# Image Compression of Grayscale Images using Principal Component Analysis

Laila Alhimale[1], Wedad Taieb Arebe[2], Kamal Ali Albashiri[3]
[1]Al Jabal Al Gharbi University
[2]Al Jabal Al Gharbi University[3], Al Jabal Al Gharbi University

*Abstract: Image compression is a particular method of compressing data to digital images in the aims of reducing the cost for storing such data or transmitting these data via telecommunication channels. Normally, image compression algorithms employ statistical characteristics or the visual perception of the image data to compress the image data without huge loss in the quality of the recovered image data after the application of such compression method[1]. In this study, we analyse the effect of compressing grayscale images using principal component analysis. The principal component analysis method is employed to decrease firstly the dimension of the image data but still keeping the important information of the image data with minimal loss of useful information. The PCA method was employed in this research in order to assess the quality of retrieved images according to the extracted principal components and also to analyse the veracity of the principal component analysis method as a good image compression method. Therefore, in this research we will observe the effect of principal component analysis on the quality of recovered images that available or commonly used in Matlab. The quality between the original and the recovered image was assessed using peak signal to noise ratio. Also the compression ratio of the recovered image was computed to compare the recovered image file size to the original image file size.In some analysed images, a low number of principal components such as PC=10 or PC=15) were sufficient to produce good quality recovered images. Also, a low number of principal components produced a highly compressed original image which thus created a low file size (in terms of bytes) of the recovered image. In our research, we demonstrated that principal component analysis method is really a powerful technique in terms of compressing a digital image, and also it helps in reducing storage costs as well as transferring information costs via communication channels.*

**Keywords:**Image compression, lossy compression, Principal Component Analysis, Grayscale images

## 1. Introduction

Digital image compression may be lossless or lossy[2]. Lossless compression are preferred methods for medical imaging, technical drawings and for archival purposes while lossy compression methods that are generally utilised at low bit rates create compression artifacts. Lossy compression methods are appropriate for natural images such as photographs were little loss of fidelity of the original image data is acceptable at the expense of reducing substantially data storage capacity[3].In this study, we study a lossy data compression method which is the principal component analysis data compression method to see its effect on the image quality of the recovered images.Next, we will describe principal component analysis method.

### 1.1 Definition and Importance of PCA

Principal component analysis is a statistical method and is used extensively in data analysis. PCA is an orthogonal linear transformation that converts the data to a new coordinate system such that the largest variance by any projection of the data rests on the first coordinate which is also called as the first principle component and the second largest variance lies on the second coordinate which is the second principal component and so forth[4]. PCA is employed in machine learning and also in signal processing as well as image compression. The main function of PCA is to reduce the dimensionality of a dataset while it retains as much important information of the dataset as possible or without losing much information. This research study uses the concept of PCA to compress grayscale image data of size 512x512 and to assess quality of recovered images based on the extracted principal components. PCA is

employed to extract the most important characteristic of the test images and these principal components are then used to compress the image and used to recover the original image.

## 1.2 Principal Component Analysis steps

The general key steps that are required to do principal component analysis on a set of data are as follows: **step 1** - it is important to ensure that the data input is in matrix form and our image data is a 2D matrix actually which conforms with step 1, and then in **step 2**, mean of the dataset is subtracted from the data itself. In order for the principal component analysis to function properly, the mean of the data should be subtracted from image data which eventually produces an image dataset whose mean is zero. The next **step 3**, the covariance matrix is calculated (See 1.3(ii))followed by **step 4** which consists of the computation of the eigenvalues and the eigenvectors of the covariance matrix. As the covariance matrix is a square matrix, one is able to determine the eigenvalues and the eigenvectors for this particular matrix (1.3(iv)). In so doing, this gives us meaningful information about the image data. However, the most important aspect of this step 4 is that it provides us with information about the patterns in the image data. By computing the process of determining the eigenvectors of the covariance matrix, one is able to extract the lines or outlines of the image data characteristics. The next **step 5**, the principal components are chosen to form a feature vector whereby the idea of data compression and reduction in data dimensionality comes into play here in this step. Actually, theeigenvector whose eigenvalue is highest is the first Principal component of the image dataset. It is in fact the most significant correlation between the data dimensions. Generally, eigenvectors are determined from the covariance matrix followed by the **last step 6** which is to sort the eigen values from highest eigenvalue to lowest eigenvalue which in the end produces the principal component in order of significance. Then, in the principal component analysis, one can ignore those components that are of least significance at the expense of information lost. Subsequently, the final image data will have fewer dimensions than the original image data[5].

## 1.3 Background of PCA and Equations

Suppose we have some attributes that are $A_1$ and $A_2$, and we have $n$ training examples. $x$'s denote values of $A_1$ and $y$'s denote values of $A_2$.

**(i)** Variance of an attribute is:

$$\mathrm{var}(A_1) = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{(n-1)} \qquad \text{...Equ.1}$$

where $\bar{x}$ is the mean of $x$

**(ii)** Mean of x

$$= \frac{\sum_{i=1}^{n} X_i \bar{X}}{n} \qquad \text{...Equ.2}$$

**(iii)** Covariance of two attributes is:

$$\mathrm{cov}(A_1, A_2) = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{(n-1)} \qquad \text{...Equ.3}$$

For instance, to interpret the covariance of two attributes, if covariance is positive, both dimensions increase together and if covariance is negative, as one increases, the other decreases whereas covariance zero means the two attributes are independent of each other.

**(iv)** Covariance matrix

Suppose we have $n$ attributes, $A_1$, ...,$A_n$, the covariance matrix is:

$$C^{n \times n} = (c_{i,j}), \text{where } c_{i,j} = \mathrm{cov}(A_i, A_j) \qquad \text{.... Equ.4}$$

**(v)** Eigenvectors

a.  Let **M** be an $n \times n$ matrix.

i.  v is an *eigenvector* of M if $M \times v = \lambda v$

ii.  $\lambda$ is called the *eigenvalue* associated with v

b.  For any eigenvector v of Mand scalar $a$,

c. Thus you can always choose eigenvectors of length 1:

$$\sqrt{v_1^{2.}+...+v_n^{2}}=1 \qquad .... \qquad \textbf{Equ.5}$$

If **M** has any eigenvectors, it has *n* of them, and they are orthogonal to one another. Thus eigenvectors can be used as a new basis for a *n*-dimensional vector space[6][7].

**(vi)** Order eigenvectors by eigenvalue, highest to lowest.

In general, one gets *n* components. To reduce dimensionality to *p*, one ignores *n−p* components at the bottom of the list.

**(vii)** Construct new feature vector. Feature vector = $(v_1,\ v_2,\ ...v_p)$ or reduced dimension feature vector of $(v_1,\ v_2... \langle v_p )$

**(viii)** Derive the new data set

TransformedData=RowFeatureVector× RowDataAdjust    … **Equ.6**

This gives original data in terms of chosen components (eigenvectors)

**(ix)** Reconstructing original data

a. RowDataAdjust=RowFeatureVector$^{-1}$×TransformedData    …**Equ.7**

b. RowDataOriginal=RowDataAdjust + OriginalDataMean    … **Equ.8**

**1.4 Peak Signal-to-Noise ratio (PSNR)**

Normally, peak signal to noise ratio is the ratio between the maximum power of a signal and the power of noise corrupting that signal quality. PSNR is expressed using logarithmic decibel scale (dB). PSNR is commonly utilised to measure the quality of reconstruction of lossy compression codecs such as for image compression. In this particular case, the original data is the original image and the noise is the "errors" introduced by that compression method.

PSNR is an approximation to our perception as human in the reconstruction quality. A higher PSNR generally indicates that the reconstruction is of higher quality and one has to be careful while comparing the results [8]. PSNR is defined using the mean square error (MSE) and given a noise free *m*×*n* monochrome image *I* and its noisy approximation *K*,

*MSE* is defined as: RMSE = √MSE    … **Equ.9**

where RMSE is a root mean squared error

The PSNR (in dB) is defined as:

$$PSNR = 10\log_{10}\left(\frac{MAX^2}{MSE}\right) = 20\log_{10}\left(\frac{MAX}{RMSE}\right)$$

…
**Equ.10**

Here, $MAX_i$ is the maximum possible pixel value of the image. When the pixels are represented using 8 bits per sample, this is 255[9][10].

**1.5 Compression Ratio**

According to Castro[11][12], low-loss compression which can be achieved by this principal component analysis method can be expressed in terms of the compression ratio (ρ) and of the mean squared error (MSE) committed in the approximation of *A* (original image) by *Ã* (image obtained from the disposal of some of the components). The compression factor is defined by:

$$\rho = \frac{\text{Unit of memory required to represent } \tilde{A}}{\text{Unit of memorial required to represent } A} \qquad … \textbf{ Equ. 11}$$

**2. Research Methodology**

In this study, we investigate the quality of image compression using different number of principal components from 5 different grayscale images imported from Matlab software. The 5 images were (i) Baboon.png, (ii) Peppers.png, (iii) Fruits.png, (iv) Barbara.png, and (v) Lena.png. They were all of size 512 x 512 pixel and were in PNG format. All the RGB images were converted to its corresponding grayscale images for principal component analysis. The quality of compression was assessed using Peak signal-to-noise ratio (PSNR) of the original image and the compressed image based on a distinct number of principal components of that particular image.
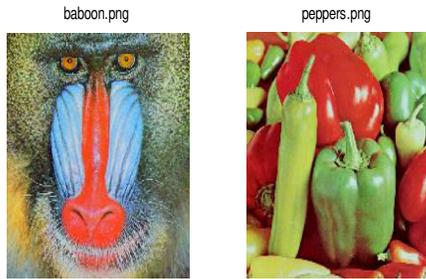
### 2.1 Original images



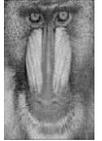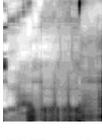Figure 1a: Baboon and Peppers (RGB images were of size 512x512 pixels)



Figure 1b: RGB Fruits.png, Grayscale Barbara.png, RGB lena.png all of size 512x512 pixels.

## 3. Results

### 3.1 Grayscale images and their corresponding histogram plots

All the 512x512x3 RGB original images were converted to 512x512 grayscale images. The image intensity for each grayscale image was computed. The intensity of an image is the average of the intensities of all the pixels found in that particular image. The lesser the value of the image intensity, it means the image pixel becomes darker with 0 being black and 255 being white in terms of image pixel intensity. The computed average image intensity for each image was as follows: for baboon the image intensity was found to be 129.6, that for peppers was 120.2, for fruits was 164.8, for Barbara was 117.3 and finally for Lena was 95.6.

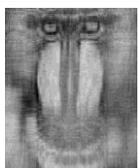### 3.2 Recovered images from using a certain number of principal components (PCs).

| PC= 5 | PC=10 | PC=15 | PC=20 | PC=25 |
|---|---|---|---|---|
|  |  |  |  |  |
| PSNR = 18.3 dB | PSNR = 19.0 dB | PSNR = 19.5 dB | PSNR = 19.9 dB | PSNR = 20.3dB |
|  |  |  |  |  |
| PSNR = 18.3 dB | PSNR=21.0 dB | PSNR=22.8 dB | PSNR = 24.3 dB | PSNR = 25.4 dB |
|  |  |  |  |  |
| PSNR = 20.7 dB | PSNR = 22.5dB | PSNR = 23.9 dB | PSNR = 25.0 dB | PSNR = 25.9 dB |

| | | | | |
|---|---|---|---|---|
| PSNR = 18.2 dB | PSNR = 20.1 dB | PSNR = 21.2 dB | PSNR = 22.0 dB | PSNR = 22.7 dB |
| PSNR = 20.7 dB | PSNR = 23.0 dB | PSNR = 24.4 dB | PSNR = 25.5dB | PSNR = 26.3 dB |

**Figure 2: Compressed and Recovered images using a varying number of PCs (5 to 25)**

Figure 2 shows the recovered images from the five original images using a minimum number of principal components of 5 to a maximum number of principal components of 25 and incrementing by 5. Some recovered images visually can be easily recognised even at low principal components used to reconstruct the original images. Also for each recovered image, the peak signal to noise ratio was noted and this was displayed alongside each recovered image.

### 3.3 Original images, Compressed images and Recovered images

| Original Grayscale Image | Compressed Image (Number of Principal Components=10) | Recovered grayscale image |
|---|---|---|
|  File size =637,192 bytes | File size = 198,943 bytes Compression ratio =31.2% | File size = 388,197 bytes PSNR=19.0 dB |



| | | |
|---|---|---|
| File size = 453,846 bytes | File size=131,247 bytes Compression ratio =28.9% | File size = 279,133 bytes PSNR = 23.0 dB |

**Figure 3: Compression ratio of the compressed image using a certain number of principal components and the recovered image reconstructed from the compressed image.**

Figure 3 depicts the original grayscale images as well as the corresponding compressed image using 10 principal components (for instance) and the recovered or reconstructed grayscale image. Then the compression ratio was computed in terms of new file size (bytes) to original file size (bytes). It was found that using 10 principal components of the original grayscale such as image baboon, the compression ratio was 31.2% and that of Lena, the compression ratio was 28.9% which is really a great reduction in storage capacity.

### 4. Discussion

This research described the principal analysis method firstly on a population of digital image data and its possibility to the compression of digital images. The importance of the application of the technique in communication engineering and also as a storage capacity method in compressing digital images was emphasized. The principal component analysis was

applied to 5 grayscale images of the same size of 512x512. The average image pixel intensity of the analysed images varied from 95.6 to 164.8. From the PCA analysis, it was observed that while recovering baboon image from its compressed image with 15 principal components, the recovered baboon image is practically visually similar to that of the quality of images recovered with greater number of principal components employed that were 20 and 25.Same observations were found with other digital images whatever their image intensity. Therefore, even with low number of principal number of components, the original grayscale image is recovered successfully with good peak signal to noise ratio. In our study, the PSNR value proved to be a good indicator of image quality of recovered image as compared to the original image. The advantage of using low number of principal components to compress the original image is that it takes much lower storage space as shown by the compression ratio values which were very promising and thus fast transfer of image data information via communication links could be achieved with minimal loss. The relevance of this research work is in the performance evaluation of the PCA formulation in compressing digital images from the measurement of the degree of compression, and the degree of information loss that the PCA introduces into the compressed images in discarding some principal components.

## 5. Conclusion:

A low number of principal components is required to compress the original images and also to recover the original images after the compression process. Principal component analysis is a powerful technique in terms of taking the main key components of a particular image. This key feature is very important for the extraction of the important characteristics of a digital image. This in turn allows the fast transfer of images through communication channels and also it requires less storage capacity and hence lower storage costs and also communication bandwidth cost of transferring data images. Future works should investigate thoroughly and incrementally the effect of principal component on peak signal to noise ratio on smaller image size (less than 512x512) as well as larger image sizes (greater than 512x512). Also

future works should see the application of PCA directly to coloured (RGB) images.

## References

[1] Li, Drew, Lossy Compression Alghorithms. Prentice Hall, 2003.
[2] ] Subramanya A, "Image Compression Techniques ".Potentials IEEE, Vol.20, Issue 1, Feb-March 2001.
[3] ChristophStamm , A new progressive file format forlossy and lossless image compression, 2002.
[4] https://coolstatsblog.com, principal component analysis, Retrieved on 16/03/2017.
[5] Ashok et al. Principal Component Analysis Based Image Recognition International Journal of Computer Science and Information Technologies, Vol. 1 (2) , 2010, 44-50.
[6] Jeff Jauregui, Principal component analysis with linear algebra, August, 2012.
[7] Liton P, Abdulla S, Face Recognition Using Principal Component Analysis Method, ISSN: 2278 – 1323 International Journal of Advanced Research in Computer Engineering & Technology (IJARCET) Volume 1, Issue 9, November 2012.

[8]Poobathy D, ManickaR,"Edge Detection Operators: Peak Signal to Noise Ratio Based Comparison", IJIGSP, vol.6, no.10, pp.55-61, 2014.DOI: 10.5815/ijigsp.2014.10.07
[9] Huynh-Thu, Q,Ghanbari, M. (2008). "Scope of validity of PSNR in image/video quality assessment". *Electronics Letters*. **44** (13): 800. doi:10.1049/el:20080522.
[10] Salomon, David (2007). *Data Compression: The Complete Reference* (4 ed.). Springer. p. 281. ISBN 978-1846286025.

[11] Castro M, Algoritmoherbiano generaliza doparaextração dos componentesprincipais de um conjunto de dados no domíniocomplexo [dissertação]. Porto Alegre: Pontifícia Universida de Católica do Rio Grande do Sul; 1996.    [ Links ]

[12] Castro M, Castro FC. Codificação de sinais. 2008.
Disponívelem: http://www.ee.pucrs.br/~decastro/download.html. 2008.