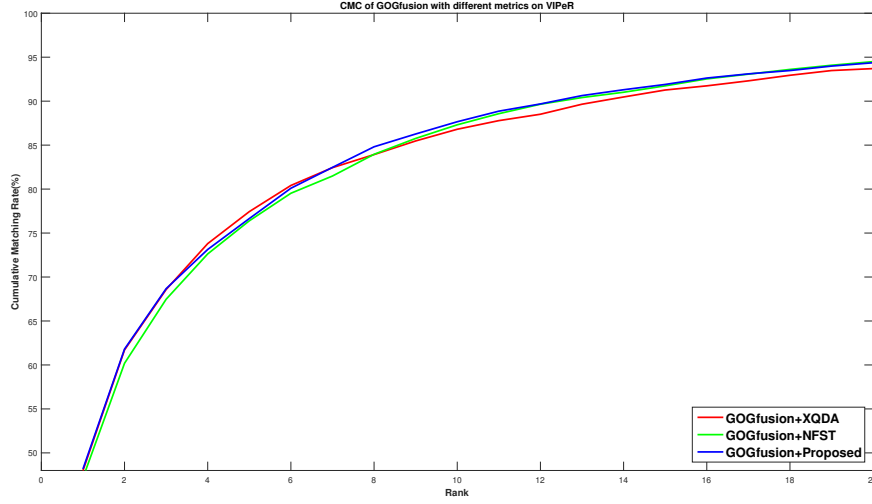


FIGURE 5.6: CMC curves on VIPeR comparing different metric learning

TABLE 5.9: Performance of different metrics on CUHK1

Methods	Rank(%)				
	1	5	10	15	20
$GOG_{rgb}+NFST$	55.60	83.02	89.07	91.98	93.56
$GOG_{rgb}+XQDA$	50.51	80.06	87.10	90.99	93.21
$GOG_{rgb}+Proposed$	55.91	84.24	90.41	93.15	94.67
$GOG_{fusion}+NFST$	56.26	83.66	89.63	92.22	93.70
$GOG_{fusion}+XQDA$	52.10	81.85	88.81	91.98	94.01
$GOG_{fusion}+Proposed$	56.67	84.49	90.51	93.31	94.84

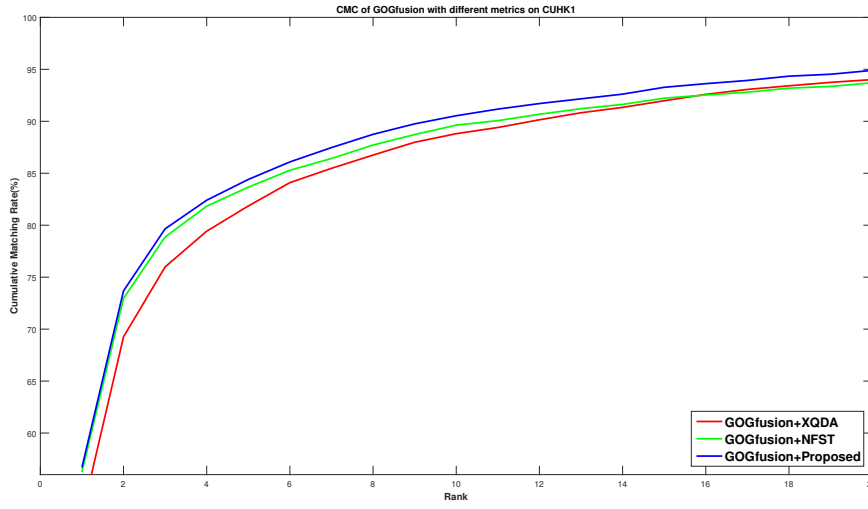


FIGURE 5.7: CMC curves on CUHK1 comparing different metric learning

TABLE 5.10: Performance of different metrics on prid_2011

Methods	Rank(%)				
	1	5	10	15	20
GOG _{rgb} +NFST	26.60	53.80	62.90	71.30	75.40
GOG _{rgb} +XQDA	31.10	55.70	66.10	72.40	76.10
GOG _{rgb} +Proposed	23.80	52.20	63.50	70.20	73.50
GOG _{fusion} +NFST	34.10	58.30	67.60	73.80	78.30
GOG _{fusion} +XQDA	38.40	61.30	70.80	75.60	79.30
GOG _{fusion} +Proposed	32.30	57.40	66.30	73.40	78.00

Prid_2011 The rank 1, rank5, rank 10, rank 15, rank 20 score of GOG_{fusion} combined with proposed metric are 6.1%, 3.9%, 4.5%, 2.2% and 1.3% lower than GOG_{fusion} combined with XQDA. The performance of NFST is slightly better than proposed metric. Also in terms of GOG_{rgb} XQDA and NFST has better performance than the proposed one. So in this dataset the proposed metric has worse performance than XQDA and NFST.

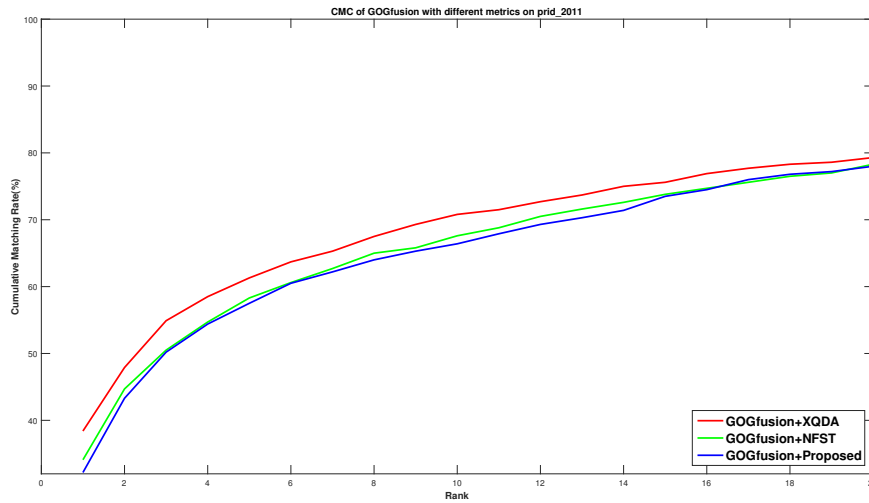


FIGURE 5.8: CMC curves on prid_2011 comparing different metric learning

TABLE 5.11: Performance of different metrics on prid_450s

Methods	Rank(%)				
	1	5	10	15	20
GOG _{rgb} +NFST	61.96	84.98	90.53	94.09	96.09
GOG _{rgb} +XQDA	65.29	85.02	91.13	94.76	96.49
GOG _{rgb} +Proposed	60.71	84.53	91.29	94.13	96.27
GOG _{fusion} +NFST	64.53	86.62	92.93	95.78	97.42
GOG _{fusion} +XQDA	68.40	87.42	93.47	95.69	97.02
GOG _{fusion} +Proposed	62.80	86.58	92.36	95.29	96.89

Prid_450s In this dataset, we can find the rank 1 score of XQDA and NFST is higher than proposed metric, but they have almost the same rank 5, rank 10, rank 15, and rank 20 scores with respect to both kinds of descriptors.

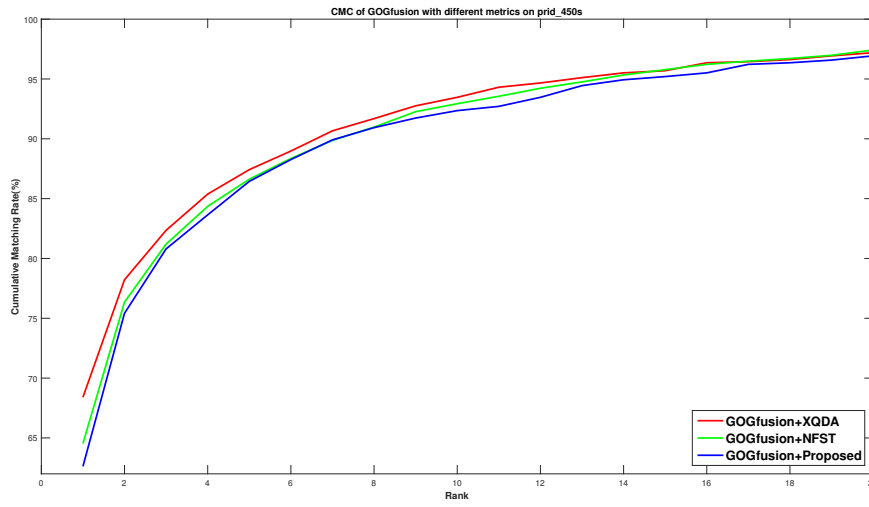


FIGURE 5.9: CMC curves on prid_450s comparing different metric learning

TABLE 5.12: Performance of different metrics on GRID

Methods	Rank(%)				
	1	5	10	15	20
GOG _{rgb} +NFST	21.84	41.28	50.96	57.44	62.88
GOG _{rgb} +XQDA	22.64	43.92	55.12	61.12	66.56
GOG _{rgb} +Proposed	22.64	43.68	52.00	59.04	65.04
GOG _{fusion} +NFST	23.04	44.40	54.40	61.84	66.56
GOG _{fusion} +XQDA	23.68	47.28	58.40	65.84	69.68
GOG _{fusion} +Proposed	23.92	44.64	54.88	62.32	66.40

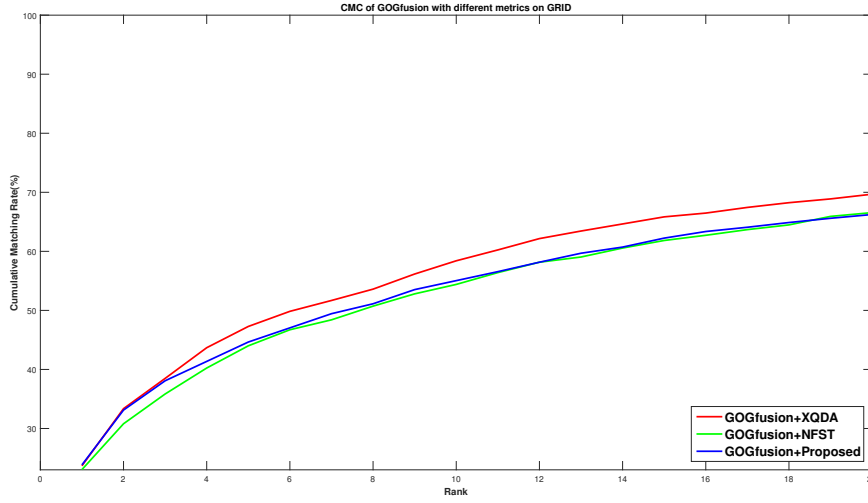


FIGURE 5.10: CMC curves on GRID comparing different metric learning

GRID We can see that the rank 1 score of proposed metric are 0.24% higher than XQDA and 0.88% higher than NFST in terms of GOG_{fusion} , but XQDA outperforms proposed metric on rank 5, rank 10, rank 15 and rank 20 scores. Besides, proposed metric outperforms NFST on rank 5, rank 10, rank 15 scores.

In summary, the Re-ID performance is improved in VIPeR, CUHK01 dataset, and has almost the same performance with NFST and XQDA on prid_450s dataset. Specifically, proposed metric learning has the best rank 1 score in GRID dataset and its performance is only second to XQDA. The proposed metric has superior performance for following reasons: (1) dimension reduction by KLFDA exploits the nonlinearity and the loss of discriminant information between classes are minimized. (2) the simplified relative distance limitation optimization helps to confine the Mahalanobis distance matrix M to discriminate different classes.

Chapter 6

Conclusion

In this thesis KLFDA is used to reduce dimension of hierarchical gaussian descriptors, gradient descent method is used to learn a Mahalanobis distance matrix on lower-dimensional space. By comparison we can find the proposed metric has better performance than NFST and XQDA on VIPeR and CUHK1 datasets, but XQDA and NFST outperforms the proposed metric learning on Prid_2011 and Prid_450s, and the proposed metric learning has better rank 1 score than NFST and its performance is only second to XQDA on GRID dataset.

6.1 Contributions

There are three contributions in this thesis: (1) The metric learning on dimension-reduced hierarchical gaussian descriptor by KLFDA has been studied. The gradient descent method optimizes the Mahalanobis distance matrix on the lower-dimensional space. It has been demonstrated that extra improvements can be achieved in this lower-dimensional space. (2) The influence of background and foreground segmentation on different descriptors have been fully studied. The results are that foreground segmentation improves performance of color based descriptors but decreases performance of texture based descriptors, which is caused by imperfect segmentation on single image based foreground segmentation. (3) Some variants of hierarchical gaussian descriptor have been tested, LBP and superpixel segmentation are combined with hierarchical gaussian descriptor but those variants have worse performance than original hierarchical gaussian descriptor.

6.2 Future work

6.2.1 Improve Hierarchical Gaussian descriptor

It has been demonstrated that in the basic pixel feature f_i , the color components are more important than y coordinate and gradient components. LBP has been turned to be a worse choice for texture representation in f_i . Therefore, one strategy to improve hierarchical gaussian descriptor is to find a better texture representation to replace gradient components in basic pixel feature.

6.2.2 Influence of video-based foreground segmentation

Though the single-image-based foreground and background segmentation's influence on different descriptors has been studied, extra effort is needed to study sequence or video-based segmentation's influence on different descriptors. Video-based foreground often has better results. If the background can be well modelled by a video or a sequence of images so that less extra textural noise is created, the segmentation might improve hierarchical gaussian descriptor's performance.

6.2.3 Computational cost of gradient descent method

The last one is to reduce the computational cost of gradient descent method. It takes about average 15 iterations when training the Mahanalobis distance matrix. The training time for larger dataset like CUHK dataset takes up to one hour. Therefore, other variants of gradient descent method like stochastic gradient method, conjugate gradient method may be tested for lower computational cost.

Bibliography

- [1] Radhakrishna Achanta et al. “SLIC superpixels compared to state-of-the-art superpixel methods”. In: *IEEE transactions on pattern analysis and machine intelligence* 34.11 (2012), pp. 2274–2282.
- [2] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. “Pictorial Structures Revisited: People Detection and Articulated Pose Estimation”. In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE. Apr. 2009, pp. 1–8.
- [3] Davide Baltieri, Roberto Vezzani, and Rita Cucchiara. “SARC3D: a new 3D body model for People Tracking and Re-identification”. In: *International Conference on Image Analysis and Processing*. Springer. Mar. 2011, pp. 1–10.
- [4] Oren Barkan et al. “Fast high dimensional vector multiplication face recognition”. In: *Proceedings of the IEEE International Conference on Computer Vision*. 2013, pp. 1960–1967.
- [5] Loris Bazzani et al. “Multiple-shot person re-identification by chromatic and epitomic analyses”. In: *Pattern Recognition Letters* 33.7 (2012), pp. 898–903.
- [6] Bedagkar-Gala et al. “A Survey of Approaches and Trends in Person Re-identification”. In: *Image and Vision Computing* 32.4 (Mar. 2014), pp. 270–286.
- [7] A Bedagkar-Gala and Shishir K Shah. “Part-based spatio-temporal model for multi-person re-identification”. In: *Pattern Recognition Letters* 33.14 (Oct. 2012), pp. 1908–1915.
- [8] Apurva Bedagkar-Gala and Shishir K Shah. “Multiple Person Re-identification using Part based Spatio-Temporal Color Appearance Model”. In: *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*. IEEE. Oct. 2011, pp. 1–8.

- [9] Dapeng Chen et al. "Similarity Learning on an Explicit Polynomial Kernel Feature Map for Person Re-Identification". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Apr. 2015, pp. 1–9.
- [10] De Cheng et al. "Person Re-Identification by Multi-Channel Parts-Based CNN With Improved Triplet Loss Function". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. July 2016, pp. 1–10.
- [11] Shaogang Gong Chen Change Loy Chunxiao Liu and Xinggang Lin. "Person Re-identification: What Features Are Important?" In: *European Conference on Computer Vision*. IEEE. Nov. 2015, pp. 1–11.
- [12] Navneet Dalal and Bill Triggs. "Histograms of Oriented Gradients for Human Detection". In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Jan. 2005, pp. 886–893.
- [13] Jason V Davis et al. "Information-Theoretic Metric Learning". In: *Proceedings of the 24th international conference on Machine learning*. ACM. June 2007, pp. 209–216.
- [14] Michael Jones Ejaz Ahmed and Tim K Marks. "An Improved Deep Learning Architecture for Person Re-Identification". In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 1–9.
- [15] M. Farenzena et al. "Person Re-Identification by Symmetry-Driven Accumulation of Local Features". In: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE. Mar. 2016, pp. 1–8.
- [16] Pedro F Felzenszwalb and Daniel P Huttenlocher. "Pictorial Structures for Object Recognition". In: *International Journal of Computer Vision* 61.1 (Jan. 2005), pp. 1–42.
- [17] Xiaofei He and Partha Niyogi. "Locality Preserving Projections". In: *NIPS*. Nov. 2003, pp. 1–8.
- [18] Nebojsa Jojic et al. "Stel component analysis: Modeling spatial correlations in image class structure". In: *Computer Vision and Pattern Recognition*. Apr. 2009, pp. 1–8.
- [19] Arif Khan, Jian Zhang, and Yang Wang. "Appearance-Based Re-identification of People in Video". In: *2010 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*. IEEE, Nov. 2010, pp. 357–362.

- [20] Martin Kostinger et al. “Large Scale Metric Learning from Equivalence Constraints”. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE. Apr. 2012, pp. 2288–2295.
- [21] Peihua Li, Qilong Wang, and Lei Zhang. “A Novel Earth Mover’s Distance Methodology for Image Matching with Gaussian Mixture Models”. In: *Proceedings of the IEEE International Conference on Computer Vision*. IEEE, Dec. 2013, pp. 1689–1696.
- [22] Shengcai Liao et al. “Person Re-identification by Local Maximal Occurrence Representation and Metric Learning”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Apr. 2015, pp. 1–10.
- [23] Miroslav Lovric, Maung Min-Oo, and Ernst A Ruh. “Multivariate Normal Distributions Parametrized as a Riemannian Symmetric Space”. In: *Journal of Multivariate Analysis* 74.1 (June 2000), pp. 1–15.
- [24] David G Lowe. “Object Recognition from Local Scale-Invariant Features”. In: *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*. IEEE. June 1999, pp. 1–8.
- [25] Tetsu Matsukawa et al. “Hierarchical Gaussian Descriptor for Person Re-Identification”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Dec. 2016, pp. 1363–1372.
- [26] Niall McLaughlin, Jesus Martinez del Rincon, and Paul Miller. “Recurrent Convolutional Network for Video-Based Person Re-Identification”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. July 2016, pp. 1–10.
- [27] Alexis Mignon and Frederic Jurie. “PCCA: A New Approach for Distance Learning from Sparse Pairwise Constraints”. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE. Apr. 2012, pp. 1–7.
- [28] Baback Moghaddam, Tony Jebara, and Alex Pentland. “Bayesian Face Recognition”. In: *Pattern Recognition* 33.11 (Jan. 2000), pp. 1771–1782.
- [29] T. Ojala, M. Pietikäinen, and D. Harwood. “A Comparative Study of Texture Measures with Classification Based on Feature Distributions”. In: *Pattern recognition* (1996), pp. 51–59.

- [30] T. Ojala, M. Pietikäinen, and D. Harwood. "Performance evaluation of texture measures with classification based on Kullback discrimination of distributions". In: *ICPR*. IEEE. 1994, pp. 582–585.
- [31] Takumi Kobayashi Otsu and Nobuyuki. "Image Feature Extraction Using Gradient Local Auto-Correlations". In: *European conference on computer vision*. Aug. 2008, pp. 1–13.
- [32] Sakrapee Paisitkriangkrai, Chunhua Shen, and Anton van den Hengel. "Learning to Rank in Person Re-Identification With Metric Ensembles". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Jan. 2015, pp. 1–10.
- [33] Sateesh Pedagadi et al. "Local Fisher Discriminant Analysis for Pedestrian Re-identification". In: *2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jan. 2013, pp. 3318–3325.
- [34] Bryan Prosser et al. "Person Re-Identification by Support Vector Ranking". In: *British Machine Vision Conference 2010*. British Machine Vision Association, July 2010, pp. 21.1–21.11.
- [35] Peter M Roth et al. "Mahalanobis Distance Learning for Person Re-Identification". In: *Person Re-Identification*. Springer, July 2014, pp. 1–21.
- [36] Riccardo Satta. "Appearance Descriptors for Person Re-identification: a Comprehensive Review". In: *Computing Research Repository* (July 2013), pp. 1–18.
- [37] Zhiyuan Shi, Timothy M Hospedales, and Tao Xiang. "Transferring a Semantic Representation for Person Re-Identification and Search". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Apr. 2015, pp. 1–10.
- [38] Masashi Sugiyama. "Local Fisher Discriminant Analysis for Supervised Dimensionality Reduction". In: *Proceedings of the 23rd international conference on Machine learning*. ACM. Dec. 2006.
- [39] Jirí Trefný and Jirí Matas. "Extended set of local binary patterns for rapid object detection". In: *Computer Vision Winter Workshop*. 2010, pp. 1–7.

- [40] Oncel Tuzel, Fatih Porikli, and Peter Meer. “Region Covariance: A Fast Descriptor for Detection and Classification”. In: *European conference on computer vision*. Springer. Dec. 2016, pp. 1–14.
- [41] Weinberger et al. “Distance Metric Learning for Large Margin Nearest Neighbor Classification”. In: *Journal of Machine Learning Research* 10 (Feb. 2009), pp. 207–244.
- [42] Max Welling. “Kernel Canonical Correlation Analysis”. In: (Mar. 2005), pp. 1–3.
- [43] Tong Xiao et al. “Learning Deep Feature Representations with Domain Guided Dropout for Person Re-identification”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Apr. 2016, pp. 1–10.
- [44] Fei Xiong et al. “Person Re-Identification using Kernel-based Metric Learning Methods”. In: *European conference on computer vision*. Springer. July 2014, pp. 1–16.
- [45] Yang Yang et al. “Salient color names for person re-identification”. In: *European Conference on Computer Vision*. 2014.
- [46] Jinjie You et al. “Top-push Video-based Person Re-identification”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Jan. 2016, pp. 1–9.
- [47] Li Zhang, Tao Xiang, and Shaogang Gong. “Learning a Discriminative Null Space for Person Re-identification”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Mar. 2016, pp. 1239–1248.
- [48] Guoying Zhao and Matti Pietikainen. “Dynamic texture recognition using local binary patterns with an application to facial expressions”. In: *IEEE transactions on pattern analysis and machine intelligence* 29.6 (2007).
- [49] Rui Zhao, Wanli Ouyang, and Xiaogang Wang. “Learning Mid-level Filters for Person Re-identification”. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, May 2014, pp. 144–151.
- [50] WeiShi Zheng, Shaogang Gong, and Tao Xiang. “Person Re-identification by Probabilistic Relative Distance Comparison”. In: *Computer vision and pattern recognition (CVPR), 2011 IEEE conference on*. IEEE. Nov. 2016, pp. 1–8.