

Double JPEG Compression Detection Based on Extended First Digit Features of DCT Coefficients

Wei Hou, Zhe Ji, Xin Jin, and Xing Li

Abstract—Double JPEG compression detection is an important research topic for digital forensics. In this paper, we propose a powerful recompression detection method by extending the first digit features. Based on the analysis of the distribution of the first digits of quantized DCT coefficients, we extract the joint probabilities of the mode based first digits of the quantized DCT coefficients including value zero as the classifying features to distinguish between singly and doubly compressed images. Extensive experiments and comparisons with prior state-of-the-art demonstrate that the proposed scheme can detect the double JPEG compression effectively and outperforms the existing algorithms significantly. Moreover, our method can achieve a satisfactory classification accuracy even for the double JPEG compression with quality factor 95 followed by 50 or 55, while many previous works fail in the detection.

Index Terms—Double compression detection, digital forensics, first digit, JPEG.

I. INTRODUCTION

With the development of image processing software, common users can easily edit or improve the image content nowadays. These widely available image editing tools bring convenience to our daily life, but they also pose new challenges to digital forensics techniques. The transmitted images on the Internet can be easily tampered malevolently. Thus, it is necessary to identify the authenticity and integrity of a given image. Due to the high compression ratio and good quality, JPEG image format has been widely used in image acquisition devices and image processing software packages. Considering that double JPEG compression usually occurs when tampering JPEG images, it is of great significant to study the detection method for double JPEG compression.

If one JPEG image is decompressed to the spatial domain and then resaved with a different quantization matrix but with the same alignment of the 8×8 grid, we say that the final image is a double JPEG compressed image in this paper. The first quantization matrix Q^1 is called the primary quantization matrix and the second matrix Q^2 is the secondary quantization matrix. We say that the Discrete Cosine Transform (DCT) coefficient $F_q(u, v)$, $u, v = 0, 1 \dots 7$ is

double quantized if the two quantization steps $Q^1(u, v) \neq Q^2(u, v)$, where (u, v) represents the position of the quantized DCT coefficient.

Recently, some effective approaches have been presented for detecting double JPEG compression. Based on the abnormal artifacts such as “missing values” and “double peak” of the DCT coefficient histogram in the double compressed image, Lukáš and Fridrich [1] proposed three different methods that can estimate the primary quantization matrix from the double compressed image. Popescu [2] presented a double compression detection method based on the periodic property of the coefficient histogram and the peak artifact of the Fourier transform of the coefficients. Then, this method is improved by Mahdian *et al.* [3] via considering only the Alternating Current (AC) coefficients. Based on the generalized Benford’s law, Fu *et al.* [4] used the distribution of the first digits of all the quantized DCT coefficients to detect double compression. The detailed implementation process and the experimental results of this method were given by Li *et al.* [5], and then he proposed an improved scheme on the basis of the distribution of the first digits of the quantized DCT coefficients from individual AC modes. Chen *et al.* [6] used Markov random process to model the JPEG coefficient 2-D arrays and extracted the transition probability matrix as features to detect double compression. Based on the discontinuity and periodic artifact of the coefficient histogram, Feng *et al.* [7] extracted three kinds of features from individual mode histograms to detect JPEG recompression. Dong *et al.* [8] proposed to model the distribution of the mode based first digits of DCT coefficients using Markov transition probability matrix and utilized its stationary distribution as features for double compression detection. Based on the blocking artifacts in the spatial domain and the periodic characteristics in the frequency domain, Chen *et al.* [9] designed a method that can detect either block-aligned or misaligned JPEG recompression.

Inspired by the above works, we propose a more powerful JPEG recompression detection method based on the extended first digit features of DCT coefficients. In this paper, we first assume that the value 0 is the first digit of the coefficient with value zero, and then use the probabilities of the first digits of quantized DCT coefficients including value 0 from individual AC modes to detect doubly compressed JPEG images. Experimental results show that our method can detect double compression effectively and outperforms the existing algorithms.

The rest of this paper is organized as follows. Section II firstly analyzes the influence of double quantization on the statistical distribution of the extended first digits of AC coefficients, and then gives the feature extraction method.

Manuscript received March 24, 2013; revise June 26, 2013. This work was financially supported by the National High Technology Research and Development Program of China (“863” Program, No. 2011AA010601, No. 2011AA010603, No. 2011AA010605).

Wei Hou, Zhe Ji, and Xin Jin are with the National Computer Network and Information Security Administration Center, 100029, Beijing, China (e-mail: hw@cert.org.cn, jz@cert.org.cn, jinjin@cert.org.cn).

Xing Li is with the National Digital Switching System Engineering and Technological Research Center, 450002, Zhengzhou Henan, China (tel.: +86 371 81632704; e-mail: listarc@163.com).

The experimental results and comparisons with prior state-of-the-art schemes are shown in Section III. Finally, Section IV concludes this paper.

II. THE PROPOSED METHOD

As observed by Fu *et al.* [4], the distribution of the first digits of all quantized DCT coefficients of a singly compressed JPEG image follows the parametric logarithmic model, called generalized Benford's law. This model can be formulated as follows:

$$p(d) = N \log_{10} \left(1 + \frac{1}{s + d^q} \right), d \in \{1, \dots, 9\} \quad (1)$$

where d denotes the first digit of a quantized DCT coefficient, N is a normalization factor that makes $p(d)$ a probability distribution, and s and q denote model parameters which precisely describe the distributions for different images with different quality factors. It is well known that the DCT coefficient histogram of a certain image presents special recompression characteristics such as "missing values" and "double peak" when the image passes through double compression with different quality factors. In such a case, the distribution of the first digits no longer follows the generalized Benford's law.

The further comprehensive investigation shows that double quantization with different steps not only changes the distribution of the first digits of DCT coefficients but also changes the statistic characteristics of those coefficients with value 0 greatly. Fig. 1 shows the mean values of the probabilities of first digits of quantized DCT coefficients in the 9th frequency mode for 1000 singly compressed images and their corresponding doubly compressed images. A detailed description about the images is given in the experimental setup section. The singly compressed images only pass through once JPEG compression with quality factor 75, and the recompression images pass through double JPEG compression with quality factors 50 and 75 consecutively. It can be observed that the probability distribution of first digits 1-9 no longer obeys the generalized Benford's law. In addition, the statistic characteristic of those coefficients with value 0 changes a lot.

In order to include statistical feature information as much as possible, we extract the joint probabilities of the first digits of nonzero DCT coefficients and the coefficients with value 0 as the recompression detection features. For an arbitrary DCT coefficient x (including 0), its first digit can be calculated by the following formulation:

$$d = \left\lfloor \frac{|x|}{10^{\lfloor \log_{10}(|x| + \varepsilon) \rfloor}} \right\rfloor, d \in \{0, 1, \dots, 9\} \quad (2)$$

where $\lfloor \cdot \rfloor$ denotes downward rounding operator and ε denotes an infinitesimal value for calculating the first digit of the coefficient with value zero conveniently.

The quantized DCT coefficients in the position (u, v) for singly compressed images and the corresponding doubly

compressed images are denoted as $F_q^1(u, v)$ and $F_q^2(u, v)$, respectively. And the twice quantization steps are denoted as $Q^1(u, v)$ and $Q^2(u, v)$, respectively. As the study in [5], if $Q^1(u, v) \neq Q^2(u, v)$ and $Q^2(u, v)$ is not an integer multiple of $Q^1(u, v)$, then the first digit distribution of $F_q^1(u, v)$ will be different from that of $F_q^2(u, v)$ obviously. We call the AC mode (u, v) as distinguishable mode at the moment; otherwise, we call it as indistinguishable mode. In order to improve the performance of recompression detection as much as possible, the classifying features should include the statistical information of the coefficients in the distinguishable modes as much as possible. Therefore, we extract the joint probabilities of the mode based first digits of quantized DCT coefficients from the first 20 AC modes according to zig-zag scanning order as the final recompression detection features which are shown as follows:

$$F = \{p_i(d) | d \in \{0, 1, \dots, 9\} \text{ and } i \in \{1, 2, \dots, 20\}\} \quad (3)$$

The total dimensions of the final classifying features are $20 \times 10 = 200$. Based on the proposed extended first digit features, we can build a two-class classifier implemented using the machine-learning approaches to distinguish doubly compressed images from singly compressed ones.

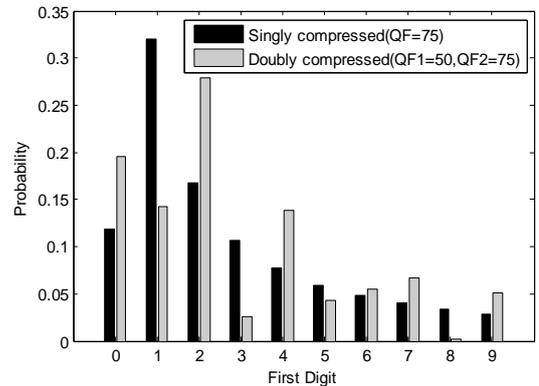


Fig. 1. Mean values of the probabilities of first digits in 9th mode for 1000 singly compressed images and their corresponding doubly compressed images.

III. EXPERIMENTAL RESULTS AND DISCUSSION

A. Experimental Setup

In order to evaluate the performance of the proposed method, we randomly select 1000 images from the image datasets BOWS2 [10], including the content of figures, landscapes, architectures, animals and plants. Some of the image samples used in the experiment is shown in Fig. 2. These images are first center-cropped into small blocks with size of 512×512 , and then converted into gray-scale images. Doubly compressed JPEG images are generated by consecutively compressing the images by a primary quality factor QF_1 and a secondary quality factor QF_2 . The corresponding singly compressed JPEG images have the same quality factor as QF_2 . Both QF_1 and QF_2 range from 50 to 95 with a step size of 5. In this way, we obtain 90 groups of

images, wherein each group including 1000 singly compressed images with QF_2 and 1000 doubly compressed images with QF_1 followed by QF_2 . Prior to the experiments, each group of images is randomly divided into two equal parts, one for training and the other for testing. In this way, it is ensured that images in the testing set used to evaluate the classifier were not used in any form during training.

It is noted that we have to construct a classifier for each secondary quality factor QF_2 due to the unknown of the primary quality factor QF_1 . Thus, we obtain 10 different two-class classifiers corresponding to each value of QF_2 in total in our experiment.



Fig. 2. Some of the image samples used in the experiment.

The support vector machine (SVM) classifier with the Gaussian kernel $k(x, y) = \exp(-\gamma \|x - y\|_2^2)$, $\gamma > 0$ is used to distinguish the doubly compressed images from the singly compressed ones. The best penalization parameter C and the kernel parameter γ are chosen by five-fold cross-validation on the training set in the following grid

$$\begin{cases} C \in \{2^i | i = -5, \dots, 15\} \\ \gamma \in \{2^j | j = -15, \dots, 3\} \end{cases} \quad (4)$$

The performance of the classifiers is evaluated by the classification accuracy, which is computed as (True Positive + True Negative)/2. To ensure the stability of the evaluation

of the proposed method, the experiments are repeated 20 times by randomly selecting the training and testing sets and the detection results are averaged over 20 times of experiments.

B. Recompression Detection Results

TABLE I: CLASSIFICATION ACCURACY (%) OF THE PROPOSED METHOD

QF ₁	QF ₂									
	5	5	6	6	7	7	8	8	9	9
	0	5	0	5	0	5	0	5	0	5
5	—	100	100	100	100	100	100	100	100	100
0	—	100	100	100	100	100	100	100	100	100
5	9	—	100	100	100	100	100	100	100	100
5	8	—	100	100	100	100	100	100	100	100
6	100	9	—	100	100	100	100	100	100	100
0	100	9	—	100	100	100	100	100	100	100
6	100	100	9	—	100	100	100	100	100	100
5	100	100	9	—	100	100	100	100	100	100
7	100	100	100	100	—	100	100	100	100	100
0	100	100	100	100	—	100	100	100	100	100
7	9	100	100	100	100	—	100	100	100	100
5	9	100	100	100	100	—	100	100	100	100
8	100	100	100	100	100	100	—	100	100	100
0	100	100	100	100	100	100	—	100	100	100
8	9	9	100	100	100	100	100	—	100	100
5	9	9	100	100	100	100	100	—	100	100
9	100	100	100	100	100	100	100	100	—	100
0	100	100	100	100	100	100	100	100	—	100
9	8	9	9	9	9	100	100	100	100	—
5	5	1	6	9	8	100	100	100	100	—

TABLE II: CLASSIFICATION ACCURACY (%) OF THE METHOD IN [5]

QF ₁	QF ₂									
	50	55	60	65	70	75	80	85	90	95
50	—	99	100	100	100	100	100	100	100	100
55	98	—	99	100	100	100	100	100	100	100
60	99	99	—	100	100	100	100	100	100	100
65	100	99	99	—	100	100	100	100	100	100
70	99	100	99	99	—	100	100	100	100	100
75	95	100	100	100	100	—	100	100	100	100
80	100	99	100	100	100	100	—	100	100	100
85	99	99	99	100	100	100	100	—	100	100
90	98	99	99	99	99	100	100	100	—	100
95	78	79	90	96	94	98	99	100	100	—

TABLE III: CLASSIFICATION ACCURACY (%) OF THE METHOD IN [7]

QF ₁	QF ₂									
	50	55	60	65	70	75	80	85	90	95
50	—	99	100	100	100	100	100	100	100	100
55	85	—	99	98	100	100	100	100	100	100
60	96	87	—	98	99	100	100	100	100	100
65	99	99	97	—	95	100	100	100	100	100
70	98	100	100	98	—	100	100	100	100	100
75	95	98	100	100	100	—	100	100	100	100
80	96	98	98	94	99	99	—	100	100	100
85	89	76	99	99	98	99	98	—	100	100
90	94	95	97	97	97	98	99	100	—	100
95	71	75	92	95	89	93	95	100	100	—

In order to show the effectiveness of the proposed method

comprehensively, we compare our method with the homogenous scheme [5] which is also designed on the basis of first digit features and the newly proposed approach [7], respectively. All these algorithms are performed on the same experimental setup. The detection results of our method and comparison methods are shown in Table I to Table III, respectively. It can be observed that all the three algorithms can detect double JPEG compression effectively when $QF_2 > QF_1$, with the accuracy of nearly 100%. However they exhibit different performances when $QF_2 < QF_1$, corresponding to the lower-left-triangle portion of the result tables. In such cases, the proposed method outperforms the other methods significantly, especially for the situation of $QF_1 = 95$.

Comparing with the method [5], the proposed method retains more statistical features by including those DCT coefficients with value zero when extracting classifying features, which resulting in a better performance. The features used in [7] are extracted from the histograms of the DCT coefficients, while that of our method are from the first digit feature space. Moreover, the first digit feature space can be seen as a nonlinear and compact map from the histogram feature space. Thus, our method contains much more information and outperforms the work [7]. Additionally, it can be seen from the detection results that all the performances of the three schemes could be further improved when $QF_1 = 95$. This is because that many quantization steps for quality factor 95 are small values, which may make Q^1 be a divisor of Q^2 easily. In such a case, some distinguishable modes will become indistinguishable ones and hence it is difficult to detect double compression.

IV. CONCLUSION

In this paper, we propose an effective double JPEG compression detection algorithm by extending the first digit features. By investigating the distributions of the first digits of quantized DCT coefficients for singly and doubly compressed images, we utilize the joint probabilities of the mode based first digits of quantized DCT coefficients including value zero to detect JPEG recompression. Extensive experiments and comparisons with prior art demonstrate that the proposed scheme outperforms the existing recompression detection methods significantly. Especially for the images compressed by quality factor 95 followed by 50 or 55, our method can still achieve a relative high classification accuracy while most of the previous works have reported to be powerless. In the future, we would like to study the application of the extended first digit features in other digital forensics purposes.

ACKNOWLEDGMENT

The authors would like to thank all the anonymous reviewers for their valuable comments.

REFERENCES

- [1] J. Luk  and J. Fridrich, "Estimation of primary quantization matrix in double compressed JPEG images," in *Proc. Digital Forensic Research Workshop*, Cleveland, Ohio, Aug. 2003.
- [2] A. C. Popescu, "Statistical tools for digital image forensics," Ph.D. Thesis, Department of Computer Science, Dartmouth College, Hanover, New Hampshire, Dec. 2004.
- [3] B. Mahdian and S. Saic, "Detecting double compressed JPEG images," *Crime Detection and Prevention (ICDP 2009), 3rd International Conference on Digital Object Identifier*, pp. 1-6, 2009.
- [4] D. Fu, Y. Q. Shi, and W. Su, "A generalized Benford's law for JPEG coefficients and its applications in image forensics," in *Proc. SPIE, Security, Steganography and Watermarking of Multimedia Contents IX*, San Jose, USA, vol. 6505, 2007, pp. 1L1-1L11.
- [5] B. Li, Y. Q. Shi, and J. Huang, "Detecting doubly compressed JPEG images by using mode based first digit features," in *Proc. MMSP, Cairns*, 2008, pp. 730-735.
- [6] C. Chen, Y. Q. Shi, and W. Su, "A machine learning based scheme for double jpeg compression detection," in *Proc. International Conference on Pattern Recognition*, 2008, pp. 1-4.
- [7] X. Feng and G. Doerr, "JPEG re-compression detection," in *Proc. SPIE*, vol. 7541, pp. 0J1-0J12, 2010.
- [8] L. Dong, X. Kong, B. Wang, and X. You, "Double compression detection based on Markov model of the first digit of DCT coefficient," in *Proc. 6th International Conference on Image and Graphics*, 2011, pp. 234-237.
- [9] Y. L. Chen and C. T. Hsu, "Detecting recompression of JPEG images via periodicity analysis of compression artifacts for tampering detection," *IEEE Trans. Inf. Forensics Security*, vol. 6, no. 2, pp. 396-406, 2011.
- [10] P. Bas and T. Furon. (July 2007). BOWS-2. [Online]. Available: <http://bows2.gipsa-lab.inpg.fr>.



Wei Hou received the B.S. degree from Harbin Engineering University, China in 2003, the M.S. degree from Harbin Institute of Technology, China in 2005, and the Ph.D. degree from Tsinghua University, in 2011.

He is currently an engineer at CNCERT/CC (National Computer Network Emergency Response Technical Team/Coordination Center of China). His research interests include cognitive radio networks and multimedia signal processing.



Zhe Ji received the B.S. degree from Harbin Institute of Technology, China, in 2006, the Ph.D. degree from Tsinghua University, in 2011.

She is currently an engineer at CNCERT/CC (National Computer Network Emergency Response Technical Team/Coordination Center of China). Her research interests include speech signal processing and network information security.



Xin Jin received the B.S. degree from USTB (University of Science and Technology Beijing) in 2005, and the Ph.D. degree from ICT/CAS (Institute of Computing Technology, Chinese Academy of Sciences) in 2010.

Now, he is the senior engineer at CNCERT/CC (National Computer network Emergency Response Technical Team/Coordination Center of China). His research interests include 4G/LTE, MIMO, wireless communication.



Xing Li received the M.S. degree in Signal and Information Processing from Zhengzhou Information Science and Technology Institute in 2012.

He is currently an assistant lecturer at National Digital Switching System Engineering & Technological Research Center. His research interests include multimedia information processing and pattern recognition.