

Rotation-invariant Binary Representation of Sensor Pattern Noise for Source-Oriented Image and Video Clustering

Xufeng Lin¹ and Chang-Tsun Li^{1,2}

¹ School of Computing and Mathematics, Charles Sturt University, Australia

² Department of Computer Science, University of Warwick, UK

{xlin, chli}@csu.edu.au

Abstract

Most existing source-oriented image and video clustering algorithms based on sensor pattern noise (SPN) rely on the pairwise similarities, whose calculation usually dominates the overall computational time. The heavy computational burden is mainly incurred by the high dimensionality of SPN, which typically goes up to millions for delivering plausible clustering performance. This problem can be further aggravated by the uncertainty of the orientation of images or videos because the spatial correspondence between data with uncertain orientations needs to be reestablished in a brute-force search manner. In this work, we propose a rotation-invariant binary representation of SPN to address the issue of rotation and reduce the computational cost of calculating the pairwise similarities. Results on two public multimedia forensics databases have shown that the proposed approach is effective in overcoming the rotation issue and speeding up the calculation of pairwise SPN similarities for source-oriented image and video clustering.

1. Introduction

The ease of creating and sharing images and videos has radically changed the way we communicate and interact with each other. People capture photos or video clips of memorable moments with their smartphones or digital cameras and share them instantly on social media with family, friends or even strangers. However, the difficulty in determining the provenance of multimedia data on social media also gives rise to cyber crimes such as Internet defamation and child pornography. With a set of images and videos at hand, e.g. collected from social media websites, a forensic investigator is often faced with the task of clustering them into a number of groups, each including the data acquired by the same source camera. This task is referred to as *source-oriented image and video clustering*.

One effective approach to source-oriented image and

video clustering is based on sensor pattern noise (SPN) [14], which has been proved to be a unique and reliable fingerprint of each individual camera. The underlying rationale is that if the presence of the same SPN signal is detected in two images or videos, they are deemed to be captured by the same device and thus can be clustered into the same group. Many SPN-based image clustering algorithms have been proposed, e.g. iterative updating algorithm [3], Markov random field based methods [9, 10], spectral clustering based method [13], hierarchical clustering based method [4], normalized cuts based clustering algorithm [1], sliding window based algorithm [5], correlation clustering based algorithm [15], and the hybrid algorithm specifically designed for large-scale databases [12].

Most of the above-mentioned algorithms rely on the pairwise similarities between SPNs. However, SPN is a feeble noise-like signal, therefore, to deliver plausible clustering performance, the dimensionality of SPN has to be very high, typically one million for images undergo little or no compression. The high dimensionality imposes a heavy burden on the calculation of the pairwise similarities. Moreover, when the rotation issue presents in the dataset, which is fairly common when images and videos with different orientations are downloaded from social media, the situation can be further aggravated. This is because the similarity between two SPNs is expressed in terms of pixel-wise normalized cross-correlation (NCC) and the rotation will desynchronize the spatial correspondence between SPNs. Consequently, the images or videos taken by the same camera will be considered to be taken by different cameras. One solution to this problem is to exhaustively search for the maximal similarity among different orientations, which, however, considerably increases the computational cost because it needs to be done for every pair of images or videos. To address the above issues, a rotation-invariant and compact representation of SPN is certainly more desirable. Therefore, in this work, we propose a rotation-invariant binary representation of SPN to reduce the computational cost for the SPN-based source-oriented image and video clustering.

The remainder of this manuscript is organized as follows. In Section 2, we will describe the proposed representation of SPN in detail. Experimental results and conclusions are provided in Section 3 and Section 4, respectively.

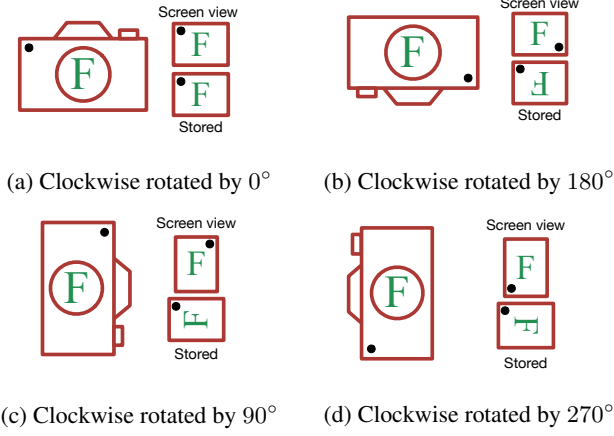


Figure 1: Four most common orientations of a digital camera (sensor) when capturing the image of “F”. For each orientation, the image stored in the camera and the image correctly displayed on screen are shown in the bottom-right panel and the top-right panel, respectively. The position corresponding to the left-upper corner of the sensor is marked by a black dot.

2. Proposed method

2.1. Rotation-Invariant Representation of SPN

Modern digital cameras are usually equipped with an orientation sensor that can sense which way the camera is held. The orientation information is recorded in the “Orientation” (for images) or “Rotation” (for videos) tag of metadata, so the image-viewing or video-playing software can later interpret the metadata and display the image or video correctly. Since the orientation of image and video is handled in a similar way in camera, we will only use the image as the example to demonstrate the rotation issue. Four most common orientations of images are shown in Fig. 1. An interesting fact about Fig. 1 is that for most consumer digital cameras, no actual in-camera image rotation is performed and the image would be stored in the landscape orientation regardless of camera orientation. As a result, the spatial correspondence between images (and thus SPNs) taken in different orientations is actually unaffected, which is favorable to SPN-based applications. Taking SPN-based image clustering for example, it is often based on the NCC similarity:

$$\rho(\mathbf{x}, \mathbf{y}) = \left\langle \frac{\mathbf{x} - \bar{\mathbf{x}}}{\|\mathbf{x} - \bar{\mathbf{x}}\|}, \frac{\mathbf{y} - \bar{\mathbf{y}}}{\|\mathbf{y} - \bar{\mathbf{y}}\|} \right\rangle, \quad (1)$$

where \mathbf{x} and \mathbf{y} are two SPNs (the mean value is denoted with a bar), $\|\cdot\|$ is the L_2 norm and $\langle \cdot, \cdot \rangle$ is the inner product. Since the calculation in Eq. (1) is pixel-wise, the consistent spatial correspondence of \mathbf{x} and \mathbf{y} is essential for accurately measuring the similarity. However, things can be messed up if the images are rotated, e.g. manually rotated by users using desktop image editing software or automatically rotated by image editing tools on social media[†], which can result in various image orientations and destroy the spatial correspondence between images.

The uncertainty of orientation can be reduced by simply rotating the vertically-oriented images clockwise by 90° to make all the images horizontally oriented. By doing so, the images can either be synchronized or desynchronized by 180° (see Fig. 1a and 1b). Assuming all the images are horizontally oriented, we illustrate how to reconstruct SPN in a rotation-invariant manner in Fig. 2. From the outside in, the pixels are decomposed into rectangular “rings” around the edge of the image. Each rectangular ring is divided into the upper-right (denoted by x'_i in Fig. 2) and bottom-left (denoted by x''_i) parts. Rotating the image by 180° will swap the positions of x'_i and x''_i , so to make the SPN rotation-invariant, we add up x'_i and x''_i in corresponding positions and scaled by $\frac{1}{\sqrt{2}}$, i.e. $x_i = \frac{1}{\sqrt{2}}(x'_i + x''_i)$. Suppose the rotation-sensitive SPN has been standardized to zero mean and unit variance, i.e. $E(x'_i) = E(x''_i) = 0, D(x'_i) = D(x''_i) = 1$, then the scaling factor $\frac{1}{\sqrt{2}}$ ensures that the rotation-invariant SPN $E(x_i) = 0, D(x_i) = 1$, which simplifies Eq. (1) to the inner product $\langle \mathbf{x}, \mathbf{y} \rangle$. If denoted by ρ^{ri} is the similarity between rotation-invariant SPNs and by ρ^{rs} is the similarity between rotation-sensitive SPNs, we have

$$\begin{aligned} \rho^{ri} &= \frac{2}{d} \sum_{i=1}^{d/2} x_i y_i = \frac{1}{d} \sum_{i=1}^{d/2} (x'_i + x''_i) (y'_i + y''_i) \\ &= \frac{1}{d} \sum_{i=1}^{d/2} (x'_i y'_i + x''_i y''_i) + \frac{1}{d} \sum_{i=1}^{d/2} (x'_i y''_i + x''_i y'_i) \\ &= \rho^{rs} + \frac{1}{d} \sum_{i=1}^{d/2} (x'_i y''_i + x''_i y'_i). \end{aligned} \quad (2)$$

x'_i and y''_i , as well as x''_i and y'_i , are independent as SPN is pixel-dependent. According to Central Limit Theorem (CLT), $\frac{1}{d} \sum_{i=1}^{d/2} (x'_i y''_i + x''_i y'_i) \rightarrow \mathcal{N}(0, 1/d)$. Therefore, we have

$$\rho^{ri} - \rho^{rs} \sim \mathcal{N}(0, 1/d). \quad (3)$$

To see how well ρ^{rs} is preserved in ρ^{ri} , let us consider two rotation-sensitive SPNs with a length of $d=2 \times 10^6$, then the probability that ρ^{ri} falls within the range of $[\rho^{rs} \pm 0.0014]$ is 95%.

[†]This is to ensure the correct image orientation after the removal of metadata for the purpose of privacy protection.

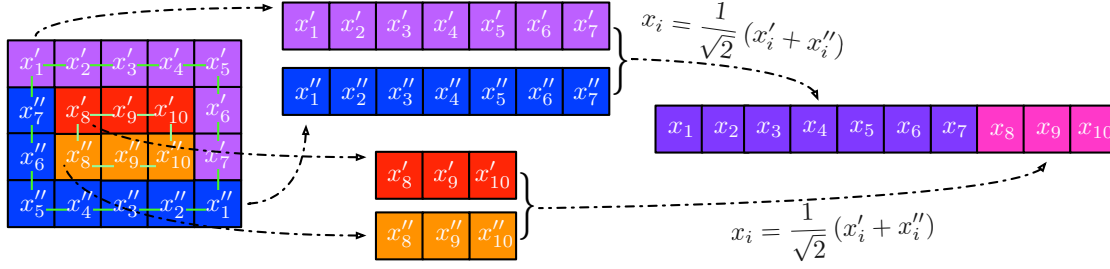


Figure 2: Demonstration of constructing a rotation-invariant SPN.

2.2. Binary Representation of SPN

Further speed up can be achieved by applying various SPN dimensionality reduction methods [2, 7, 11, 12, 16, 18], among which binarization [2] is very attractive due to its simplicity and effectiveness. However, its effect on the similarity measurement as well as the performance of source-oriented image and video clustering has not yet been well studied. Following the work of Bayram *et. al* [2], we represent two SPNs $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ in their binary forms $\hat{\mathbf{x}}, \hat{\mathbf{y}} \in \{+1, -1\}^d$ with

$$\hat{x}_i = \begin{cases} +1, & x_i > 0 \\ -1, & x_i \leq 0, \end{cases} \quad \hat{y}_i = \begin{cases} +1, & y_i > 0 \\ -1, & y_i \leq 0. \end{cases} \quad (4)$$

We are particularly interested in the effect of binarization on the similarity in Eq. (1). For simplicity, we assume that the entries $(x_i, y_i), i = 1, 2, \dots, d$ are d samples independently drawn from a bivariate normal distribution with mean $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$ and covariance matrix $\begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix}$, i.e.

$$p(x, y) = \frac{1}{2\pi\sqrt{1-r^2}} e^{-\frac{x^2+y^2-2rxy}{2(1-r^2)}}, -1 < r < 1, \quad (5)$$

where r^\S is the underlying correlation of random variables x and y . Typically, $r = 0$ for the SPNs of two different cameras and $r > 0$ for the SPNs of the same camera. Let us first examine the expectation and variance of xy :

$$\begin{cases} E(xy) = \iint xy p(x, y) dx dy = r \\ D(xy) = \iint (xy)^2 p(x, y) dx dy - E^2(xy) = 1 + r^2 \end{cases} \quad (6)$$

Therefore, according to CLT, when $d \rightarrow \infty$,

$$\rho = \frac{1}{d} \sum_{i=1}^d x_i y_i \sim \mathcal{N}\left(r, \frac{1+r^2}{d}\right). \quad (7)$$

[§] r is the true correlation value of x and y , and may not equal to the correlation value calculated from two specific SPNs using Eq. (1).

Similarly, for the binarized version $\hat{x}\hat{y}$

$$\begin{aligned} E(\hat{x}\hat{y}) &= \iint \hat{x}\hat{y} p(x, y) dx dy \\ &= 2 \int_0^{+\infty} \int_0^{+\infty} p(x, y) dx dy - 2 \int_0^{+\infty} \int_{-\infty}^0 p(x, y) dx dy \\ &= \frac{2 \arcsin(r)}{\pi}. \end{aligned} \quad (8)$$

$$\begin{aligned} D(\hat{x}\hat{y}) &= E(\hat{x}^2 \hat{y}^2) - E^2(\hat{x}\hat{y}) \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x, y) dx dy - E^2(\hat{x}\hat{y}) \\ &= 1 - \frac{4 (\arcsin(r))^2}{\pi^2}. \end{aligned} \quad (9)$$

When $d \rightarrow \infty$, we have

$$\hat{\rho} = \frac{1}{d} \sum_{i=1}^d \hat{x}_i \hat{y}_i \sim \mathcal{N}\left(\mu(r), \sigma^2(r)\right), \quad (10)$$

where

$$\begin{cases} \mu(r) = \frac{2 \arcsin(r)}{\pi} \\ \sigma^2(r) = \frac{\pi^2 - 4 (\arcsin(r))^2}{\pi^2 d}. \end{cases} \quad (11)$$

Eq. (7) and Eq. (10) show how the binarization changes the distribution of similarities. It is worth mentioning that Eq. (10) gives an analytical solution for the mean and a different variance compared to the numerical result $\hat{\rho} \sim \mathcal{N}(r/1.57, (1-r^2)/2.46d)$ in [2]. We will show their difference in the experiments in Section 3.2. If the original rotation-invariant SPN is represented by 32-bit floating-point values, binarization immediately reduces the length of SPN by a factor of 32. Furthermore, the correlation similarity of two binarized SPNs can be further simplified to their Hamming distance [2]. In this work, we binarize the rotation-invariant SPNs and pack the binary bits into 64-bit unsigned integers to speed up the calculation of pairwise similarities for source-oriented image and video clustering.

3. Experimental Results

3.1. Verification of Theoretical Results

We will first verify the theoretical results in Eq. (3) and Eq. (10). To this end, we prepared a dataset consisting of 1000 images randomly chosen from those taken by 25 cameras in the Dresden Image Database [6], with each accounting for 40 images. These 25 cameras cover nearly all the popular camera brands including Canon, Casio, FujiFilm, Kodak, Nikon, Olympus, Panasonic, Pentax, Samsung, and Sony. For each image, we only considered the central 1536×2048 pixels of the green channel, which results in original rotation-sensitive SPNs of length $d=3145728$. For simplicity, we will use the abbreviations “RS”, “RI” and “BRI” for “rotation-sensitive”, “rotation-invariant” and “binarized rotation-invariant”, and denote the pairwise NCC similarities calculated using RS, RI and BRI SPNs as ρ^{rs} , ρ^{ri} and $\hat{\rho}^{ri}$, respectively.

Fig. 3a and 3b show the scatter plot of ρ^{rs} and ρ^{ri} and the estimated probability density function (PDF) of $\rho^{rs} - \rho^{ri}$. As can be seen, ρ^{ri} shows a good agreement with ρ^{rs} and the PDF of their difference fits well with the theoretical distribution in Eq. (3). While for Eq. (10), r varies across different cameras, so we only showed the case of the binarized inter-camera similarities $\hat{\rho}^{ri}$ for $d=1572864$ (Fig. 3c) and $d=393216$ (Fig. 3d). In both cases, r is estimated as the average of the unbinarized inter-camera similarities ρ^{ri} . Clearly, Eq. (10) fits more accurately with the estimated PDF than the distribution given in [2].

3.2. Source-oriented Image Clustering

We then evaluate the proposed SPN representation on the task of source-oriented image clustering. We used the same 1000 images in Section 3.1, but to simulate the scenario of image rotation, we intentionally rotated 20 of the 40 images of each camera by 180° . Fig. 4a shows a series of ROC curves, each obtained by comparing the thresholds varying from -1 to 1 with the similarities calculated using SPNs of specific type and length. The corresponding areas under the ROC curve (AUC) are presented in the legend text. As can be seen, the conventional RS SPN, by its nature, is sensitive to image rotation, with an AUC of only 0.74 even for $d=3145728$. By contrast, the RI SPN accurately measures the similarities between unrotated and rotated images and gives a high AUC of 0.95 with SPNs of length 1572864. Furthermore, the binarization has little effect on the similarity accuracy when the length of SPN is sufficiently large, e.g. the AUC only drops by 2% when $d=1572864$, but binarization tends to have an increasingly negative effect on the accuracy as d decreases.

Based on the pairwise similarities, we applied our algorithm in [10] to cluster the 1000 images. The clustering quality was measured by the Adjusted Rand Index (ARI)

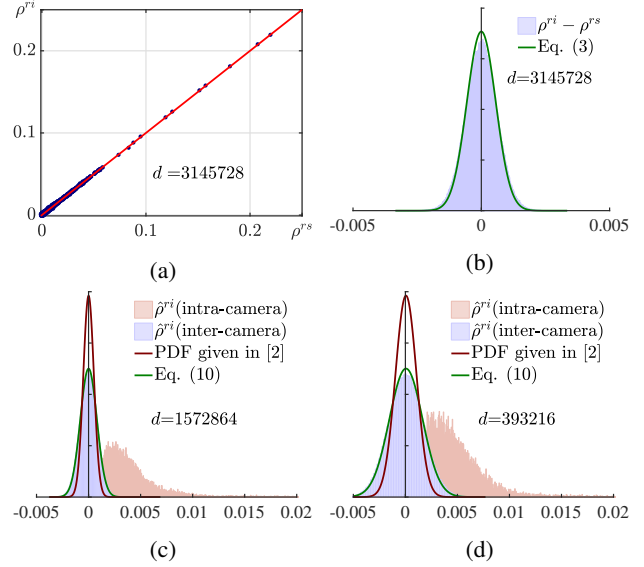


Figure 3: (a) Scatter plot of ρ^{rs} vs. ρ^{ri} for $1000 \times (1000 - 1)/2 = 499500$ pairwise similarities. (b) Estimated PDF of $\rho^{ri} - \rho^{rs}$ and the theoretical distribution in Eq. (3). (c) and (d) show the estimated PDF of the 480000 binarized inter-camera similarities and the theoretical distribution in Eq. (10).

[8], which takes on a value in $[-1, 1]$ with a higher value indicating better performance. The ARIs averaged over 50 runs were given in Table 1, which shows a good agreement with the results Fig. 4a. When $d = 1572864$, a substantially lower ARI can be observed for the conventional SPN than for the RI SPN: 0.57 versus 0.91. Similar to the trend shown in Fig. 4a, the performance degradation is more severe as d decreases, especially for the binarized SPN.

We also compared the running times of the pairwise similarity calculation for the 32-bit floating-point SPNs and the binarized SPNs, which are packed into 64-bit unsigned integers and the NCC similarity is simplified as the Hamming distance. We developed multi-threaded implementations with the C-MEX of MATLAB and reused the code as much as we can for fair comparison. The experiment was conducted on a laptop with 2.6 GHz Intel Core i7-4960HQ processor (4 cores and 8 threads) and 16 GB RAM. The running times averaged over 10 runs are shown in Fig. 4b, where the running time for binarized SPNs in Fig. 4b includes the time used for packing sign bits into 64-bit unsigned integers. As can be seen, the binarization considerably speeds up the calculation by about 40 times, which is even higher than the expected 32 times speed up. The extra speed up is benefited from the more efficient calculation of Hamming distance between integers over the inner product of floating-point values.

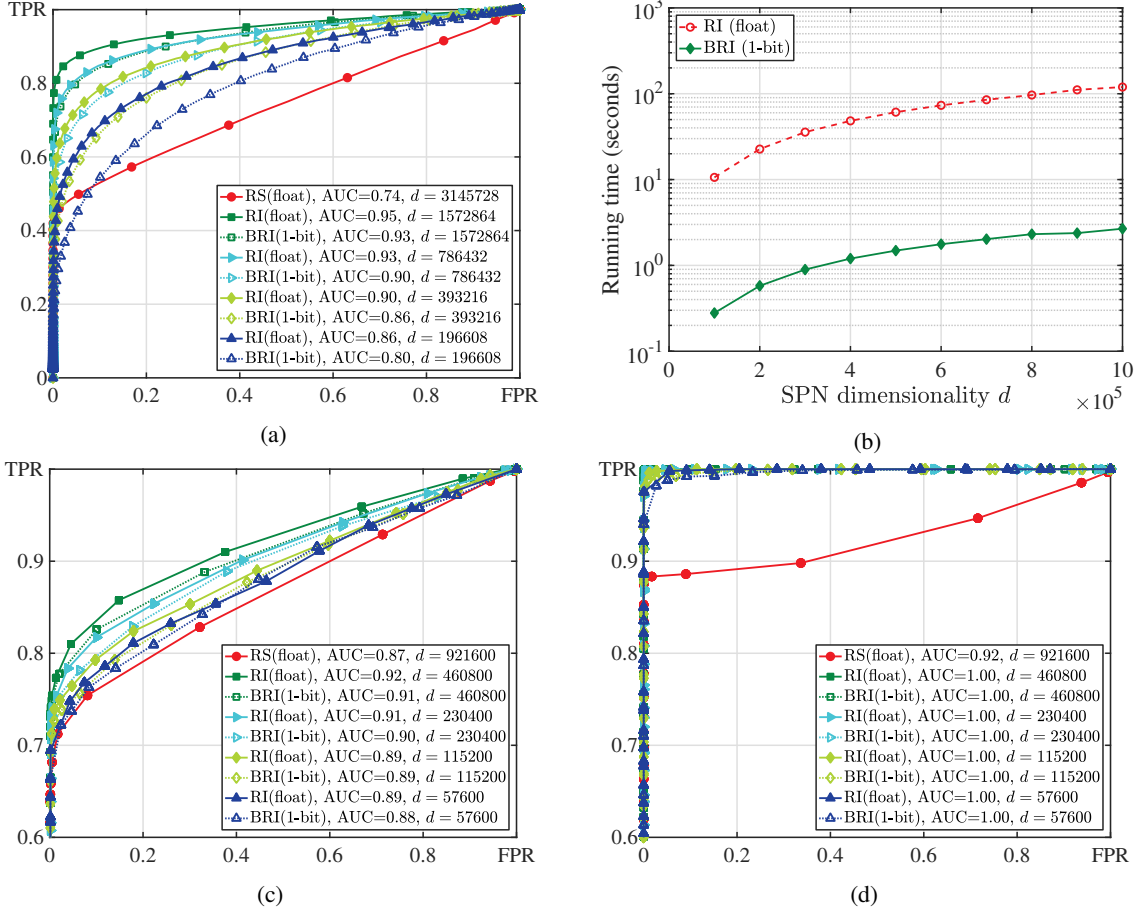


Figure 4: (a) ROC curves for images. (b) Comparison of running time (in seconds) of the floating-point and binarized SPNs. (c) ROC curves for stabilized videos. (d) ROC curves for non-stabilized videos.

Table 1: Clustering performance on images.

ARIs (images)	SPN dimensionality $d = 3145728$				
	d	$d/2$	$d/4$	$d/8$	$d/16$
RS (float)	0.57	0.57	0.54	0.48	0.40
RI (float)	—	0.91	0.88	0.85	0.75
BRI (1-bit)	—	0.90	0.84	0.68	0.41

Table 2: Clustering performance on stabilized and non-stabilized videos.

ARIs (all videos)	SPN dimensionality $d = 921600$				
	d	$d/2$	$d/4$	$d/8$	$d/16$
RS (float)	0.42	0.45	0.46	0.41	0.42
RI (float)	—	0.51	0.52	0.56	0.47
BRI (1-bit)	—	0.51	0.49	0.43	0.42

Table 3: Clustering performance on non-stabilized videos.

ARIs (videos w/o DStab)	SPN dimensionality $d = 921600$				
	d	$d/2$	$d/4$	$d/8$	$d/16$
RS (float)	0.88	0.88	0.86	0.88	0.87
RI (float)	—	0.95	0.93	0.93	0.93
BRI (1-bit)	—	0.95	0.95	0.93	0.95

3.3. Source-oriented Video Clustering

In this experiment, we aim to evaluate the proposed SPN representation on the task of source-oriented video clustering. We used 389 native indoor and outdoor videos from 32 portable devices in the VISION dataset [17]. The SPN of each video is estimated as the average of the SPNs extracted from the central 1280×720 pixels in the green channel of all video frames, which results in an original rotation-sensitive SPN dimensionality of $d = 921600$. 19 of the 389 videos have a portrait orientation and thus offer a realistic scenario for evaluating the performance of our rotation-invariant binary representation of SPN.

The ROC curves and the ARIs averaged over 50 runs for source-oriented video clustering are shown in Fig. 4c and Table 2, respectively. We note that the RI SPN delivers better performance than the RS SPN, but with only 19/389 rotated videos in the dataset, the performance gain is not as remarkable as in the image clustering scenario. Surprisingly, binarization only introduces a negligible impact on the performance even when $d = 57600$. Further investigation reveals that by averaging the noise residuals from thousands of video frames (almost all the videos in the VISION dataset last more than 1 minute), the average intra-camera similarities of most *correctly identified clusters* are higher than 0.05 and some even reach 0.2, which make them easy to be distinguished from the inter-camera similarities even after binarization (refer to Eq. (10) for the effect of binarization). It is well known that digital stabilization (DStab) can disturb the pixel-to-pixel correspondence of video frames and thus severely degrade the performance of source-oriented video clustering based on SPN. With this in mind, we excluded the stabilized videos and repeated the same experiments on the remaining 223 non-stabilized videos. Compared to the results in Fig. 4c and Table 2, significant improvements, i.e. $\sim 10\%$ AUC increase and more than 40% ARI increase on average, can be observed in Fig. 4d and Table 3.

4. Conclusion

We have proposed a rotation-invariant binary representation of sensor pattern noise to reduce the computational cost of constructing the pairwise similarities for source-oriented image and video clustering. SPN is reconstructed to overcome the rotation issue without exhaustively searching for the matching orientation between images or videos. Further speedup is achieved by binarizing the floating-point SPNs and packing the binary bits into 64-bit unsigned integers. Experiments on two public forensic image and video databases show that the proposed rotation-invariant binary representation effectively addresses the rotation issue and brings about 40 times speed up with only a slight performance degradation.

Acknowledgment

This work is partially funded by the EU project, Computer Vision Enabled Multimedia Forensics and People Identification (Project no. 690907; Acronym: IDENTITY).

References

- [1] I. Amerini, R. Caldelli, P. Crescenzi, A. Del Mastio, and A. Marino. Blind image clustering based on the normalized cuts criterion for camera identification. *Signal Processing: Image Communication*, 29(8):831–843, 2014.
- [2] S. Bayram, H. T. Sencar, and N. Memon. Efficient sensor fingerprint matching through fingerprint binarization. *IEEE Transactions on Information Forensics Security*, 7(4):1404–1413, Aug. 2012.
- [3] G. J. Bloy. Blind camera fingerprinting and image clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(3):532–534, 2007.
- [4] R. Caldelli, I. Amerini, F. Picchioni, and M. Innocenti. Fast image clustering of unknown source images. In *Proceedings of IEEE International Workshop Information Forensics Security*, pages 1–5, Dec. 2010.
- [5] L. Debiasi and A. Uhl. Blind Biometric Source Sensor Recognition Using Advanced PRNU Fingerprints. In *Proceedings of European Signal Processing Conference*, pages 779–783, Nice, France, 2015.
- [6] T. Gloe and R. Böhme. The ‘Dresden Image Database’ for Benchmarking Digital Image Forensics. *Journal on Digital Forensic Practice*, 3(2-4):150–159, 2010.
- [7] Y. Hu, C.-T. Li, and Z. Lai. Fast source camera identification using matching signs between query and reference fingerprints. *Multimedia Tools and Applications*, 74(18):7405–7428, Sep. 2015.
- [8] L. Hubert and P. Arabie. Comparing partitions. *Journal on Classification*, 2(1):193–218, 1985.
- [9] C.-T. Li. Unsupervised classification of digital images using enhanced sensor pattern noise. In *Proceedings of IEEE International Symposium Circuits and Systems*, pages 3429–3432, May 2010.
- [10] C.-T. Li and X. Lin. A fast source-oriented image clustering method for digital forensics. *EURASIP Journal Image Video Processing: Special Issues on Image and Video Forensics for Social Media analysis*, 1:69–84, Oct. 2017.
- [11] R. Li, C.-T. Li, and Y. Guan. Inference of a compact representation of sensor fingerprint for source camera identification. *Pattern Recognition*, 74(2):556–567, 2018.
- [12] X. Lin and C.-T. Li. Large-scale image clustering based on camera fingerprints. *IEEE Transactions on Information Forensics and Security*, 12(4):793–808, Apr. 2017.
- [13] B. Liu, H.-K. Lee, Y. Hu, and C.-H. Choi. On classification of source cameras: A graph based approach. In *Proceedings of IEEE International Workshop Information Forensics Security*, pages 1–5, Dec. 2010.
- [14] J. Lukas, J. Fridrich, and M. Goljan. Digital camera identification from sensor pattern noise. *IEEE Transactions on Information Forensics and Security*, 1(2):205–214, 2006.
- [15] F. Marra, G. Poggi, C. Sansone, and L. Verdoliva. Blind PRNU-Based Image Clustering for Source Identification. *IEEE Transactions on Information Forensics and Security*, 12(9):2197–2211, Sep. 2017.
- [16] T. F. Miroslav Goljan, Jessica Fridrich. Managing a large database of camera fingerprints. In *Proceedings of SPIE*, volume 7541, pages 1–12, 2010.
- [17] D. Shullani, M. Fontani, M. Iuliani, O. A. Shaya, and A. Piva. Vision: a video and image dataset for source identification. *EURASIP Journal on Information Security*, 2017(1):15, Oct 2017.
- [18] D. Valsesia, G. Coluccia, T. Bianchi, and E. Magli. Compressed fingerprint matching and camera identification via random projections. *IEEE Transactions on Information Forensics and Security*, 10(7):1472–1485, July 2015.