

Analysis of Benford's Law in Digital Image Forensics

Abhay Kumar

12011

abhayk@iitk.ac.in

Abhay Agarwal

12010

asonu@iitk.ac.in

Abhishek Kumar Singh

12032

akshrt@iitk.ac.in

Abstract-In this term paper, first we have described Benford's law in general. Then we have analyzed its application in various aspects of Digital Image Forensics. We have established the results for single and doubled compressed JPEG and JPEG2000. We have also analyzed for computer generated images. Amount of deviation from the actual Benford's curve can be used as a qualitative indicator of forgery.

I. INTRODUCTION

Image Forgery is not new. Digital Image Forgery does not differ very much in nature compared to conventional Image Forgery. Instead of using photograph, digital image forgery deals with digital image. The process of creating fake image has been tremendously simple with the introduction of powerful computer graphics editing software such as Adobe Photoshop, GIMP, and Corel Draw, some of which are available for free. There are many cases of digital image forgery. From the tabloid magazines to the fashion industry and in mainstream media outlets, scientific journals, political campaigns, courtrooms, and the photo hoaxes that land in our e-mail in-boxes, doctored photographs are appearing with a growing frequency and sophistication. Over the past few years, the field of digital forensics has emerged to help restore some trust to digital

Digital image forensics can be broadly classified into two categories:-

1.) Passive v/s Active Forensics:- In passive forensic, image generation process can't be traced as the image is only in 'read only' mode. Whereas, In active forensics the generation is purposely modified so that it can leave behind the traces.

2.) Blind v/s Non-Blind Forensics:- In Blind forensics, the tampered image is investigated without having any prior knowledge about its generation process and its original scene. Whereas, the original scene can be easily detected and the intermediate steps are also known to the investigator in case of Non-Blind forensics.

Digital watermarking has been proposed as a means by which an image can be authenticated. The drawback of this approach is that a watermark must be inserted at the

time of recording, which would limit this approach to specially equipped digital cameras. In contrast to these approaches, passive techniques for image forensics operate in the absence of any watermark or signature.

II. BENFORD'S LAW

Benford's law, also known as the **first digit law** or **significant digit law**, is an empirical law. It was first discovered by Newcomb in 1881 and rediscovered by Benford in 1938. Hill gave a statistical explanation of this law. It states that the probability distribution of the first digits, x ($x = 1, 2, \dots, 9$), in a set of natural numbers is logarithmic. More specifically, if a data set satisfies Benford's law, its significant digits will have the following distribution:

$$p(x) = \log_{10} \left(1 + \frac{1}{x} \right) \quad x = 1, 2, \dots, 9$$

where $p(x)$ stands for probability of x .

Since, it's probability distribution function, it must sum up to unity.

$$\begin{aligned} \sum_{x=1,2,\dots,9} p(x) &= \log_{10} \left(1 + \frac{1}{1} \right) + \log_{10} \left(1 + \frac{1}{2} \right) + \dots + \log_{10} \left(1 + \frac{1}{9} \right) \\ &= \log_{10} \left(\frac{2}{1} \right) + \log_{10} \left(\frac{3}{2} \right) + \dots + \log_{10} \left(\frac{10}{9} \right) \\ &= \log_{10} \left(\frac{2}{1} \times \frac{3}{2} \times \dots \times \frac{10}{9} \right) \\ &= \log_{10} (10) \\ &= 1 \end{aligned}$$

HISTORY OF BENFORD'S LAW

Astronomer-mathematician *Simon Newcomb* was the first to observe this digit bias. He published the article as follows:-

"That the significant digits don't occur with equal frequency must be evident to anyone using logarithmic tables more often, and noticing how much faster the first pages wear out than the last ones. Digit '1' appears as the first significant digit more frequently than any other digit."

Frank Benford provided some justification that why this problem is worth to study. He argued that:-

“It has been observed that the pages of a much used logarithms tables shows evidences of a selective use of natural numbers. The pages containing the logarithms of the low numbers 1 and 2 are apt to be more stained and frayed by use than those of the higher numbers 8 and 9. Of course, no one could be expected to be greatly interested in the condition of logarithm table, but the matter may be considered more worthy of study when we recall that the table is used in the building up of our scientific, engineering, and general factual literature.”

INTERSTING FACTS ABOUT BENFORD’S LAW

Scale Invariance

If a set of numbers follows Benford’s law, multiplying the numbers by any possible constant will create another set of numbers that also follows Benford’s law. A system that remains unchanged when multiplied by a constant is called scale invariant

Base Invariance

If a group of numbers follows Benford’s law in one base, it will also follow Benford’s law if converted to another base. Any universal law should apply whether it is being observed by humans, with 10 fingers apiece, or by ducks with six toes. There should be nothing special about the base 10 number system. If a data set of base ‘d’ number system satisfies Benford’s law, its significant digits will have the following generalised distribution:

$$p(x) = \log_d \left(1 + \frac{1}{x} \right) \quad x = 1, 2, \dots, d - 1$$

JUSTIFICATION OF BENFORD’S LAW

It is an empirical law and there is no as such proper mathematical proof. Many writers have come to the conclusion that Benford’s law is a mysterious law of nature, for which a true explanation lies with the gods. However, one possible explanation is that the fractional part of the logarithm of the data is uniformly distributed between 0 and 1. May be nature prefers logarithmic scale and hence there is uniform distribution of fractional part of mantissa in logarithmic scale rather than in natural number system. Alternative formulation of Benford’s first significant digit distribution is:-

$$p(x) = \log_{10} (x + 1) - \log_{10} (x) \quad x = 1, 2, \dots, d - 1$$

This also suggests a systematic method to check whether a random variable X satisfies Benford’s law iff:

$$\log X - \lfloor \log X \rfloor \sim U(0,1) \text{ or equivalently if}$$

$$\log X \bmod 1 \sim U(0,1)$$

GENERAL SIGNIFICANT DIGIT LAW

$$prob(D_1 = d_1, \dots, D_k = d_k) = \log_{10} \left[1 + \left(\frac{1}{\sum_{i=1}^k d_i \times 10^{k-i}} \right) \right]$$

for all positive integers k , and all $d_i \in \{1, 2, \dots, 9\}$ and all $d_i \in \{0, 1, 2, \dots, 9\}; j = 2, \dots, k$

And the probability that d ($d = 0, 1, \dots, 9$) is encountered as the n -th ($n > 1$) digit is

$$prob(D_n = d) = \sum_{k=10^{n-2}}^{k=10^{n-1}-1} \log_{10} \left[1 + \left(\frac{1}{10k + d} \right) \right]$$

Why does one set of numbers follow Benford’s Law, while another set of numbers does not?

There is no as such general check to determine if the data-set will follow Benford’s law. Most of the dataset, free of any bias and influence, produced as a result of mathematical combination of numbers, whose mean is greater than the median and the skew is positive will follow Benford’s law. Most of the common distributions we have learned do not follow Benford’s law. One of the primary mysteries of Benford’s law has been this seemingly unpredictable behaviour.

Benford’s law curve for first significant digit:

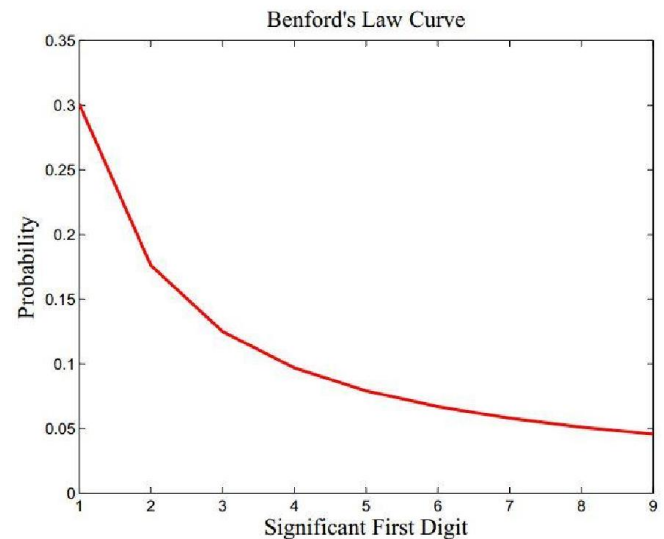


Fig 1: Benford’s Law Curve for first significant digit

In the table below, we have populated the probability of the digit appearing as the first significant digit and as the second significant digit. After the fourth significant digit onward, all digits are almost equally probable.

| Digit,d | Probability as first digit | Probability as second digit |
|---------|----------------------------|-----------------------------|
| 0 | Not applicable | 0.1197 |
| 1 | 0.3010 | 0.1139 |
| 2 | 0.1761 | 0.1088 |
| 3 | 0.1249 | 0.1043 |
| 4 | 0.0969 | 0.1003 |
| 5 | 0.0792 | 0.0967 |
| 6 | 0.0669 | 0.0934 |
| 7 | 0.0580 | 0.0904 |
| 8 | 0.0512 | 0.0876 |
| 9 | 0.0458 | 0.0850 |
| Total | 1 | 1 |

Table-1: Probability of d as first and second digit

APPLICATIONS OF BENFORD'S LAW

Benford's law has been widely exploited for various types of fraud detection. Some of its wide uses include-Accounting fraud detection, election data fraud detection, analysing tax-fraud, scientific fraud detection, digital image forensics etc.

III. APPLICATION OF BENFORD'S LAW IN DIGITAL IMAGE FORENSICS

In general, the gray scale values of pixels of images don't follow the Benford's law as it is limited to the range of [0,255]. There is limitations on the values of gray scale values. But, the DCT coefficients are better suited to follow Benford's law, as it is as a result of mathematical combinations. We have used the following standard images (cameraman.jpg and lena.jpg) for simulation purpose.

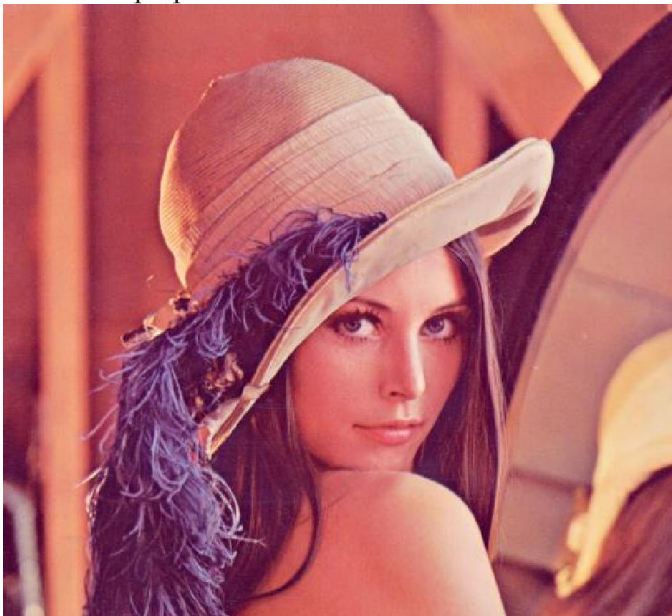


Fig 2: Lena.jpg (<http://www.eecs.qmul.ac.uk/~phao/CIP/Images>)



Fig 3: Cameraman.jpg (<http://www.eecs.qmul.ac.uk/~phao/CIP/Images>)

The system model for analysis of JPEG images is explained in the following figure:

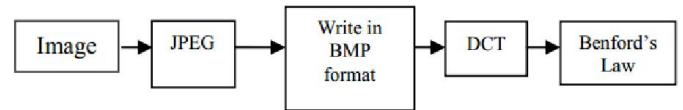


Fig 4: System model for analysis of JPEG images

IV. SIMULATION RESULTS AND OBSERVATIONS

CASE-1: Uncompressed JPEG Images

Probability distribution of first significant digit of the DCT coefficients of the images are as shown :

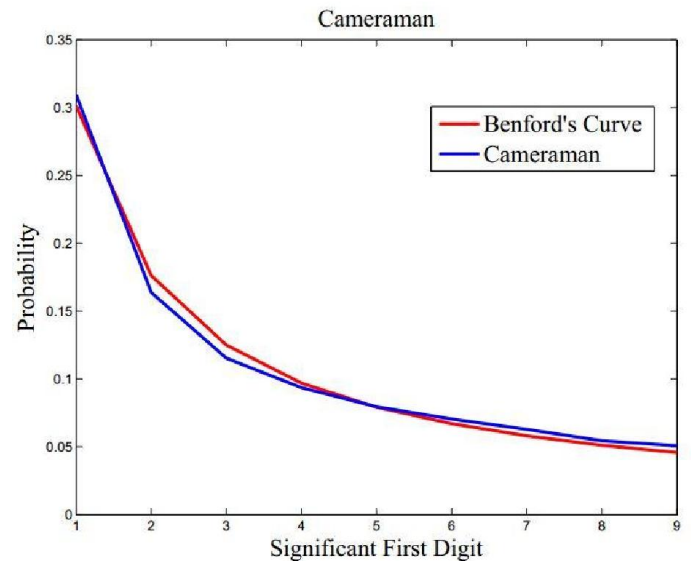


Fig 5: Probability Curve of Cameraman.jpg

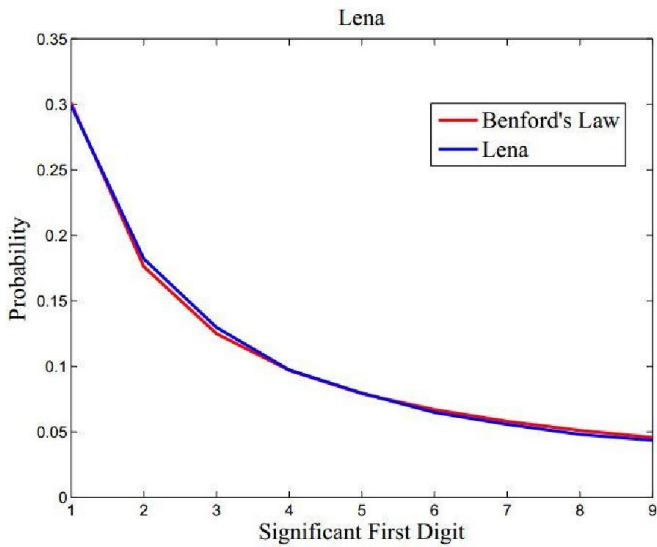


Fig 6: Probability Curve of Lena.jpg

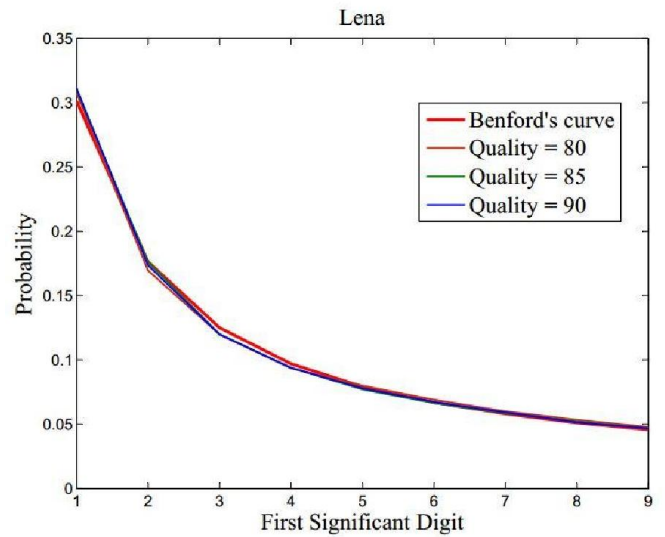


Fig 8: Probability Curves of Lena.jpg at different Quality factors

OBSERVATION & CONCLUSION:

Both the images follows Benford's curve in its original uncompressed form.

CASE-2: Single compressed JPEG Images

We have compressed the images using 'imwrite' function in Matlab with different 'Quality' factor and analyzed Benford's law for the single compressed images.

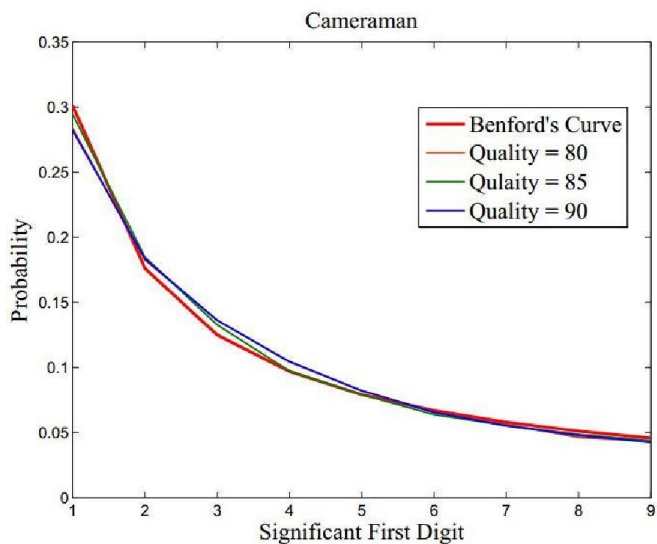


Fig 7: Probability Curves of Cameraman.jpg at different Quality factors

OBSERVATION & CONCLUSION:

Single compressed images with different 'Quality' factor almost follow Benford's Law (Fig 7 and Fig 8). Since the deviation of these curves from the actual Benford's curve is not significant enough, Benford's law can't be exploited to identify the manipulations (single compression) applied to the images.

CASE-3: Double compressed JPEG Images

We have compressed the images twice with same or different Quality factors using the 'imwrite' function in Matlab consecutively. The significant first digit of the DCT coefficients of the double compressed images have been analyzed. It is observed that double compressed images don't follow Benford's Curve at all. But, the distribution is still logarithmic in nature. There have been significant deviation from Benford's curve. We have compressed the images with all three possible cases i.e. first compression with high Quality factor and then with lower in next compression, vice versa and same Quality factors. The deviation of curves depend on the Quality factors used at different compression stages. One of the deviation measure could be Mean Squared error between probabilities of original and double compressed images

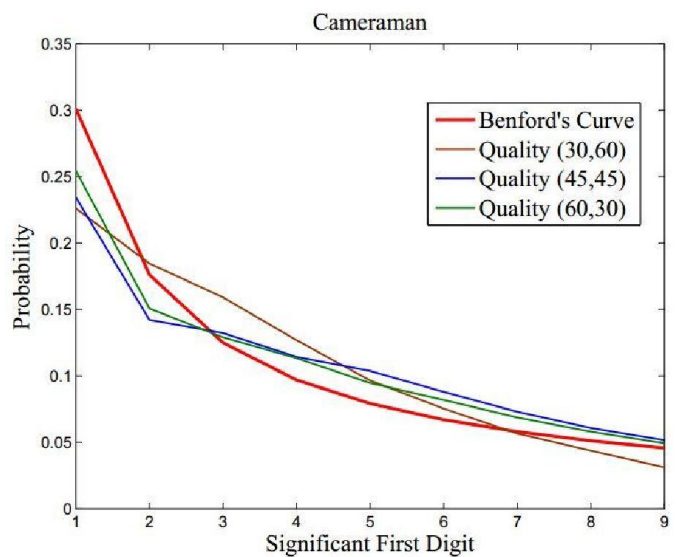


Fig9: Probability curves of double compressed Cameraman.jpg

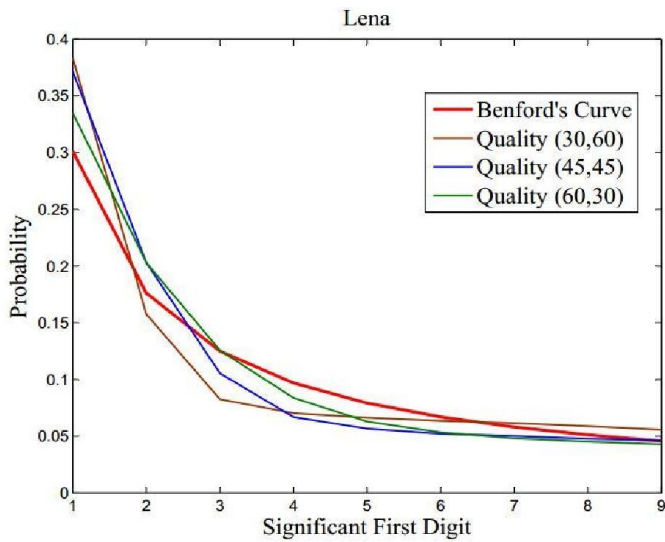


Fig10: Probability curves of double compressed Lena.jpg

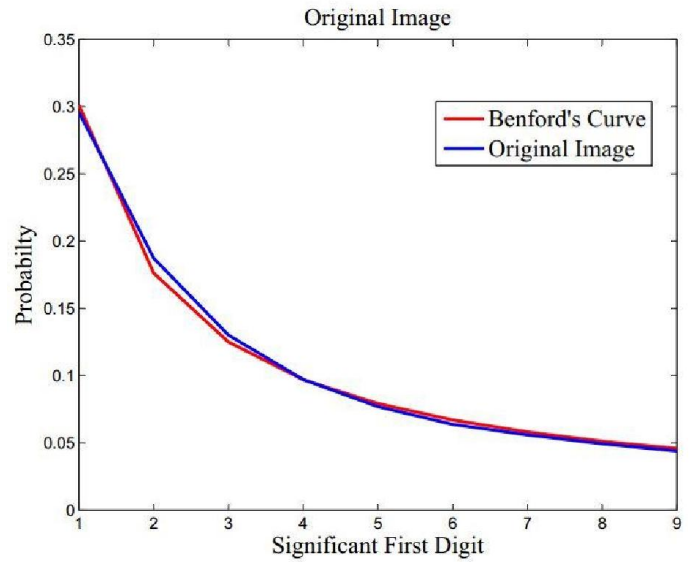


Fig12: Probability curves of original image

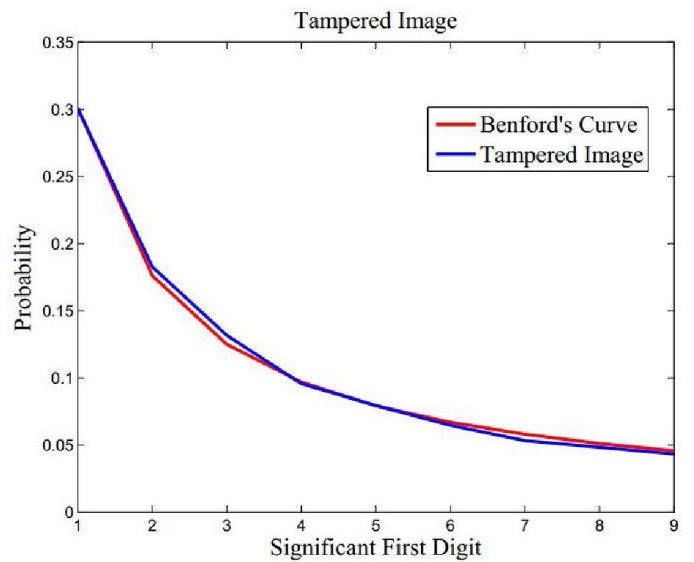


Fig13: Probability curves of original image

OBSERVATION & CONCLUSION:

Double compressed images don't follow Benford's Curve. We can't use this law for double compression detection in raw form. But, we can use deviation measure to differentiate between single and double compressed images. In fact, we can use measure of deviation as a parameter to estimate the compression.

CASE-4: Images with copy-paste forgery



Fig1: Original Image (top), Tampered Image (bottom)

OBSERVATION & CONCLUSION:

Both original and tampered images follow Benford's Curve. We can't use this law for detecting copy-paste forgery. However, taking into consideration next few digits may help detecting copy move forgery.

CASE-5: Computer Generated JPEG Images



Fig. 14. CGI1 (left) and CGI2 (right)



Fig. 15. CGI3

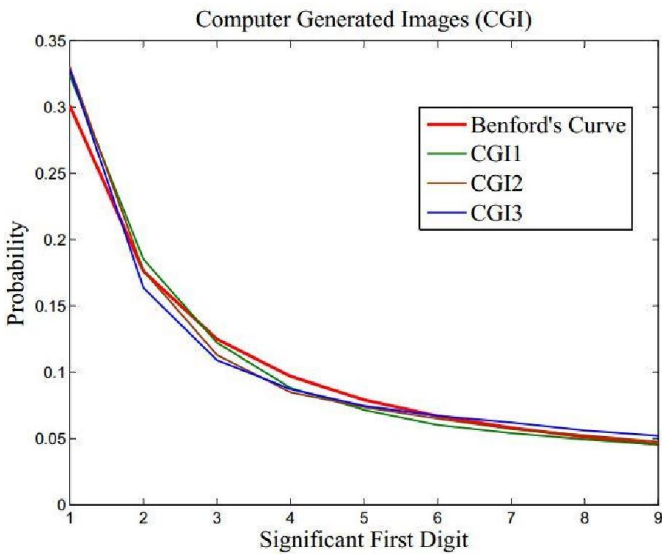


Fig. 16. Probability Curves for Computer Generated Images

OBSERVATION & CONCLUSION:

From above figure, it can be concluded that Computer Generated Images (CGI) also follow Benford's law in general. So, we can't differentiate between natural and computer generated image using Benford's law.

CASE-6: Single compressed JPEG2000 Images

JPEG2000 is a new compression standard which is based on Discrete Wavelet Transform while JPEG coding is based on the Discrete Cosine Transform. In the literature review, authors have argued that the probability of significant first digit of DWT coefficients follows Benford's law more precisely than that of DCT coefficients. But, we will use DCT coefficient itself, because the author have not clearly mentioned which out of the four coefficients or all four coefficients are used to get the probability distribution of first significant digit.

`[[cA(:,:,j),cH(:,:,j),cV(:,:,j),cD(:,:,j)]=dwt2(R(:,:,j),'haar');]. We tried using different combinations but we didn't get expected result or somewhat similar to the expected. So, we will be using DCT coefficients itself. Firstly we will write image in .jpeg2000 format using following Matlab command :-`
`imwrite(R1,'imagenam.jp2','jp2','Mode','lossy','CompressionRatio',z(i));`

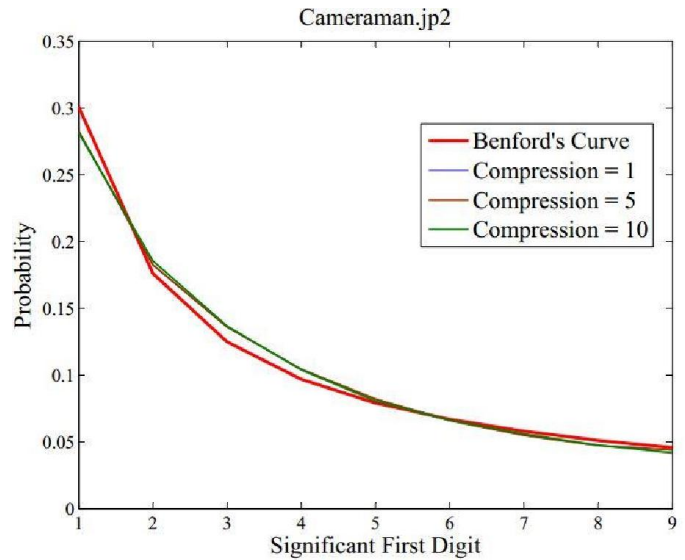


Fig 17: Probability Curves of Cameraman.jp2 at different Compression Ratios

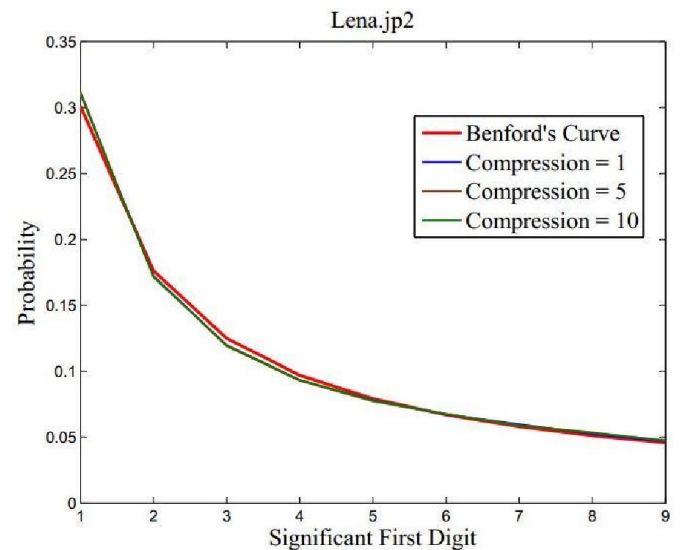


Fig 18: Probability Curves of Lena.jp2 at different Compression Ratios

OBSERVATION & CONCLUSION:

From above figure, it can be concluded that single compressed JPEG2000 images also follow Benford's law in general. So, we can't differentiate between original and single compressed images using Benford's law. However, if the deviation will be significantly high for higher compression ratio. So, we may differentiate qualitatively for different compression ratios.

CASE-7: Double compressed JPEG2000 Images

Firstly, JPEG Images are written into JPEG2000 format and then compressed twice with either same compression ratio at both stages of with different compression ratios. Then, probability distribution of the first significant digit is observed for the .jp2 images.

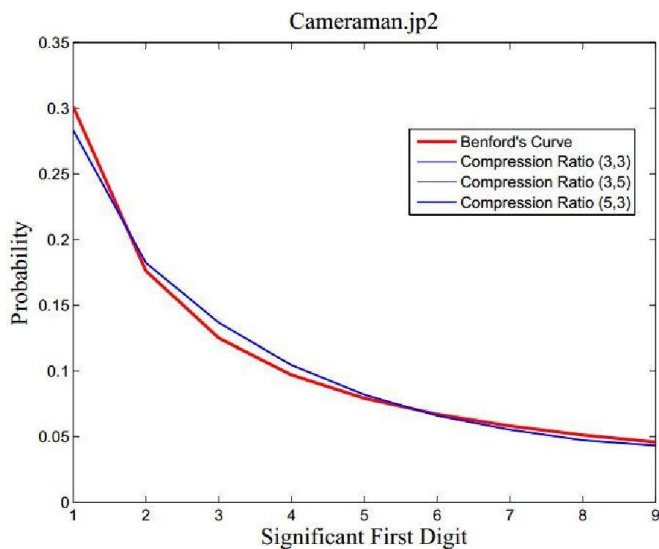


Fig 19: Probability Curves of double compressed Cameraman.jp2.

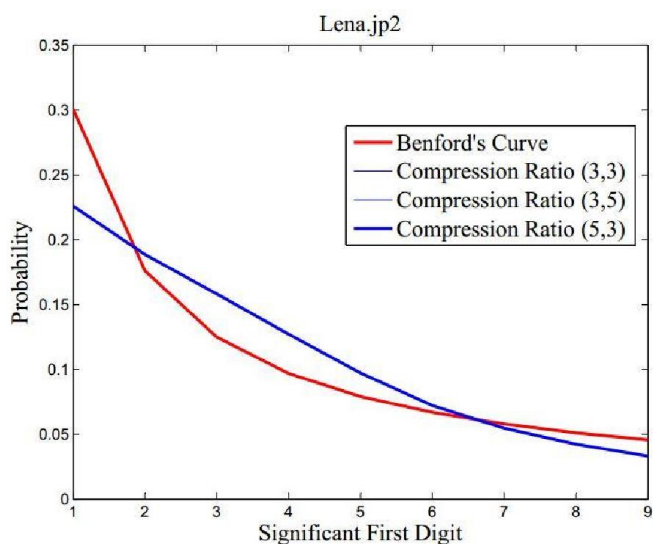


Fig 20: Probability Curves of double compressed Lena.jp2.

OBSERVATION & CONCLUSION:

From above figure, it can be concluded that there is not similar trends for compressed JPEG2000 images. It follow Benford's law for Cameraman.jp2 but not in case of Lena.jp2. So, we can't really differentiate between original and double compressed images using Benford's law. (However, the authors have given some ambiguous conclusion for double compressed .jp2 images. Please see the highlighted portion in the pdf attached)

V. CONTEMPORARY RESEARCH, CHALLENGES AND CONCLUSION

A great deal of research on Digital image forensic methods is underway as there are numerous unmet challenges. Because of the unprecedented rise in use of forged images, we need to have robust and cost-efficient

forensic methodologies. Benford's law based forensic methods seems promising, atleast in determining qualitative forgery. Generalized Benford's law for first few significant digits (rather than only the first significant digit) need to be taken into consideration to get better accuracy. Some authors have argued that Benford's Law is sometimes essentially useless as a forensic indicator of fraud. Deviations from either the first or second digit version of that law can arise regardless of whether an image is original or forged. In fact, fraud can move data in the direction of satisfying that law and thereby occasion wholly erroneous conclusions. So, an extensive study on the applicability of Benford's Law need to be done.

REFERENCES

- [1] Neetu Singh, Rishab Bansal "Analysis of Benford's Law in Digital Image Forensics"
- [2] https://books.google.co.in/books?id=J_NnBgAAQBAJ&pg=PP1&dq=Benford's%20Law%3A%20Theory%20and%20Applications&pg=PA26#v=onepage&q=Benford's%20Law:%20Theory%20and%20Applications&f=false
- [3] https://www.stat.auckland.ac.nz/~fewster/RFewster_Benford.pdf
- [4] <http://www.imperial.ac.uk/~nadams/classificationgroup/Benford's-Law.pdf>
- [5] Qadir, Ghulam, Xi Zhao, Anthony TS Ho, and Matthew Casey. "Image forensic of glare feature for improving image retrieval using Benford's Law."
- [6] https://en.wikipedia.org/wiki/Benford's_law
- [7] <http://www.eecs.qmul.ac.uk/~phao/CIP/Images/>
- [8] Fu, Dongdong, Yun Q. Shi, and Wei Su. "A generalized Benford's law for JPEG coefficients and its applications in image forensics."
- [9] A. Berger, T.P. Hill, "A basic theory of Benford's Law"
- [10] Ghulam Qadir, Xi Zhao and Anthony TS Ho, on "Estimating JPEG2000 compression For Image Forensics Using The Benford's Law"
- [11] H. Farid, "Image forgery detection"
- [12] Benford's Law Applications for Forensic, Accounting, Auditing, and Fraud Detection by MARK J. NIGRINI
- [13] "Detecting Copy-Paste Forgery in Images Using Statistical Fingerprints" Amandeep Kaur, Surbhi Gupta, and Parvinder S. Sandhu
- [14] "The Irrelevance of Benford's Law for Detecting Fraud in Elections" Joseph Deckert, Mikhail Myagkov and Peter C. Ordeshook

NOTE:

- We have included all Matlab files used for getting our simulation results. Please read "ReadMeFile.txt" to know the naming of those files.
- We have also included all the standard images used for simulations.
- We have attached the sample code in this report
- We have attached the original paper with discrepancies highlighted and commented. (This will explain more clearly the DIFFERENCES b/n OUR and THEIRS Simulations.)