

One-pass Person Re-identification by Sketch Online Discriminant Analysis

Wei-Hong Li, Zhuowei Zhong, and Wei-Shi Zheng*

Abstract—Person re-identification (re-id) is to match people across disjoint camera views in a multi-camera system, and re-id has been an important technology applied in smart city in recent years. However, the majority of existing person re-id methods are not designed for processing sequential data in an online way. This ignores the real-world scenario that person images detected from multi-cameras system are coming sequentially. While there is a few work on discussing online re-id, most of them require considerable storage of all passed data samples that have been ever observed, and this could be unrealistic for processing data from a large camera network. In this work, we present an one-pass person re-id model that adapts the re-id model based on each newly observed data and no passed data are directly used for each update. More specifically, we develop an Sketch online Discriminant Analysis (SoDA) by embedding sketch processing into Fisher discriminant analysis (FDA). SoDA can efficiently keep the main data variations of all passed samples in a low rank matrix when processing sequential data samples, and estimate the approximate within-class variance (i.e. within-class covariance matrix) from the sketch data information. We provide theoretical analysis on the effect of the estimated approximate within-class covariance matrix. In particular, we derive upper and lower bounds on the Fisher discriminant score (i.e. the quotient between between-class variation and within-class variation after feature transformation) in order to investigate how the optimal feature transformation learned by SoDA sequentially approximates the offline FDA that is learned on all observed data. Extensive experimental results have shown the effectiveness of our SoDA and empirically support our theoretical analysis.

Index Terms—Online learning, Person re-identification, Discriminant feature extraction

I. INTRODUCTION

Person re-identification (re-id) [51], [1], [13], [22], [31], [20], [54] is crucially important for successfully tracking people in a large camera network. It is to match the same person's images captured at non-overlapping camera views at different time. Person re-id by visual matching is inherently challenging because of the existence of many visually similar persons and dramatic appearance changes of the same person caused by the serious cross-camera-view variations such as illumination, viewpoint, occlusions and background clutter. Recently, a large number of works [22], [23], [3], [16], [27], [30], [35], [44], [53] have been reported to solve this challenge.

Wei-Hong Li is with the School of Electronics and Information Technology, Sun Yat-sen University, Guangzhou, China. E-mail: li-weih3@mail2.sysu.edu.cn

Zhuowei Zhong is with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China. E-mail: zhongzhw6@gmail.com

Wei-Shi Zheng is with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China. E-mail: wszheng@ieee.org/zhwshi@mail.sysu.edu.cn

* Corresponding author.

However, it is largely unsolved to perform online learning for person re-identification, since most person re-id models except [25], [39], [29], [37] are only suitable for offline learning. On one hand, the offline learning mode cannot enable a real-time update of person re-id model when a large amount of persons are detected in a camera network. An online update is important to keep the cross-view matching system work on recent mostly interested persons, that is to make the whole re-id system work on sequential data. On the other hand, online learning is helpful to alleviate the large scale learning problem (either with high-dimensional feature, or on large-scale data set, or both) nowadays. By using online learning, especially the one-pass online learning, it is not necessary to always store (all) observed/passed data samples.

In this paper, we overcome the limitation of offline person re-id methods by developing an effective online person re-id model. We proposed to embed the sketch processing into Fisher discriminant analysis (FDA), and the new model is called Sketch online Discriminant Analysis (SoDA). In SoDA, the sketch processing preserves the main variations of all passed data samples in a low-rank sketch matrix, and thus SoDA enables selecting data variation for acquiring discriminant features during online learning. SoDA enables the newly learned discriminant model to embrace information from a new coming data sample in the current round and meanwhile retain important information learned in previous rounds in a light and fast manner without directly saving any passed observed data samples and keeping large-scale covariance matrices, so that SoDA is formed as an one-pass online adaptation model. While no passed data samples are saved in SoDA, we propose to estimate the within-class variation from the sketch information (i.e. a low-rank *sketch matrix*), and thus in SoDA an approximate within-class covariance matrix can be derived. We have provided in-depth theoretical analysis on how sketch affects the discriminant feature extraction in an online way. The rigorous upper and lower bounds on how SoDA approaches its offline model (i.e. the classical Fisher Discriminant Analysis [41]) are presented and proved.

Compared to existing online models for person re-id [25], [39], [29], [37], SoDA is succinct, but it is theoretically guaranteed and effective. While most existing online re-id models have to retain all observed passed data samples, the proposed SoDA relies on the sketch information from historical data without any explicit storage of passed data samples, and sketch information will assist our online model in preventing one-pass online model from being biased by a new coming data. While a more conventional way for online learning of FDA is to update both within-class and between-

class covariance matrices directly [33], [48], [38], [26], [34], [15], we introduce a novel approach to realize online FDA by mining any within-class information from a sketch data matrix, and this provides a lighter, more efficient and effective online learning for FDA. We also find that an extra benefit of embedding sketch processing in SoDA is to simultaneously embed dimension reduction as well, so that no extra learning task on learning dimension reduction technology (e.g. PCA) is required and SoDA is more flexible when learning on some high dimensional data [22], [5] in an online manner.

We have conducted extensive experiments on three largest scale person re-identification datasets in order to evaluate the effectiveness of SoDA for learning person re-identification model in an online way. Extensive experiments are also included for comparing SoDA with related online learning models, even though they were not applied to person re-identification before.

The rest of the paper is organized as follows. In Sec. II, the related literatures are first reviewed. We elaborate our online algorithm and analyze the space and time complexity of SoDA in Sec. III. Then we present theoretical analysis on the relationship between our SoDA and the offline FDA in Sec. IV. Experimental results for evaluation and verification of our theoretical analysis are reported in Sec. V and finally we conclude the work in Sec. VI.

II. RELATED WORK

Online Person re-identification. While person re-identification has been investigated in a large number of works [51], [1], [13], [31], [20], [54], [22], [23], [3], [16], [27], [30], [35], [44], [53], [32], [28], [47], the majority of them only address by offline learning. That is person re-id model is learned on a fixed training dataset. This ignores the increase demand of data from a visual surveillance system, since thousands of person images are captured day by day and it is demanded to train a person re-id system on streaming data so as to keep the system update to date.

Recently, only a few works [37], [25], [39], [29] have been developed towards online processing for person re-identification. The most related work is the incremental distance metric based online learning mechanism (OL-IDM) proposed in [37]. For updating the KISSME metric [17], the OL-IDM utilizes the modified Self-Organizing Incremental Neural Network (SOINN) [8] to produce two pairwise sets: a similar pairs set and a dissimilar pairs set. Although SOINN enables learning KISSME [17] on sequential data, SOINN has to compare the newly observed sample with all the preserved nodes and adds the newly observed sample as a new node if it does not appear in the network. This would be costly as sequential data increase and when feature dimension is high.

Another related work is the human-in-the-loop ones [39], [25], [29], which proposed incremental method learned with the involvement of humans' feedback. Wang et al. [39] assumes that an operator is available to scan the rank list provided by the proposed algorithm when matching a new probe sample with existing observed gallery ones, and this operator will select the true match, strong-negative match,

and weak-negative match for the probe. After having the human feedback, the algorithm is able to be update. Martinel et al. presented a graph-based approach to exploit the most informative probe-gallery pairs for reducing human efforts and developed an incremental and iterative approach based on the feedback [29].

Unlike these models, we design a sketch FDA model called SoDA for one-pass online learning, without any storage of passed observed samples, maintaining a small size sketch matrix on handling streaming data so that the discriminant projections can be updated efficiently for extracting discriminative features for identifying different individuals.

Thanks to the sketch matrix, our SoDA is capable of obtaining comparable performance with offline FDA models on streaming data or large and high dimensional datasets with very low cost on space and time. Compared to the related online person re-id models, SoDA is theoretically sounded since the bounds on approximating the offline model is provided.

In particular, compared to Wang et al.'s and Martinel et al.'s work, our work has the following distinct aspects: Firstly, the proposed SoDA is developed for the one-pass online learning, while Wang et al.'s and Martinel et al.'s work cannot work for one-pass online learning, because the former one requires human feedback between probe sample and all preserved gallery samples, and the latter one needs to store all sample pairs during interactive learning. Secondly, the proposed SoDA could be orthogonal to the human-in-the-loop work, since we discuss how to automatically update a person re-identification model on streaming data without elaborated human interaction (feedback), and thus our work and the idea of incorporating more human interaction in human-in-the-loop work can accompany each other.

SoDA vs. Incremental Fisher Discriminant Learning. SoDA is related to existing incremental/online Fisher Discriminant Analysis (FDA) methods, which aim to update within-class and between-class covariance matrix sequentially. Pang et al. proposed to directly update the between-class and within-class scatter matrices [33]. However, Pang et al.'s method has to preserve the whole scatter matrices in the memory, which becomes impractical for high dimensional data. Ye et al. [48] and Uray et al. [38] performed online learning by updating PCA components to derive an approximate update of scatter matrices. Compared to Pang's method, Ye's and Uray's can only perform online learning sample by sample, which can be time consuming for large scale data. Also, Ye's method is based on QR decomposition of between-class covariance matrix, and therefore it would increase computational cost when the number of class is large. Since, Ye's method is limited to learning discriminant projections in the range space of between-class covariance matrix but not the range space of total-class covariance matrix [46], which may lose discriminant information. Lu et al. proposed a complete model that picks up the lost discriminant information [26]. But Lu's method only can update the model sample by sample. Peng et al. alternately proposed a chunk version of Ye's method in order to process multiple data points at a time [34]. Kim et

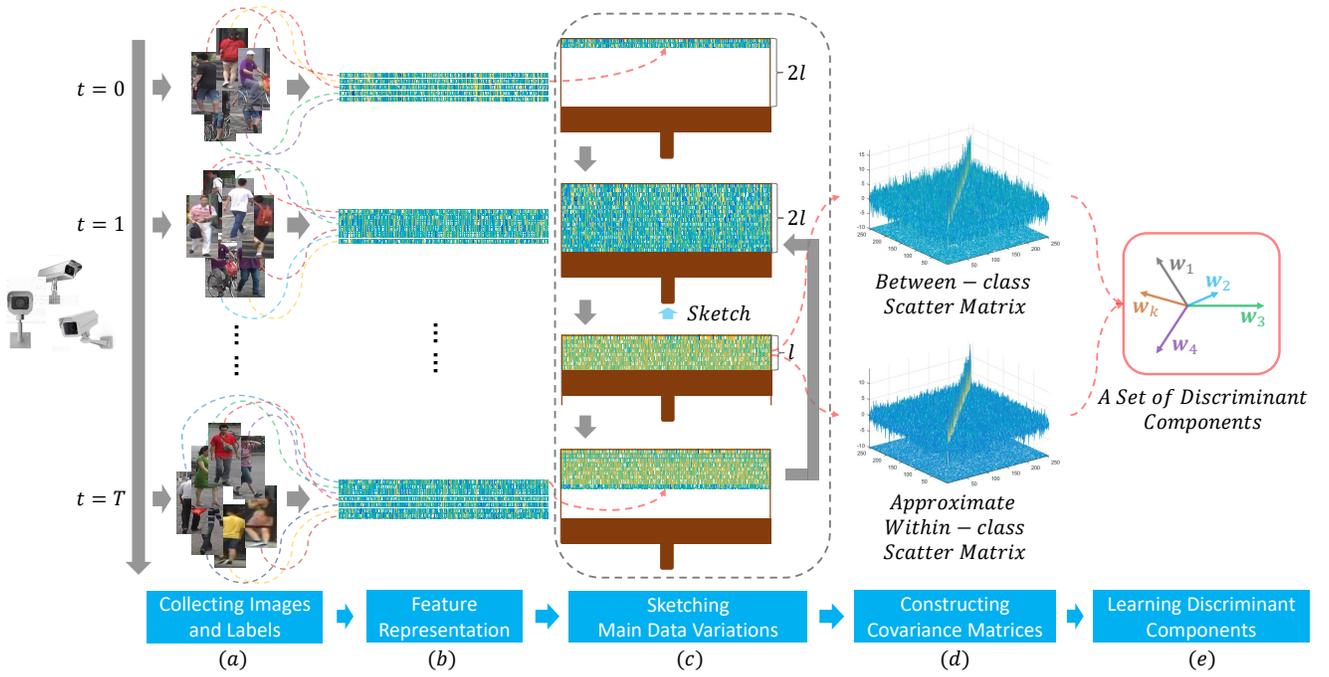


Fig. 1. Illustration of our proposed Sketch online Discriminant Analysis (SoDA) (Best viewed in color). (a) In real-world application, images are generated endlessly from visual surveillance camera network. (b) ($t = 0, 1, \dots, T$), every presented image is represented by a d -dimensional row feature vector. (c) We maintain a low rank sketch matrix to summarize all passed data by matrix sketch: 1) At the beginning, we set $\mathbf{B} \in \mathcal{R}^{2\ell \times d}$, the sketch matrix, to be a zero matrix. 2) All rows of \mathbf{B} would be filled by 2ℓ samples from top to bottom one by one. 3) we maintain the main data variations in the upper half of \mathbf{B} by sketch. 4) Each row of the lower half of \mathbf{B} is set to be all zero and will be replaced by a new sample. (d) After sketch, the between-class and within-class covariance matrices are constructed. (e) Due to the sketch, we can update a set of discriminant components efficiently only using limited space and time.

al. proposed a sufficient spanning set based incremental FDA [15] to overcome the limitations in the previous works. Since it is hard to directly update the discriminant components in FDA, Yan et al. [45] and Hiraoka et al. [10] modified FDA in order to get the discriminant components updated. They proposed iterative methods for directly updating discriminant projections.

Compared to the above mentioned incremental/online FDA methods, our proposed SoDA embeds sketch processing into FDA and therefore mines the within-class scatter information from a sketch data matrix rather than directly from samples. This gives the benefit that while the passed data samples are not necessary to be saved, SoDA is still able to extract useful within-class information from the compressed data information contained in the sketch matrix. In general, SoDA is an online version of FDA, and SoDA can not only approximate the FDA, which optimizes discriminant components on whole data directly, but also run faster with limited memory. Also, dimension reduction is naturally embedded into SoDA and no extra online model for dimension reduction is required. In-depth theoretical investigation is provided in Sec. IV to explain its rationale and to guarantee its effectiveness.

Although the proposed SoDA can be seen as embedding sketch processing into FDA, we contribute solid theoretical analysis on how SoDA will approximate the Batch mode FDA when estimating the within-class variations from sketch information, where the lower bound and upper bound are provided. The theoretical analysis guarantees SoDA to be an effective and efficient online learning method.

Online Learning. SoDA is an online learning methods. In literatures, online learning [2], [6], [12], [40], [11] is known as a light and rapid means to process streaming data or large-scale datasets, and it has been widely exploited in many real-world tasks such as Face Recognition [14], [36], Images Retrieval [21], [42] and Object Tracking [19], [18]. It enables learning a up-to-date model based on streaming data. However, most of these online leaning based models [6], [18], [19] are not suitable for person re-identification, since they are incapable of predicting labels of data samples from unseen classes which do not appear in the training stage.

III. SKETCH ONLINE DISCRIMINANT ANALYSIS (SoDA)

In this section, we start to present the Sketch online Discriminant Analysis (SoDA) for Person re-identification. In real-world scenario, samples come endlessly and sequentially from vision system (Figure 1). The number of samples received in each round is random, and the individual sample obtained is also stochastic. Suppose the t^{th} ($t = 1, 2, \dots$) new coming sample represented as a d -dimensional feature vector $\mathbf{x}_i \in \mathcal{R}^d$ is labelled with class label \mathbf{y}_i . For convenience, at the t^{th} round, we denote all passed data (i.e. N training samples collected in the current and previous rounds) as a training sample matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathcal{R}^{N \times d}$, and denote all the corresponding labels as $\mathbf{y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N]^T \in \mathcal{R}^N$ where \mathbf{y}_i is the class label of \mathbf{x}_i and $\mathbf{y}_i \in \{1, 2, \dots, C\}$.

At each round ($t = 1, 2, \dots$), the proposed SoDA maintains the main variations of all passed data ($\mathbf{X} \in \mathcal{R}^{N \times d}$) in a low rank matrix, which is named as the “sketch matrix”.

Algorithm 1: Sketch online Discriminant Analysis

Input: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathcal{R}^{N \times d}, \mathbf{y} \in \mathcal{R}^N, \lambda > 0$
 1 $\mathbf{B} \leftarrow$ zero matrix $\in \mathcal{R}^{2\ell \times d}$;
 2 **for** each data $\mathbf{x}_i \in \mathcal{R}^d$ and label \mathbf{y}_i **do**
 3 using \mathbf{x}_i^T to replace one zero row of \mathbf{B} ;
 4 **if** all samples in \mathbf{X} are processed **then**
 5 deleting all zero rows of \mathbf{B} ;
 6 **end**
 7 **if** \mathbf{B} has no zero rows **then**
 8 $[\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}] = \text{SVD}(\mathbf{B})$;
 9 setting ξ as the $(\ell + 1)^{\text{th}}$ largest element $\Sigma_{\ell+1}$ of $\mathbf{\Sigma}$;
 10 $\hat{\mathbf{\Sigma}} = \sqrt{\max(\mathbf{\Sigma}^2 - \mathbf{I}_{2\ell}\xi^2, \mathbf{O})}$;
 11 $\mathbf{B} = \hat{\mathbf{\Sigma}}\mathbf{V}^T$ (\mathbf{B} contains ℓ rows non-zero values);
 12 **end**
 13 $\mathbf{m}_c \leftarrow (N_c \mathbf{m}_c + \mathbf{x}_i) / (N_c + 1)$ ($c = 0, \mathbf{y}_i$);
 14 $N_c \leftarrow N_c + 1$ ($c = 0, \mathbf{y}_i$);
 15 **end**
 16 $\mathbf{B} \leftarrow \mathbf{B}^+, \mathbf{P} = \mathbf{V}^+$;
 17 $\mathbf{S}_b = \sum_{c=1}^C \frac{N_c}{N_0} (\mathbf{m}_c - \mathbf{m}_0)(\mathbf{m}_c - \mathbf{m}_0)^T$;
 18 $\tilde{\mathbf{S}}_t = \mathbf{B}^T \mathbf{B} / N_0 - \mathbf{m}_0 \mathbf{m}_0^T$;
 19 $\tilde{\mathbf{S}}_w = \tilde{\mathbf{S}}_t - \mathbf{S}_b$;
 20 $\hat{\mathbf{S}}_b = \mathbf{P}^T \mathbf{S}_b \mathbf{P}$;
 21 $\hat{\mathbf{S}}_w = \mathbf{P}^T \tilde{\mathbf{S}}_w \mathbf{P}$;
 22 $[\mathbf{W}, \mathbf{\Lambda}] = \text{EVD}(\hat{\mathbf{S}}_b, \hat{\mathbf{S}}_w)$;
Output: $\mathbf{B}, \mathbf{W}, \mathbf{\Lambda}$

The sketch matrix keeps a small number of selected frequent directions, which are obtained and updated by a matrix sketch technique during the whole online learning process. While sketching main data variations, the population mean and the one of each class are also updated. We further utilize these updated means and the low rank sketch matrix to estimate between-class covariance matrix and derive the approximate within-class covariance matrix after all new coming samples are compressed into the sketch matrix. Finally, we generate discriminant components by eigenvalue decomposition for simultaneously minimizing the approximate within-class variance and maximizing the between-class variance. The whole procedure of SoDA is illustrated in Figure 1 and presented in Algorithm 1. The in-depth theoretical investigation to explain why SoDA can approximate the offline FDA model by sketch and guarantee its effectiveness on extracting discriminant components is provided in Sec. IV.

A. Estimating Between-class covariance matrix

During online learning, we keep updating the population mean \mathbf{m}_0 and mean of each class \mathbf{m}_c ($c = 1, 2, \dots, C$) so as to construct the between-class covariance matrix \mathbf{S}_b . When having a new coming sample \mathbf{x}_i with class label \mathbf{y}_i , the population mean and mean of class \mathbf{y}_i are updated by

$$\mathbf{m}_c = (N_c \mathbf{m}_c + \mathbf{x}_i) / (N_c + 1), \quad c = 0, \mathbf{y}_i, \quad (1)$$

and the population number and the number of samples for class \mathbf{y}_i are also updated by:

$$N_c = N_c + 1, \quad c = 0, \mathbf{y}_i. \quad (2)$$

We then use the updated means to estimate the between-class covariance matrix as follows:

$$\mathbf{S}_b = \sum_{c=1}^C \frac{N_c}{N_0} (\mathbf{m}_c - \mathbf{m}_0)(\mathbf{m}_c - \mathbf{m}_0)^T. \quad (3)$$

B. Estimating Approximate Within-class covariance matrix

For realizing one-pass online learning, we aim to update/form the within-class covariance matrix which describes the within-class variation without using any passed observed data samples. Different from previous online FDA approaches, we embed sketch processing into FDA and derive a novel approximate within-class covariance matrix efficiently and effectively. For this purpose, we first employ the sketch technique [24] to compress the passed data samples into a sketch matrix so as to maintain the main variations of passed data. More specifically, we maintain the main variations of all passed data \mathbf{X} in a small size matrix $\mathbf{B} \in \mathcal{R}^{2\ell \times d}$, called a sketch matrix, where \mathbf{B} is initialized by a zero matrix. Each new coming sample \mathbf{x}_i^T (i.e. the i -th row of \mathbf{X}) replaces a zero row of \mathbf{B} from top to bottom until \mathbf{B} is full without any all zero rows. When \mathbf{B} is full, we apply Singular Value Decomposition (SVD) on \mathbf{B} such that $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{B}$, where $\mathbf{\Sigma}$ is a diagonal matrix with singular values on the diagonal in decreasing order. Each row in \mathbf{V}^T corresponds to a singular value in $\mathbf{\Sigma}$, and let vectors $\{\mathbf{v}_j\}$ of \mathbf{V}^T corresponding to the first half singular values denoted as *frequent directions* and the ones corresponding to lower half singular values denoted as *unfrequent directions*. By employing the sketch algorithm, the frequent directions \mathbf{v}_j are scaled by $\sqrt{\lambda_j^2 - \xi^2}$ and retained in \mathbf{B} , where ξ is the $(\ell + 1)^{\text{th}}$ largest singular value in $\Sigma_{\ell+1}$ of $\mathbf{\Sigma}$. In this way, the sketch matrix \mathbf{B} is obtained by $\hat{\mathbf{\Sigma}}\mathbf{V}^T$, where $\hat{\mathbf{\Sigma}} = \sqrt{\max(\mathbf{\Sigma}^2 - \mathbf{I}_{2\ell}\xi^2, \mathbf{O})}$ and \mathbf{O} is a zero matrix. Therefore, the sketch matrix \mathbf{B} is a $2\ell \times d$ matrix, where \mathbf{B}^+ , the upper half of \mathbf{B} , retains the main variations of passed data samples, and \mathbf{B}^- , the lower half of \mathbf{B} , is reset to zero.

Although no passed observed data samples are saved, we propose to derive an approximate within-class covariance matrix using the sketch matrix \mathbf{B} below:

$$\tilde{\mathbf{S}}_w = \tilde{\mathbf{S}}_t - \mathbf{S}_b, \quad (4)$$

where

$$\tilde{\mathbf{S}}_t = \mathbf{B}^T \mathbf{B} / N_0 - \mathbf{m}_0 \mathbf{m}_0^T. \quad (5)$$

In the above, $\tilde{\mathbf{S}}_w$ is not always the exact within-class covariance matrix but it is an approximate one. In Sec. IV, we will provide in-depth theoretical analysis of the bias of this approximation on discriminant feature component extraction.

C. Dimension Reduction and Extraction of Discriminant Components

Normally, after updating the two covariance matrices \mathbf{S}_b and $\tilde{\mathbf{S}}_w$, it is only necessary to compute the generalized eigen-vectors of $\mathbf{\Lambda}\tilde{\mathbf{S}}_w\mathbf{W} = \mathbf{S}_b\mathbf{W}$. However, in person re-identification, some kinds of features are of high dimensionality such as HIPHOP [5], LOMO [22] and etc, and the size of the two covariance matrices \mathbf{S}_b and $\tilde{\mathbf{S}}_w$ was determined by the feature dimensionality. Thus the above eigen-decomposition remains costly when the size of both \mathbf{S}_b and $\tilde{\mathbf{S}}_w$ are large.

An intuitive solution is to conduct another online learning for dimension reduction, which spends extra time and space. However, SoDA does not require such an extra learning. Due to sketch, SoDA actually maintains a set of frequent directions that describe main data variations. And thus we take these

frequent directions as basis vectors and the span of them can approximate the data space. Hence, we set $\mathbf{P} = \mathbf{V}^{T+}$, the upper half of matrix \mathbf{V}^T (Line 16 in Algorithm 1), and the dimension reduction is performed by:

$$\begin{aligned}\hat{\mathbf{S}}_b &= \mathbf{P}^T \mathbf{S}_b \mathbf{P}, \\ \hat{\mathbf{S}}_w &= \mathbf{P}^T \mathbf{S}_w \mathbf{P},\end{aligned}\quad (6)$$

where $\mathbf{P} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$ consists of k frequent directions. In this way, $\hat{\mathbf{S}}_b$ and $\hat{\mathbf{S}}_w$ become matrices in $\mathcal{R}^{k \times k}$, and computing generalized eigen-vectors will become much faster. Finally, the generalized eigen-vectors (Line 22 in Algorithm 1) are computed by $\Lambda \hat{\mathbf{S}}_w \mathbf{W} = \hat{\mathbf{S}}_b \mathbf{W}$, and they are the discriminant components we pursuit.

D. Computational Complexity

As presented above, after processing all observed samples, we maintain $\mathbf{B} \in \mathcal{R}^{\ell \times d}$, $\mathbf{P} \in \mathcal{R}^{d \times k}$, $\mathbf{m}_c \in \mathcal{R}^d$ and $N_c (c = 0, 1, 2, \dots, C)$. The time and space cost of the rest procedure is $\mathcal{O}(d\ell^2)$ (After the whole processing, N_0 is equal to N) and $\mathcal{O}((\ell + C)d)$, respectively. Therefore, the cost of time and space is $\mathcal{O}(d\ell^2)$ and $\mathcal{O}((\ell + k + C)d)$, respectively, almost the same as the cost of sketch algorithm [24].

IV. THEORETICAL ANALYSIS

In this section, we theoretically show that SoDA approximates FDA in a principled way, although SoDA is formed based on the approximate within-class covariance matrix mined from sketch data information.

A. Fisher Discriminant Analysis

Fisher discriminant analysis (FDA) aims to seek discriminant projections for minimizing within-class variance and maximizing between-class variance, which are estimated over the data matrix \mathbf{X} and its label set \mathbf{y} in an offline way. There are several equivalent criteria \mathbf{J}_F for the multi-class case. For analysis, we consider the one that maximizes the following criterion:

$$\mathbf{J}_F(\mathbf{W}) = \frac{\mathbf{W}^T \mathbf{S}_b \mathbf{W}}{\mathbf{W}^T \mathbf{S}_w \mathbf{W}}, \quad (7)$$

where \mathbf{S}_b is the between-class covariance matrix and \mathbf{S}_w is the within-class covariance matrix. They are given by

$$\mathbf{S}_b = \sum_{c=1}^C \frac{N_c}{N} (\mathbf{m}_c - \mathbf{m}_0)(\mathbf{m}_c - \mathbf{m}_0)^T, \quad (8)$$

$$\mathbf{S}_w = \sum_{c=1}^C \frac{N_c}{N} \sum_{\mathbf{y}_i=c} \frac{1}{N_c} (\mathbf{x}_i - \mathbf{m}_c)(\mathbf{x}_i - \mathbf{m}_c)^T, \quad (9)$$

where \mathbf{m}_c and N_c are the data mean and the number of samples of the c^{th} class, respectively, and N and \mathbf{m}_0 are the population number and population mean, respectively. And the total covariance matrix is

$$\mathbf{S}_t = \mathbf{S}_w + \mathbf{S}_b = \frac{1}{N} \sum_{c=1}^C \sum_{\mathbf{y}_i=c} (\mathbf{x}_i - \mathbf{m}_0)(\mathbf{x}_i - \mathbf{m}_0)^T. \quad (10)$$

Generally, the analysis seeks a set of feature vectors $\{\mathbf{w}_j\}$ that maximize the criterion subject to the normalization constraint $\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) = 1$, where \mathbf{W} is the matrix whose columns are $\{\mathbf{w}_j\}$. This leads to the computation of generalized

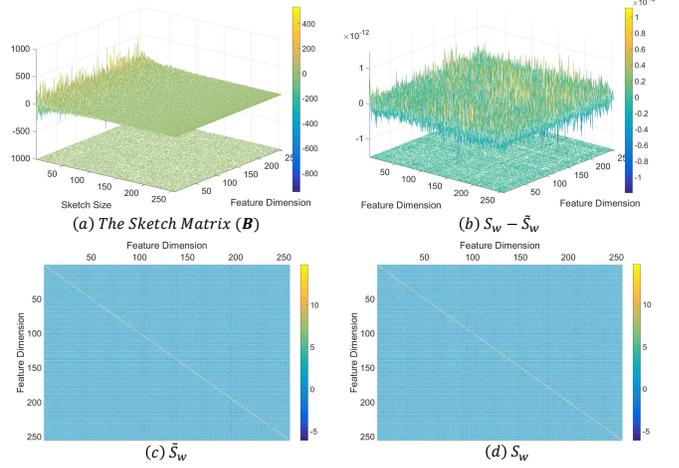


Fig. 2. (a) is the sketch matrix (\mathbf{B}). (c) is the approximate within-class covariance matrix ($\tilde{\mathbf{S}}_w$) generated by SoDA while (d) is the groundtruth one (\mathbf{S}_w) produced by FDA. (b) is the difference ($\mathbf{S}_w - \tilde{\mathbf{S}}_w$) of the groundtruth within-class covariance matrix and the approximate one. It is noteworthy that the distinction between \mathbf{S}_w and $\tilde{\mathbf{S}}_w$ is less than 1×10^{-12} , which indicates that $\tilde{\mathbf{S}}_w$ estimated by SoDA can approximate the groundtruth one (Best viewed in color).

eigenvectors, that is $\Lambda \mathbf{S}_w \mathbf{W} = \mathbf{S}_b \mathbf{W}$ and Λ is a diagonal matrix with generalized eigenvalues on the diagonal. Here, the eigenvectors corresponding to the largest eigenvalues are used to compress a high dimensional data vector to a low dimensional feature representation.

B. Relation between SoDA and FDA

Before presenting the theoretical analysis, we first define the following notations. Let

$$\begin{aligned}\mathbf{J}_F^1(\mathbf{W}) &= \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}, \\ \mathbf{J}_F^2(\mathbf{W}) &= \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \tilde{\mathbf{S}}_w \mathbf{W})},\end{aligned}\quad (11)$$

where $\mathbf{J}_F^1(\mathbf{W})$ is the conventional FDA criterion and $\mathbf{J}_F^2(\mathbf{W})$ is SoDA criterion by replacing \mathbf{S}_w with $\tilde{\mathbf{S}}_w$ that is mined from sketch data information.

Let the largest Fisher scores in the above equations be

$$\begin{aligned}\mathbf{J}_F^1(\mathbf{W}^1) &= \max_{\mathbf{W} \in \mathcal{R}^{d \times k}} \mathbf{J}_F^1(\mathbf{W}) = \mu_1, \\ \mathbf{J}_F^2(\mathbf{W}^2) &= \max_{\mathbf{W} \in \mathcal{R}^{d \times k}} \mathbf{J}_F^2(\mathbf{W}) = \mu_2.\end{aligned}\quad (12)$$

Since for optimizing Eqs. (12), we can form a Lagrangian function by imposing the constraint $\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) = 1$ for both criteria [41] We define $\mathcal{D} = \{\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in \mathcal{R}^{d \times k} | \text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) = 1\}$, and thus we can reform Eqs. (12) by:

$$\begin{aligned}\mu_1^{-1} &= \min_{\mathbf{W}^1 \in \mathcal{D}} \{\mathbf{J}_F^1(\mathbf{W}_1)\}^{-1} = \text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W}), \\ \mu_2^{-1} &= \min_{\mathbf{W}^2 \in \mathcal{D}} \{\mathbf{J}_F^2(\mathbf{W}_2)\}^{-1} = \text{tr}(\mathbf{W}^T \tilde{\mathbf{S}}_w \mathbf{W}).\end{aligned}\quad (13)$$

In the following sections, we first discuss the relationship between μ_1 and μ_2 . And then this relationship will be used to present a bound for $\mathbf{J}_F^1(\mathbf{W}^2)$. Note that $\mathbf{J}_F^1(\mathbf{W}^2)$ is to measure how well the optimal projection learned by our SoDA approximates the optimal solution that maximizes $\mathbf{J}_F^1(\mathbf{W})$. Note

that our analysis will not take any dimension reduction before extracting discriminant components below for discussion. Our analysis can be extended if the same dimension reduction is applied to all methods discussed below.

C. Relationship Between the Maximum Fisher Score of FDA and that of SoDA

We first present the relationship between the maximum Fisher score of FDA and the one of SoDA, i.e. the relationship between μ_1 and μ_2 . Suppose that matrix $\mathbf{X} \in \mathcal{R}^{N \times d}$ is the totally training sample set consisting of samples acquired at each time step.

However, it is not intuitive to obtain the relationship between the maximum Fisher score of FDA and the one of SoDA based on the covariance matrices inferred in Eq. (5). In order to exploit such a relationship, we first investigate the Fisher score obtained by \mathbf{S}_b and the approximate within-class covariance matrix $\tilde{\mathbf{S}}_w$ as follows:

$$\tilde{\mathbf{S}}_w = \tilde{\mathbf{S}}_t - \mathbf{S}_b = \mathbf{B}^T \mathbf{B} / N - \mathbf{m}_0 \mathbf{m}_0^T - \mathbf{S}_b. \quad (14)$$

Let \mathbf{S}_w be the within-class covariance matrix computed in batch mode (i.e. for offline FDA). Since it is known that $\mathbf{S}_w = \mathbf{S}_t - \mathbf{S}_b = \mathbf{X}^T \mathbf{X} / N - \mathbf{m}_0 \mathbf{m}_0^T - \mathbf{S}_b$, it can be verified that

$$\begin{aligned} & \mathbf{S}_w - \tilde{\mathbf{S}}_w \\ &= (\mathbf{X}^T \mathbf{X} / N - \mathbf{m}_0 \mathbf{m}_0^T - \mathbf{S}_b) - (\mathbf{B}^T \mathbf{B} / N - \mathbf{m}_0 \mathbf{m}_0^T - \mathbf{S}_b) \quad (15) \\ &= (\mathbf{X}^T \mathbf{X} - \mathbf{B}^T \mathbf{B}) / N. \end{aligned}$$

By combining Eq. (25) as stated in the Appendix, it is not hard to have the following theorem about the relation between \mathbf{S}_w and $\tilde{\mathbf{S}}_w$, and we visualize the approximation between the groundtruth within-class covariance matrix and our approximate one in Figure 2. We assume that \mathbf{S}_w , $\tilde{\mathbf{S}}_w$ and $\hat{\mathbf{S}}_w$ are not singular in the following analysis¹.

Theorem 1. $\tilde{\mathbf{S}}_w \preceq \mathbf{S}_w$, and $\|\mathbf{S}_w - \tilde{\mathbf{S}}_w\| \leq 2\|\mathbf{X}\|_f^2 / (N\ell)$, where $\|\cdot\|$ is the induced norm of a matrix and $\|\cdot\|_f$ is the Frobenius norm.

Based on the above theorem, we particularly consider the two-class classification case.

Theorem 2. Considering the two criteria in Eq. (13) when the discriminant feature transformation is a one-dimensional vector, i.e. $\mathbf{W}^1 = \mathbf{w}^1 \in \mathcal{R}^d$ and $\mathbf{W}^2 = \mathbf{w}^2 \in \mathcal{R}^d$, the relationship between μ_1 and μ_2 is as follow:

$$\mu_1^{-1} - 2(s_0 r_b)^{-\frac{1}{2}} \|\mathbf{X}\|_f^2 / (N\ell) \leq \mu_2^{-1} \leq \mu_1^{-1}, \quad (16)$$

where s_0 is the smallest (non-zero) singular value of matrix \mathbf{S}_b and $r_b = \text{rank}(\mathbf{S}_b)$.

Proof. Let $\mathbb{D} = 2\|\mathbf{X}\|_f^2 / (N\ell)$. From the Theorem 1, we have for any nonzero $\mathbf{w} \in \mathcal{R}^d$, $0 \leq \frac{\mathbf{w}^T (\mathbf{S}_w - \tilde{\mathbf{S}}_w) \mathbf{w}}{\|\mathbf{w}\|_2} \leq \mathbb{D}$. That is $\forall \mathbf{w} \in \mathcal{R}^d$, $\mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} \leq \mathbf{w}^T \mathbf{S}_w \mathbf{w}$, $\mathbf{w}^T \mathbf{S}_w \mathbf{w} \leq \mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} + \mathbb{D} \|\mathbf{w}\|_2$. (17)

Let \mathbf{w}^1 and \mathbf{w}^2 be the discriminant vectors that minimize the Criterion in Eq. (13) under the constraints $\mathbf{w}^{1T} \mathbf{S}_b \mathbf{w}^1 = 1$ and $\mathbf{w}^{2T} \mathbf{S}_b \mathbf{w}^2 = 1$, respectively. That is $\mathbf{w}^{1T} \mathbf{S}_w \mathbf{w}^1 = \mu_1^{-1}$ and $\mathbf{w}^{2T} \tilde{\mathbf{S}}_w \mathbf{w}^2 = \mu_2^{-1}$, i.e. \mathbf{w}^1 and \mathbf{w}^2 would minimize

$\mathbf{w}^T \mathbf{S}_w \mathbf{w}$ and $\mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w}$ when constraining $\mathbf{w}^T \mathbf{S}_b \mathbf{w} = 1$. In addition, since $\mathbf{w}^{2T} \mathbf{S}_b \mathbf{w}^2 = 1$, we have $s_0 r_b \|\mathbf{w}^2\|_2^2 \leq 1$, i.e. $\|\mathbf{w}^2\|_2 \leq (s_0 r_b)^{-\frac{1}{2}}$. Therefore, based on Theorem 1, we have

$$\begin{aligned} \mu_2^{-1} &= \mathbf{w}^{2T} \tilde{\mathbf{S}}_w \mathbf{w}^2 \leq \mathbf{w}^{1T} \tilde{\mathbf{S}}_w \mathbf{w}^1 \leq \mathbf{w}^{1T} \mathbf{S}_w \mathbf{w}^1 = \mu_1^{-1}, \\ \mu_1^{-1} &= \mathbf{w}^{1T} \mathbf{S}_w \mathbf{w}^1 \leq \mathbf{w}^{2T} \mathbf{S}_w \mathbf{w}^2 \\ &\leq \mathbf{w}^{2T} \tilde{\mathbf{S}}_w \mathbf{w}^2 + \mathbb{D} (s_0 r_b)^{-\frac{1}{2}} = \mu_2^{-1} + \mathbb{D} (s_0 r_b)^{-\frac{1}{2}}. \end{aligned} \quad (18)$$

Then $\mu_1^{-1} - 2(s_0 r_b)^{-\frac{1}{2}} \|\mathbf{X}\|_f^2 / (N\ell) \leq \mu_2^{-1} \leq \mu_1^{-1}$. \square

From the theorem above, we can claim that the largest Fisher score $\mathbf{J}_F^2(\mathbf{w}^2)$ is always greater than or equal to the original one $\mathbf{J}_F^1(\mathbf{w}^1)$ after sketch. From another aspect, the inequalities “ $\mu_1^{-1} - 2(s_0 r_b)^{-\frac{1}{2}} \|\mathbf{X}\|_f^2 / (N\ell) \leq \mu_2^{-1} \leq \mu_1^{-1}$ ” means when more rows are set in the sketch matrix \mathbf{B} , (i.e. much larger ℓ is set), μ_2 becomes μ_1 , and thus SoDA becomes exactly the FDA.

For the multi-class case, we can generalize the above proof below.

Theorem 3. Considering the two criteria in Eq. (13), when the discriminant feature transformation is a d -dimensional transformation where $d > 1$, we have $\mu_1 \leq \mu_2$.

Proof. Note that \mathbf{W}^1 and $\mathbf{W}^2 (\in R^{d \times k})$ make the two criteria minimized in Eq. (13), respectively. Let $\mathbf{W}^1 = [\mathbf{w}_1^1, \dots, \mathbf{w}_k^1]$ and $\mathbf{W}^2 = [\mathbf{w}_1^2, \dots, \mathbf{w}_k^2]$. Since for any $\mathbf{w} \in R^d$, $\mathbf{w}^T \mathbf{S}_w \mathbf{w} \geq \mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w}$ by Theorem 1, we have

$$\begin{aligned} \mu_2^{-1} &= \text{tr}(\mathbf{W}^{2T} \tilde{\mathbf{S}}_w \mathbf{W}^2) \leq \text{tr}(\mathbf{W}^{1T} \tilde{\mathbf{S}}_w \mathbf{W}^1) \\ &= \sum_{i=1}^k \mathbf{w}_i^{1T} \tilde{\mathbf{S}}_w \mathbf{w}_i^1 \leq \sum_{i=1}^k \mathbf{w}_i^{1T} \mathbf{S}_w \mathbf{w}_i^1 \\ &= \text{tr}(\mathbf{W}^{1T} \mathbf{S}_w \mathbf{W}^1) = \mu_1^{-1}. \end{aligned}$$

Hence, the theorem is proved. \square

D. How Does the Projection Learned by SoDA Optimize the Original Fisher Criterion Approximately?

In the above, we analyze the quotient values between $\frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}$ and $\frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \tilde{\mathbf{S}}_w \mathbf{W})}$. However, in SoDA, our within-class covariance matrix is estimated by sketch and is not the exact within-class covariance matrix. In the following, we will present the effect of the learned discriminant component using SoDA on minimizing the ground-truth within-class covariance. For this purpose, the following theorems are presented.

Theorem 4. For any $\mathbf{w} \in \mathcal{Q} = \{\mathbf{w} \in R^d | \mathbf{w}^T \mathbf{w} = 1\}$, we have

$$\mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} \leq \mathbf{w}^T \mathbf{S}_w \mathbf{w} \leq \mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} + \frac{2}{N} \|\mathbf{X}\|_f^2 / \ell. \quad (19)$$

Proof. While the inequality $\mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} \leq \mathbf{w}^T \mathbf{S}_w \mathbf{w}$ is obvious by using Theorem 1, we focus on the latter one. Since $\mathbf{S}_w = \tilde{\mathbf{S}}_w + (\mathbf{X}^T \mathbf{X} - \mathbf{B}^T \mathbf{B}) / N$ in Eq. (15), by applying Eq. (25), we have $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = \mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} + \frac{1}{N} \mathbf{w}^T (\mathbf{X}^T \mathbf{X} - \mathbf{B}^T \mathbf{B}) \mathbf{w} \leq \mathbf{w}^T \tilde{\mathbf{S}}_w \mathbf{w} + \frac{2}{N} \|\mathbf{X}\|_f^2 / \ell$. \square

Theorem 5. Considering the two criteria in Eq. (13), we define $\mathcal{D} = \{\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in R^{d \times k} | \text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) = 1\}$, denote the smallest non-zero singular value of \mathbf{S}_b as s_0 , and

¹The analysis can be generalized to the case when $\tilde{\mathbf{S}}_w$ is not invertible if the same regularization is imposed on both \mathbf{S}_w , $\tilde{\mathbf{S}}_w$ and $\hat{\mathbf{S}}_w$

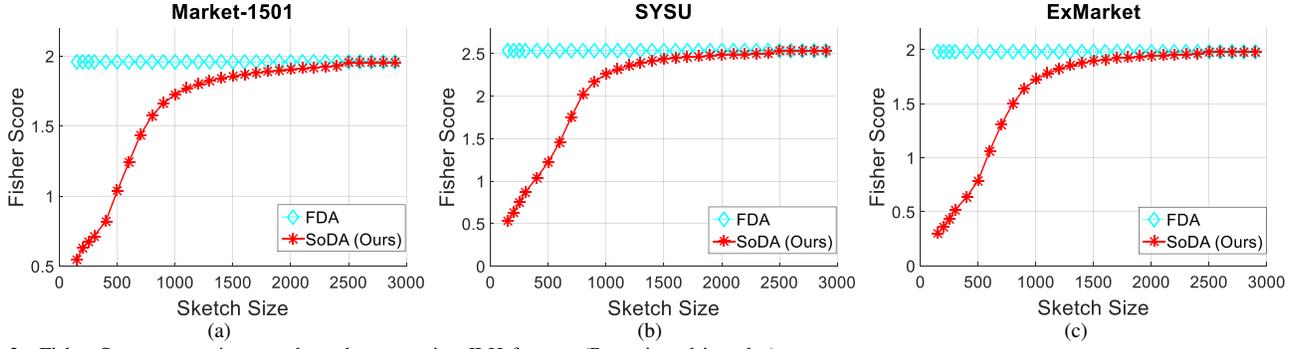


Fig. 3. Fisher Score comparison on three datasets using JLH feature. (Best viewed in color).

TABLE I
COMPARISON AMONG DIFFERENT ONLINE/INCREMENTAL APPROACHES

Approaches	IFDA [15]	Pang's IFDA [33]	IDR/QR [48]	OL-IDM [37]	Wang et al. [39]	Martinel et al. [29]	SoDA (Ours)
Save within-class scatter matrix?	✓	✓	✓	--	--	--	✗
Save between-class scatter matrix?	✓	✓	✗	--	--	--	✗
Is an one-pass algorithm?	✗	✓	✓	✓	✗	✗	✓
Human feedback	✗	✗	✗	✗	✓	✓	✗
Can the model be trained on streaming data?	✓	✓	✓	✓	✗	✗	✓
Is the model embedded with dimension reduction?	✗	✗	✗	✗	✗	✗	✓
time complexity	$\mathcal{O}(d^3)$	$\mathcal{O}(nd^2)$	$\mathcal{O}(ndc)$	--	--	--	$\mathcal{O}(\min(\ell, d)^2 \max(\ell, d))$
space complexity	$\mathcal{O}(d^2)$	$\mathcal{O}(d^2)$	$\mathcal{O}(d^2)$	--	--	--	$\mathcal{O}((\ell + k + C)d)$

let $r_b = \text{rank}(\mathbf{S}_b)$. Suppose the norm of each data vector \mathbf{x}_i (i.e. each row of the data matrix $\mathbf{X} \in R^{N \times d}$) is bounded by M , that is $\|\mathbf{x}_i\|_2 \leq M$. Then we have

$$\frac{1}{\mu_1^{-1} + \frac{2k}{s_0 r_b} M / \ell} \leq \mathbf{J}_F^1(\mathbf{W}^2) \leq \mu_1. \quad (20)$$

Proof. First, given $\mathbf{W}^2 \in \mathcal{D}$ that minimize $\{\mathbf{J}_F^2(\mathbf{W})\}^{-1}$.

$$\begin{aligned} \{\mathbf{J}_F^1(\mathbf{W}^2)\}^{-1} &= \text{tr}(\mathbf{W}^{2T} \mathbf{S}_w \mathbf{W}^2) \\ &= \sum_{i=1}^k \mathbf{w}_i^{2T} \mathbf{S}_w \mathbf{w}_i^2 \\ &= \sum_{i=1}^k \|\mathbf{w}_i^2\|_2^2 \frac{\mathbf{w}_i^{2T}}{\|\mathbf{w}_i^2\|_2} \mathbf{S}_w \frac{\mathbf{w}_i^2}{\|\mathbf{w}_i^2\|_2} \\ &\leq \sum_{i=1}^k \|\mathbf{w}_i^2\|_2^2 \frac{\mathbf{w}_i^{2T}}{\|\mathbf{w}_i^2\|_2} \tilde{\mathbf{S}}_w \frac{\mathbf{w}_i^2}{\|\mathbf{w}_i^2\|_2} \\ &\quad + \sum_{i=1}^k \frac{2}{N} \|\mathbf{w}_i^2\|_2^2 \|\mathbf{X}\|_f^2 / \ell \\ &\leq \sum_{i=1}^k \mathbf{w}_i^{2T} \tilde{\mathbf{S}}_w \mathbf{w}_i^2 \\ &\quad + \frac{2k}{N} \|\mathbf{w}_i^2\|_2^2 \|\mathbf{X}\|_f^2 / \ell. \end{aligned} \quad (21)$$

Since $\text{tr}(\mathbf{W}^{2T} \mathbf{S}_b \mathbf{W}^2) = 1$, we have $\mathbf{w}_i^{2T} \mathbf{S}_b \mathbf{w}_i^2 \leq 1$. Here, for convenience, one can further assume $\mathbf{w}_i^{2T} \mathbf{S}_b \mathbf{w}_i^2 > 0$, otherwise a much tighter bound can be inferred. And thus $s_0 r_b \|\mathbf{w}_i^2\|_2^2 \leq 1$. So we have

$$\begin{aligned} \{\mathbf{J}_F^1(\mathbf{W}^2)\}^{-1} &= \text{tr}(\mathbf{W}^{2T} \mathbf{S}_w \mathbf{W}^2) \\ &\leq \sum_{i=1}^k \mathbf{w}_i^{2T} \tilde{\mathbf{S}}_w \mathbf{w}_i^2 + \frac{2k}{N} (s_0 r_b)^{-1} \|\mathbf{X}\|_f^2 / \ell \\ &= \mu_2^{-1} + \frac{2k}{N} (s_0 r_b)^{-1} \|\mathbf{X}\|_f^2 / \ell. \end{aligned} \quad (22)$$

Note that $\mu_2^{-1} = \sum_{i=1}^k \mathbf{w}_i^{2T} \tilde{\mathbf{S}}_w \mathbf{w}_i^2$ since it is assumed that $\mathbf{W}^2 \in \mathcal{D}$ minimizes $\{\mathbf{J}_F^2(\mathbf{W})\}^{-1}$. Thus, under the constraint $\text{tr}(\mathbf{W}^{2T} \mathbf{S}_b \mathbf{W}^2) = 1$, we have

$$\frac{1}{\mu_2^{-1} + \frac{2k}{N} (s_0 r_b)^{-1} \|\mathbf{X}\|_f^2 / \ell} \leq \mathbf{J}_F^1(\mathbf{W}^2) \leq \mu_1, \quad (23)$$

where the latter equation is obvious since \mathbf{W}^2 may not be the optimal projection for maximizing $J_F^1(\mathbf{W})$. Finally, since $\mu_1 \leq \mu_2$ and $\|\mathbf{x}_i\|_2 \leq M$ that means the norm of any data vector \mathbf{x}_i (i.e. each row of the data matrix $\mathbf{X} \in R^{N \times d}$) is bounded by M , we have

$$\frac{1}{\mu_1^{-1} + \frac{2k}{s_0 r_b} M / \ell} \leq \mathbf{J}_F^1(\mathbf{W}^2) \leq \mu_1. \quad (24)$$

□

E. Discussion

1) *SoDA vs. FDA*: The above theorem indicates that 1) the learned transformation by SoDA may not be the optimal one for the FDA directly learned on all observed data since $\mathbf{J}_F^1(\mathbf{W}^2) \leq \mu_1$, which is obvious and reasonable; 2) however, there is a lower bound on $\mathbf{J}_F^1(\mathbf{W}^2)$, since $\frac{1}{\mu_1^{-1} + \frac{2k}{s_0 r_b} M / \ell} \leq$

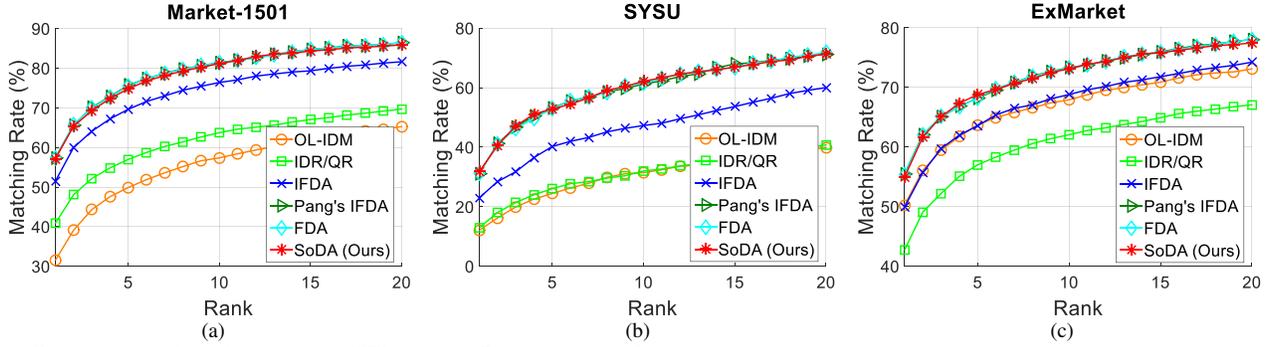


Fig. 4. Comparison on three datasets using JSTL feature. (Best viewed in color).

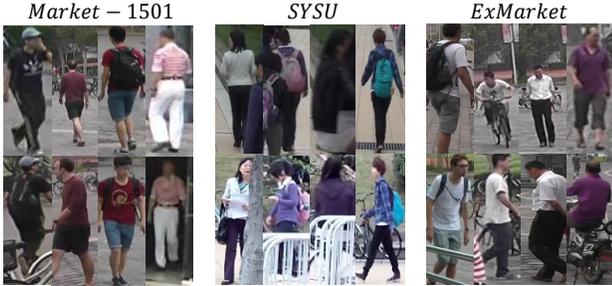


Fig. 5. Example images from different person re-id datasets. For each dataset, two images in a column correspond to the same person.

$\mathbf{J}_F^1(\mathbf{W}^2)$; 3) as long as more and more rows are set in the sketch matrix \mathbf{B} used in SoDA, i.e. ℓ is larger and larger, $\frac{2k}{s_0}M/\ell \rightarrow 0$ and so that $\mathbf{J}_F^1(\mathbf{W}^2) \approx \mu_1$ in such a case. The latter case is reasonable because although the sketch in SoDA enables selecting data variation during the online learning, more data information is kept when a much larger sketch matrix \mathbf{B} is used, and this will be verified in the experiments (see Figure 3 for example).

2) *SoDA vs. Incremental/online models*: In Table I, we compare SoDA with related incremental/online FDA models in details. A distinct and important characteristic of SoDA is that it is able to perform one-pass online learning directly only relying on sketch data information. SoDA does not have to keep within-class covariance matrix and between-class covariance matrix in memory during online learning, due to embedding sketch processing, which has not been considered for online learning of FDA before. Moreover, as compared to the others, SoDA does not need any extra online learning progress on dimension reduction, which is naturally embedded. Thus the training cost of SoDA is much lighter.

When applied SoDA to person re-id, we perform the comparison with related online person re-id models. An important distinction is that no extra human feedback is required, and SoDA is able to be applied on streaming data in an one-pass learning manner. In comparison with OL-IDM, SoDA has its merits: 1) dimension reduction is naturally embedded in SoDA; 2) embedding sketch into person re-id model learning is a more efficient and effective way to maintain the main variations of data, which has been verified by our experimental results.

V. EXPERIMENTS

A. Datasets and Evaluation Settings

1) *Datasets*: We extensively evaluated the proposed approach on three large person re-id benchmarks: Market-1501, SYSU, and ExMarket.

- **Market-1501** dataset [51] contains person images collected in front of a campus supermarket at a University. It consists of 32,643 person images of 1,501 identities.
- **SYSU** dataset contains totally 48,892 images of 502 pedestrians captured by two cameras. Similar to [4], we randomly selected 251 identities from two views as training set which contains 12308 images. And we randomly selected three images of each person from the rest 251 identities of both cameras to form the testing set, where the 753 images of the first camera were used as query images.
- **ExMarket** dataset was formed by combining the MARS dataset [50] and Market-1501 dataset. MARS was formed as a video dataset for person re-identification. All the identities from MARS are of a subset of those from Market. More specifically, for each identity, we extracted one frame for each five consecutive frames firstly and combined images extracted from MARS and the ones from Market-1501 of the same person. Therefore, ExMarket contains 237147 images of 1501 identities, the largest population size among the three benchmark datasets tested.

2) *Features*: In this work, we conducted the evaluation based on four types of feature for evaluation: 1) JSTL, 2) LOMO, 3) HIPHOP, 4) JSTL + LOMO + HIPHOP (JLH).

- JSTL is a kind of low-dimensional deep feature representation (\mathcal{R}^{256}) extracted by a deep convolutional network [43];
- LOMO is an effective handcraft feature proposed for person re-id in [22], and it is a 26960-dimensional vector;
- HIPHOP is another recently proposed person re-id feature (\mathcal{R}^{84096}) [5] that extracts more view invariant histogram features from shallow layers of a convolution network.

In addition, since person re-id can benefit from using multiple different types of appearance features as shown in [5], [7], [9], [49], [52]. we concatenated JSTL, LOMO and HIPHOP as a high dimensional feature (\mathcal{R}^{111312}), named **JLH**

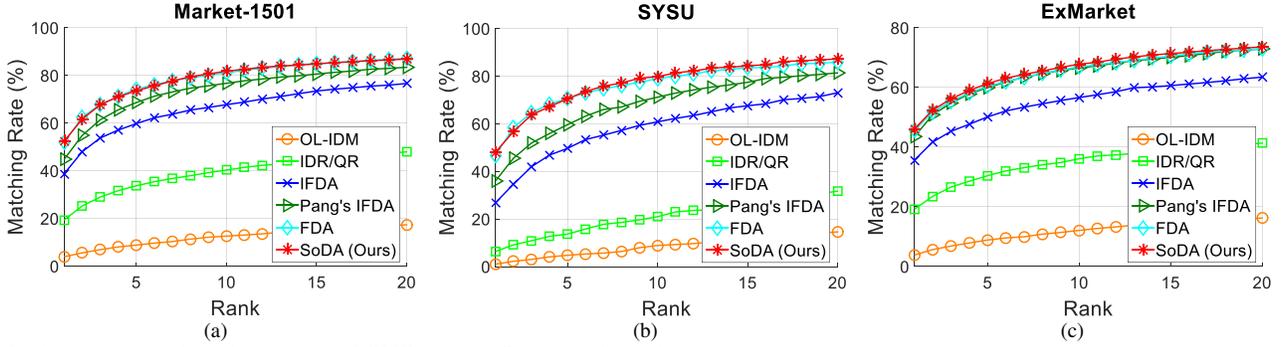


Fig. 6. Comparison on three datasets using LOMO feature. (Best viewed in color).

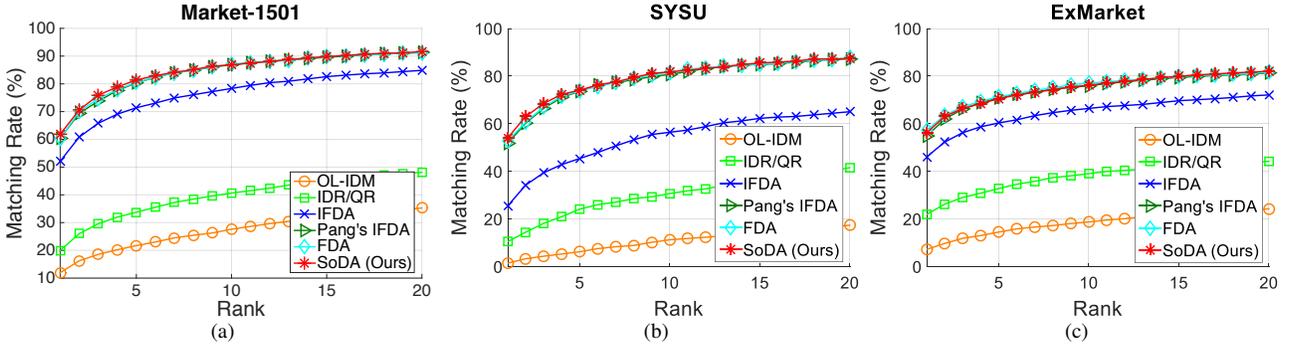


Fig. 7. Comparison on three datasets using HIPHOP feature. (Best viewed in color).

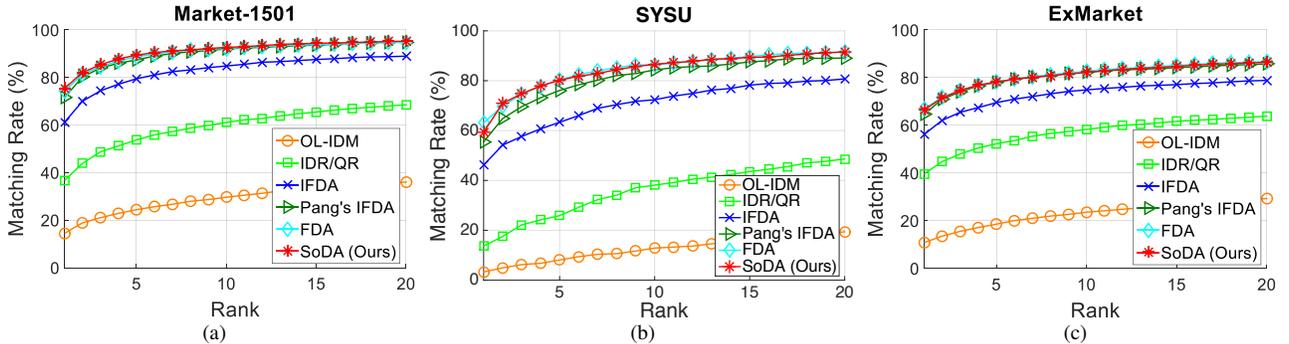


Fig. 8. Comparison on three datasets using JLH feature. (Best viewed in color).

in this work for convenience of description. On all datasets, we report experimental results of SoDA using the concatenated feature in Table VI. Since LOMO, HIPHOP, and JLH are of high dimension, for all methods except SoDA, we first reduced their feature dimension of the three types of feature to 2000, 2000 and 2500, respectively, on all datasets. For SoDA, we set the sketch size (ℓ) to the (reduced) feature dimension mentioned above on all datasets.

3) *Evaluation protocol*: On all datasets, we followed the standard evaluation settings on person re-identification, i.e. images of half of the persons were used for training and images of the rest half were used for testing, so that there is no overlap in persons between training and testing sets. More specifically, on Market-1501 dataset, we used the standard training (12936 images of 750 people) and testing (19732 images of 751 people) sets provided in [51]. On SYSU dataset,

similar to [4], we randomly picked all images of the selected 251 identities from two views to form the training set which contains 12308 images, and we randomly picked 3 images of each pedestrian of the rest 251 identities in each view for forming the gallery and query sets for testing. On ExMarket dataset, we conducted the same identity split as the Market-1501 dataset. The training set contains 112351 images, and the testing set contains 124796 images, among which 3363 images are considered as query images and the rest are considered as gallery images.

On all datasets, the cumulative matching characteristic (CMC) curves is shown to measure the performance of the compared methods on re-identifying individuals across different camera views under online setting. In addition to this, we also report results using another two performance metrics: 1) rank-1 Matching Rate, and 2) mean Average Precision

TABLE II
COMPARISON WITH FDA ON ALL BENCHMARKS.

Feature		JSTL					LOMO					HIPHOP					JLH				
Dataset	Method	rank-1	rank-5	rank-10	rank-20	mAP	rank-1	rank-5	rank-10	rank-20	mAP	rank-1	rank-5	rank-10	rank-20	mAP	rank-1	rank-5	rank-10	rank-20	mAP
Market-1501	FDA	57.30	75.53	81.38	86.49	28.57	51.90	74.26	81.12	87.14	23.60	60.27	80.52	87.05	91.18	31.45	74.20	88.75	92.19	94.80	49.01
	SoDA	57.13	74.79	81.18	85.90	28.25	52.41	73.37	81.38	87.17	23.58	61.88	81.41	86.70	91.60	33.39	75.27	89.28	92.70	95.22	49.82
SYSU	FDA	31.21	52.99	61.49	71.85	25.86	46.61	70.78	79.42	86.19	41.81	52.86	73.84	81.67	87.78	48.20	63.08	80.35	86.32	91.50	56.82
	SoDA	31.74	52.86	62.15	71.31	26.04	47.81	70.39	78.75	86.72	41.69	53.12	73.97	80.88	87.25	48.48	64.81	80.74	87.25	91.77	59.82
Ex-Market	FDA	53.89	68.11	73.13	77.97	22.71	45.64	60.42	66.86	72.89	17.98	57.24	71.38	77.11	81.74	27.20	66.86	78.18	82.63	86.70	39.00
	SoDA	54.93	68.79	73.13	77.46	22.87	46.08	61.31	67.81	73.63	17.77	55.76	70.40	76.10	81.59	24.97	66.18	78.36	82.48	86.64	37.11

TABLE III
COMPARISON WITH INCREMENTAL FDA MODELS AND ONLINE METHOD USING JSTL.

Dataset	Market-1501			SYSU			ExMarket		
Method	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)
OL-IDM	31.50	10.48	3706.84	12.08	10.29	10588.15	50.24	18.93	1646433.70
IDR/QR	41.15	13.20	803.59	12.88	10.24	247.17	42.70	11.20	6172.79
IFDA	51.45	21.21	38.22	22.97	18.12	12.40	49.91	16.58	394.31
Pang's IFDA	57.36	28.58	13.68	31.08	25.28	7.65	55.46	22.97	120.94
SoDA	57.13	28.25	7.84	31.74	26.04	4.68	54.93	22.87	50.52

TABLE IV
COMPARISON WITH INCREMENTAL FDA MODELS AND ONLINE METHOD USING LOMO.

Dataset	Market-1501			SYSU			ExMarket		
Method	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (min)
OL-IDM	3.95	0.73	736707.11	1.06	1.59	743335.02	3.86	0.33	> 1 week
IDR/QR	19.36	5.09	345181.63	6.37	5.16	83903.98	19.92	3.58	74393.24
IFDA	38.75	13.32	314470.08	26.83	22.59	67003.60	35.63	10.43	69668.26
Pang's IFDA	44.80	18.64	314461.09	35.99	31.82	66646.88	43.50	15.42	69625.84
SoDA	52.41	23.53	2127.47	47.81	41.69	3345.30	46.08	17.77	359.28

TABLE V
COMPARISON WITH INCREMENTAL FDA MODELS AND ONLINE METHOD USING HIPHOP.

Dataset	Market-1501			SYSU			ExMarket		
Method	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)
OL-IDM	11.97	2.22	277104.72	1.46	2.00	252626.33	7.24	0.54	> 1 week
IDR/QR	19.98	6.00	225226.32	10.49	9.34	86513.64	21.97	5.32	2392922.71
IFDA	52.14	21.30	185390.31	25.50	22.32	66202.88	46.08	15.50	2133499.12
Pang's IFDA	60.42	31.30	185174.97	51.79	47.51	65593.56	54.84	25.11	2135671.23
SoDA	61.88	33.39	3620.00	53.12	48.48	13849.61	55.76	24.97	83319.79

(mAP). mAP first computes the area under the Precision-Recall curve for each query and then calculates the mean of Average Precision over all query persons. All experiments were implemented using MATLAB on a machine with CPU E5 2686 2.3 GHz and 256 GB RAM, and the accumulative time of all compared methods were also computed and reported for measuring efficiency.

B. SoDA vs. FDA

In Sec. IV, we provide theoretical analysis on the relation between SoDA and FDA. In this section, we provide empirical evaluation on three datasets by the comparison on Fisher Score between SoDA and FDA in Figure 3. The figure indicates that by keeping more rows in the sketch matrix, SoDA can acquire more similar Fisher Score as the one of FDA, and this is supported by Theorem 5. We also compared SoDA with FDA on the three datasets in Table II, and the comparison shows that they work comparably. Therefore the results reported here have validated that our sketch approach approximates FDA

(i.e. the offline model) for extracting discriminant information very well, and thus the effectiveness of our model is verified both theoretically and empirically.

C. SoDA vs. Incremental FDA Model

There are existing works that are related to incremental learning of FDA, which also process sequential data and update the models online. We compared extensively our method SoDA with three related online/incremental FDA methods, including IFDA [15], IDR/QR [48] and Pang's IFDA [33]. We show CMC curve of all methods using different types of features in Figure 4, Figure 6, Figure 7 and Figure 8. The results illustrate that the proposed SoDA outperformed the compared incremental FDA. For instance, when using JLH, SoDA outperformed Pang's IFDA and achieved 75.27%, 64.81% and 66.18% rank-1 matching rate on Market, SYSU and ExMarket, respectively. We further report mAP and accumulative time in Table III, Table IV, Table V and Table VI. It suggests that SoDA has a better mAP values especially on

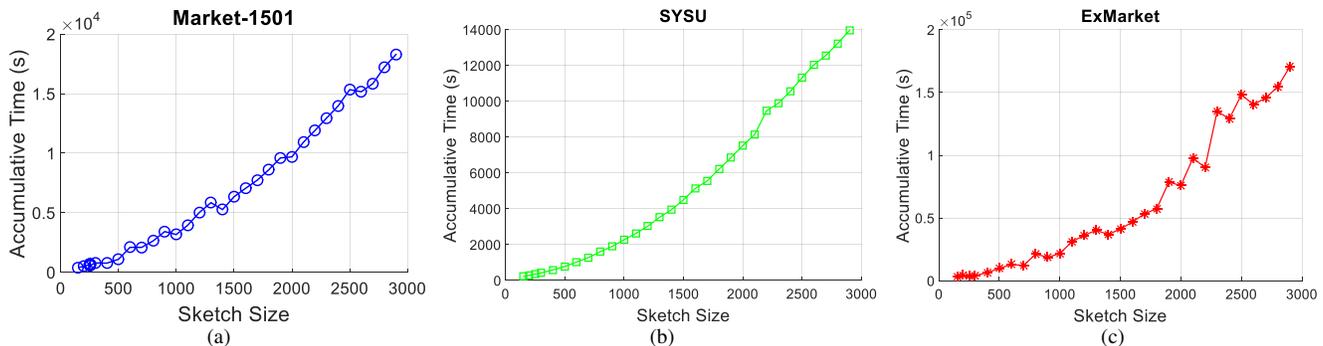


Fig. 9. Effect of the sketch size on accumulative time consumption. (Best viewed in color).

TABLE VI
COMPARISON WITH INCREMENTAL FDA MODELS AND ONLINE METHOD USING JLH.

Dataset	Market-1501			SYSU			ExMarket		
	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)	rank-1 matching rate (%)	mAP (%)	Accumulative Time (s)
OL-IDM	14.43	2.48	356136.53	3.32	4.91	554908.70	10.84	0.70	> 1 week
IDR/QR	36.70	13.73	251934.68	15.80	12.82	220962.28	39.64	10.85	2479401.99
IFDA	61.19	30.36	203537.09	21.65	18.23	189960.96	56.24	23.46	2032679.17
Pang's IFDA	71.64	45.15	204406.03	56.31	49.60	189897.24	64.64	34.80	2036601.02
SoDA	75.27	49.82	12952.07	64.81	59.82	9951.20	66.18	37.11	164475.67

TABLE IX
COMPARISON WITH OFFLINE RE-ID MODELS ON EXMARKET USING JLH(%).

Method	rank-1	rank-5	rank-10	rank-20	mAP
CRAFT	54.51	69.39	75.56	80.94	24.26
MLAPG	50.21	65.29	70.90	77.20	25.63
KISSME	57.42	69.71	74.23	78.83	30.03
XQDA	55.05	68.02	73.10	77.73	28.36
SoDA	66.18	78.36	82.48	86.64	37.11

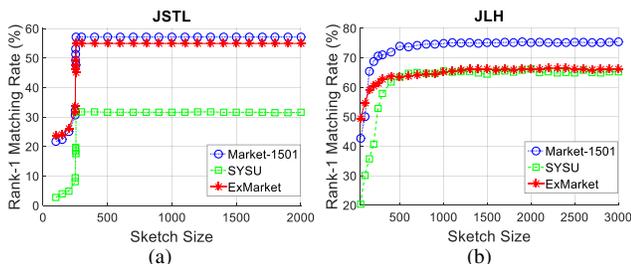


Fig. 10. Effect of the sketch size on rank-1 Matching Rate. (Best viewed in color).

TABLE VII
COMPARISON WITH OFFLINE RE-ID MODELS ON MARKET-1501 USING JLH (%).

Method	rank-1	rank-5	rank-10	rank-20	Map
CRAFT	71.20	87.35	391.69	94.39	44.24
MLAPG	69.33	85.63	90.23	93.82	46.16
KISSME	67.99	83.67	88.93	92.79	39.79
XQDA	67.96	83.91	88.95	93.14	43.89
SoDA	75.27	89.28	92.70	95.22	49.82

TABLE VIII
COMPARISON WITH OFFLINE RE-ID MODELS ON SYSU USING JLH(%).

Method	rank-1	rank-5	rank-10	rank-20	Map
CRAFT	24.70	43.03	55.11	67.73	23.31
MLAPG	18.46	35.86	47.01	58.83	18.03
KISSME	62.28	79.81	86.06	90.31	56.23
XQDA	64.14	80.88	86.85	91.90	59.12
SoDA	64.81	80.74	87.25	91.77	59.82

SYSU and spends much less time, where for instance SoDA gains around 60% reduction on the cost of computation time, as compared with Pang's ILDA.

D. SoDA vs. Related Person re-id Models

Comparison with online re-id model. We compared the online re-id method OL-IDM [37] that addresses the same setting as ours in this work. Table III, IV, V and VI tabulate the comparison results. It is noteworthy that our SoDA obtains much more stable results on rank-1 matching rate and mAP performance. Moreover, SoDA is more efficient than OL-IDM, taking 30 times smaller accumulative time.

Comparison with related subspace model and classical models. We also compared two related subspace model for person re-identification: 1) CRAFT [5]; 2) MLAPG [23], and two classical methods: 1) KISSME [16]; 2) XQDA [22], when the JLH feature was applied on all datasets. All of these methods were learned in an offline way, and the results of these methods on all benchmarks using JLH features are presented in Table VII, VIII and IX. Among all compared methods, the rank-1 matching rate and mAP of SoDA are the highest, and its accumulative time is the lowest. This indicates that SoDA achieves better or comparable performance of the related offline subspace person re-id models.

E. Further Evaluation of SoDA

We report the performance of SoDA in Figure 10 and Figure 9 when varying two key parameters ℓ .

Effect of the sketch size ℓ using low dimensional feature.

On all benchmarks, we conducted experiments using JSTL feature (256-dimensional) for evaluating the effect of the sketch size ℓ on low dimensional feature. The experimental results in Figure 10(a) indicate that the performance of our proposed SoDA can be improved when ℓ (i.e. the rank of \mathbf{B}) is larger. That is the performance is better when more variations of passed data are remained in the sketch matrix. It is reasonable because when more data variations are reserved, the estimated within-class covariance matrix from the sketch matrix \mathbf{B} can approximate the ground-truth one better. However, larger ℓ indeed increases the accumulative time since the computation complexity and memory depend on ℓ when the number of samples and the dimensionality of features are determined (Sec. III-D). Fortunately, we empirically find that good performance and low accumulative time can be achieved at the same time when setting the rank of the sketch matrix \mathbf{B} to a properly small value, i.e. $\ell = d = 256$.

Effect of the sketch size ℓ using high dimensional features.

We also show the effect of ℓ when using high dimensional features, as some recent proposed state-of-the-art person re-id features are of high dimension, such as LOMO (26960-dimensional), HIPHOP (84096-dimensional) and also the JLH (111312-dimensional) formed in this work. High dimensionality will increase the computational and space complexities (e.g., the whole training data matrix of ExMarket is a 112351×111312 matrix). Instead of conducting another online learning for dimension reduction, SoDA utilizes a set of orthogonal frequent directions maintained by the sketch matrix \mathbf{B} for reducing feature dimension. The experimental results shown in Figure 10(b) and Figure 9 again verify that increasing the sketch size ℓ can improve the performance of SoDA but also increase the accumulative time due to extra computation for dimension reduction. Also, on high dimensional feature, setting ℓ to be a properly small value (e.g. $\ell = 1000$) can gain a good balance between good performance and low accumulative computation time.

VI. CONCLUSION

We contribute to developing a succinct and effective on-line person re-identification (re-id) methods namely SoDA. Compared with existing online person re-id models, SoDA performs one-pass online learning without any explicit storage of passed observed data samples, meanwhile preserving a small sketch matrix that describes the main variation of passed observed data samples. And moreover, SoDA is able to be trained on streaming data efficiently with low computational cost, upon on no elaborated human feedback. Compared with the related online FDA models, we take a novel approach by embedding sketch processing into FDA, and we approximately estimate the within-class variation from a sketch matrix and finally derive SoDA for extracting discriminant components. More importantly, we have provided in-depth theoretical analysis on how the sketch information affects the discriminant

component analysis. The rigorous upper and lower bounds on how SoDA approaches its offline model (i.e. the classical Fisher Discriminant Analysis) are given and proved. Extensive experimental results have clearly illustrated the effectiveness of our SoDA and verified our theoretical analysis.

ACKNOWLEDGEMENT

This research was supported by the NSFC (No. 61472456, No. 61573387, No. 61522115).

REFERENCES

- [1] E. Ahmed, M. Jones, and T. K. Marks. An improved deep learning architecture for person re-identification. In *CVPR*, 2015.
- [2] G. Chechik, V. Sharma, U. Shalit, and S. Bengio. Large scale online learning of image similarity through ranking. *JMLR*, 11(Mar):1109–1135, 2010.
- [3] Y.-C. Chen, W.-S. Zheng, and J. Lai. Mirror representation for modeling view-specific transform in person re-identification. In *IJCAI*, 2015.
- [4] Y.-C. Chen, W.-S. Zheng, J.-H. Lai, and P. Yuen. An asymmetric distance model for cross-view feature mapping in person re-identification. *TCSVT*, 2016.
- [5] Y.-C. Chen, X. Zhu, W.-S. Zheng, and J.-H. Lai. Person re-identification by camera correlation aware feature augmentation. *TPAMI*, 2017.
- [6] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *JMLR*, 7(Mar):551–585, 2006.
- [7] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani. Person re-identification by symmetry-driven accumulation of local features. In *CVPR*, 2010.
- [8] S. Furoo and O. Hasegawa. An incremental network for on-line unsupervised classification and topology learning. *NN*, 19(1):90–106, 2006.
- [9] D. Gray and H. Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*, 2008.
- [10] K. Hiraoka, K.-i. Hidai, M. Hamahira, H. Mizoguchi, T. Mishima, and S. Yoshizawa. Successive learning of linear discriminant analysis: Sanger-type algorithm. In *ICPR*, 2000.
- [11] L.-K. Huang, Q. Yang, and W.-S. Zheng. Online hashing. *TNNLS*, 2017.
- [12] P. Jain, B. Kulis, I. S. Dhillon, and K. Grauman. Online metric learning and fast similarity search. In *ANIPS*, 2009.
- [13] X.-Y. Jing, X. Zhu, F. Wu, X. You, Q. Liu, D. Yue, R. Hu, and B. Xu. Super-resolution person re-identification with semi-coupled low-rank discriminant dictionary learning. In *CVPR*, 2015.
- [14] T.-K. Kim, J. Kittler, and R. Cipolla. On-line learning of mutually orthogonal subspaces for face recognition by image sets. *TIP*, 19(4):1067–1074, 2010.
- [15] T.-K. Kim, B. Stenger, J. Kittler, and R. Cipolla. Incremental linear discriminant analysis using sufficient spanning sets and its applications. *IJCV*, 91(2):216–232, 2011.
- [16] M. Koestinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof. Large scale metric learning from equivalence constraints. In *CVPR*, 2012.
- [17] M. Koestinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof. Large scale metric learning from equivalence constraints. In *CVPR*, 2012.
- [18] H. Li, Y. Li, and F. Porikli. Deeptrack: Learning discriminative feature representations online for robust visual tracking. *TIP*, 25(4):1834–1848, 2016.
- [19] X. Li, C. Shen, A. Dick, Z. M. Zhang, and Y. Zhuang. Online metric-weighted linear representations for robust visual tracking. *TPAMI*, 38(5):931–950, 2016.
- [20] X. Li, W.-S. Zheng, X. Wang, T. Xiang, and S. Gong. Multi-scale learning for low-resolution person re-identification. In *ICCV*, 2015.
- [21] J. Liang, Q. Hu, W. Wang, and Y. Han. Semisupervised online multikernel similarity learning for image retrieval. *TMM*, 19(5):1077–1089, 2017.
- [22] S. Liao, Y. Hu, X. Zhu, and S. Z. Li. Person re-identification by local maximal occurrence representation and metric learning. In *CVPR*, 2015.
- [23] S. Liao and S. Z. Li. Efficient psd constrained asymmetric metric learning for person re-identification. In *ICCV*, 2015.
- [24] E. Liberty. Simple and deterministic matrix sketching. In *SIGKDD, KDD '13*, 2013.
- [25] C. Liu, C. Change Loy, S. Gong, and G. Wang. Pop: Person re-identification post-rank optimisation. In *ICCV*, 2013.
- [26] G.-F. Lu, J. Zou, and Y. Wang. Incremental complete lda for face recognition. *PR*, 45(7):2510–2521, 2012.

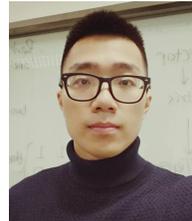
- [27] L. Ma, X. Yang, and D. Tao. Person re-identification over camera networks using multi-task distance metric learning. *TIP*, 23(8):3656–3670, 2014.
- [28] N. Martinel, A. Das, C. Micheloni, and A. K. Roy-Chowdhury. Re-identification in the function space of feature warps. *TPAMI*, 37(8):1656–1669, 2015.
- [29] N. Martinel, A. Das, C. Micheloni, and A. K. Roy-Chowdhury. Temporal model adaptation for person re-identification. In *ECCV*, 2016.
- [30] A. Mignon and F. Jurie. Pcca: A new approach for distance learning from sparse pairwise constraints. In *CVPR*, 2012.
- [31] S. Paisitkriangkrai, C. Shen, and A. van den Hengel. Learning to rank in person re-identification with metric ensembles. In *CVPR*, 2015.
- [32] R. Panda, A. Bhuiyan, V. Murino, and A. K. Roy-Chowdhury. Unsupervised adaptive re-identification in open world dynamic camera networks. In *CVPR*, 2017.
- [33] S. Pang, S. Ozawa, and N. Kasabov. Incremental linear discriminant analysis for classification of data streams. *TSMCB*, 35(5):905–914, 2005.
- [34] Y. Peng, S. Pang, G. Chen, A. Sarrafzadeh, T. Ban, and D. Inoue. Chunk incremental idr/qr lda learning. In *IJCNN*, 2013.
- [35] B. Prosser, W.-S. Zheng, S. Gong, T. Xiang, and Q. Mary. Person re-identification by support vector ranking. In *BMCV*, 2010.
- [36] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015.
- [37] Y. Sun, H. Liu, and Q. Sun. Online learning on incremental distance metric for person re-identification. In *RB*, 2014.
- [38] M. Uray, D. Skocaj, P. M. Roth, H. Bischof, and A. Leonardis. Incremental lda learning by combining reconstructive and discriminative approaches. In *BMVC*, 2007.
- [39] H. Wang, S. Gong, X. Zhu, and T. Xiang. Human-in-the-loop person re-identification. In *ECCV*, 2016.
- [40] M. K. Warmuth and D. Kuzmin. Randomized online pca algorithms with regret bounds that are logarithmic in the dimension. *JMLR*, 9(Oct):2287–2320, 2008.
- [41] A. R. Webb. *Statistical pattern recognition*. 2003.
- [42] P. Wu, S. C. Hoi, P. Zhao, C. Miao, and Z.-Y. Liu. Online multi-modal distance metric learning with application to image retrieval. *TKDE*, 28(2):454–467, 2016.
- [43] T. Xiao, H. Li, W. Ouyang, and X. Wang. Learning deep feature representations with domain guided dropout for person re-identification. In *CVPR*, 2016.
- [44] F. Xiong, M. Gou, O. Camps, and M. Szaier. Person re-identification using kernel-based metric learning methods. In *ECCV*, 2014.
- [45] J. Yan, B. Zhang, S. Yan, Q. Yang, H. Li, Z. Chen, W. Xi, W. Fan, W.-Y. Ma, and Q. Cheng. Immc: incremental maximum margin criterion. In *SIGKDD*, 2004.
- [46] J. Yang, A. F. Frangi, J.-y. Yang, D. Zhang, and Z. Jin. Kpca plus lda: a complete kernel fisher discriminant framework for feature extraction and recognition. *TPAMI*, 27(2):230–244, 2005.
- [47] H. Yao, S. Zhang, D. Zhang, Y. Zhang, J. Li, Y. Wang, and Q. Tian. Large-scale person re-identification as retrieval.
- [48] J. Ye, Q. Li, H. Xiong, H. Park, R. Janardan, and V. Kumar. Idr/qr: an incremental dimension reduction algorithm via qr decomposition. *TKDE*, 17(9):1208–1222, 2005.
- [49] L. Zhang, T. Xiang, and S. Gong. Learning a discriminative null space for person re-identification. In *CVPR*, 2016.
- [50] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, and Q. Tian. Mars: A video benchmark for large-scale person re-identification. In *ECCV*, pages 868–884. Springer, 2016.
- [51] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015.
- [52] L. Zheng, S. Wang, L. Tian, F. He, Z. Liu, and Q. Tian. Query-adaptive late fusion for image search and person re-identification. In *CVPR*, 2015.
- [53] W.-S. Zheng, S. Gong, and T. Xiang. Person re-identification by probabilistic relative distance comparison. In *CVPR*, 2011.
- [54] W.-S. Zheng, X. Li, T. Xiang, S. Liao, J. Lai, and S. Gong. Partial person re-identification. In *ICCV*, 2015.

APPENDIX

Matrix Sketch. The sketch technique we discuss in this work is related to the matrix sketch [24], which is pass-efficient to read streaming data at most a constant number of time. The sketch algorithm learns a set of frequent directions from an $N \times d$ matrix $\mathbf{X} \in \mathcal{R}^{N \times d}$ in a stream, where each row of \mathbf{X} is a d -dimensional vector. It maintains a sketch matrix $\mathbf{B} \in \mathcal{R}^{\ell \times d}$ containing ℓ ($\ell \ll N$) rows and guarantees that:

$$\mathbf{B}^T \mathbf{B} \preceq \mathbf{X}^T \mathbf{X} \quad \& \quad \|\mathbf{X}^T \mathbf{X} - \mathbf{B}^T \mathbf{B}\| \leq 2\|\mathbf{X}\|_f^2 / \ell. \quad (25)$$

Such a sketch processing is light in both processing time (bounded by $\mathcal{O}(d\ell^2)$) and space (bounded by $\mathcal{O}(\ell d)$).



Wei-Hong Li is currently a postgraduate student majoring in Information and Communication Engineering in School of Electronics and Information Technology at Sun Yat-sen University. He received the bachelor's degree in intelligence science and technology from Sun Yat-Sen University in 2015. His research interests include person re-identification, object tracking, object detection and image-based modeling.
Homepage: <https://weihonglee.github.io>.



Zhuowei Zhong is a student from Sun Yat-sen University under the joint supervision program of the Chinese University of Hong Kong. He is now graduated and received BSc degree in computer science. His research interest is in Artificial Intelligence, especially in machine learning and constraint satisfaction problem.



Wei-Shi Zheng is currently a Professor with Sun Yat-sen University. He has joined Microsoft Research Asia Young Faculty Visiting Programme. He has authored over 90 papers, including over 60 publications in main journals (TPAMI, TNN/TNNLS, TIP, TSMC-B, and PR) and top conferences (ICCV, CVPR, IJCAI, and AAAI). His recent research interests include person association and activity understanding in visual surveillance. He was a recipient of Excellent Young Scientists Fund of the National Natural Science Foundation of China, and Royal

Society-Newton Advanced Fellowship, U.K.

Homepage: <http://isee.sysu.edu.cn/~zhwshi>.