## PCA (Principal Component Analysis)

- Image compression
- Removal of uncorrelated phenomena (irregularities) in the image
- Methods for anomaly detection in the image
  - Localization of uncorrelated phenomena by filtering correlated phenomena using PCA (X-Xpca => occurrence of anomalies in the image)
  - Localization of uncorrelated phenomena by filtering correlated phenomena using PCA (Xpca => occurrence of anomalies in the image)

## Note (SVD ... Singular Value Decomposition)

The SVD method (decomposition) is analogous and is used for the same purposes as SVD in signal and image processing. SVD is recommended for sparse matrices (images), while PCA is recommended in the opposite cases.

## ∨ PCA

### Compression - Reducing the number of columns in matrix/image $\mathbf{X}$

Image

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_M \end{bmatrix}$$

$$= \begin{bmatrix} x_1(k=1) & x_2(k=1) & \cdots & x_M(k=1) \\ x_1(k=2) & x_2(k=2) & \cdots & x_M(k=2) \\ \vdots & \vdots & \cdots & \vdots \\ x_1(k=N) & x_2(k=N) & \cdots & x_M(k=N) \end{bmatrix} \ldots N \times M$$

Column Centering (typically Z-Score here)

$$\mathbf{x}_i \leftarrow \frac{\mathbf{x}_i - \overline{\mathbf{x}_i}}{\sigma_i}; i = 1 \ldots M$$

Covariance Matrix (matrix of correlation coefficients between columns of $\mathbf{X}$)

$$\mathbf{covX} = \frac{1}{N}\mathbf{X}^T \cdot \mathbf{X}.$$

Eigenvalues and Eigenvectors

$$d, \mathbf{V} = \text{eig}(\mathbf{covX}),$$

where

$$\mathbf{V} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \ldots & \mathbf{v}_M \end{bmatrix} \text{ where } \mathbf{v}_i \ldots M \times 1.$$

In the columns of $\mathbf{V}_{PCA}$ are the eigenvectors of the covariance matrix.

We choose the number of most significant eigenvectors $m << M$.

Column Compression:

$$\mathbf{X}_{compressed} = \mathbf{X} \cdot \mathbf{V}_{[:,1:m]} \ldots N \times m.$$

Decompression

$$\mathbf{X}_{PCA} = (\mathbf{X} \cdot \mathbf{V}_{[:,1:m]}) \cdot \mathbf{V}_{[:,1:m]}^T \ldots N \times M,$$

where $\mathbf{V}_{[:,1:m]}^T$ is the pseudoinverse.

Then, transpose the matrix $\mathbf{X}_{PCA}$, and similarly compress the original image through rows...

PCA pro odfiltrování nekorelovaných anomálií pomocí $n$ a $m$ nejvýznamnějších vlastních vektorů kovarianční matice (zachovávají se lineárně korelované=opakující se jevy).



Or vice versa: PCA for filtering correlated phenomena (for compression and decompression, we do not use $n$ and $m$ most significant eigenvectors of the covariance matrix).
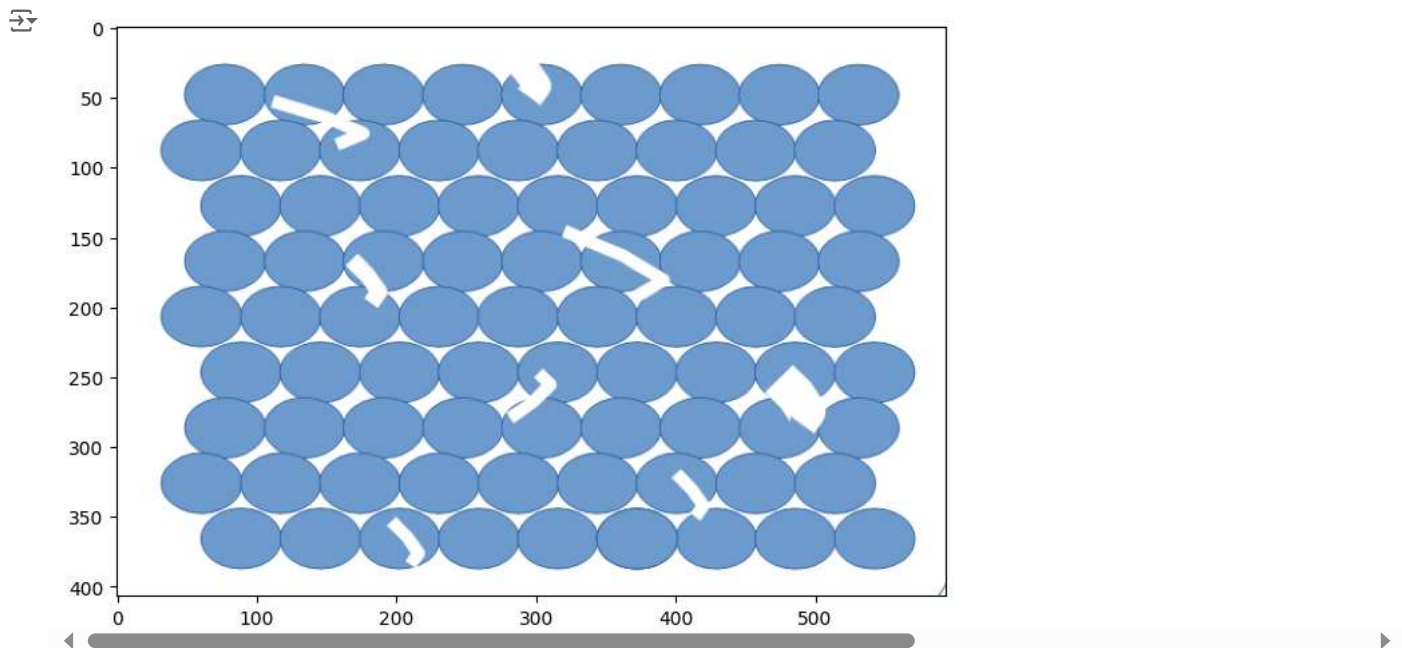
Similarly:





```
1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib.image as mpimg
5 import PIL as pil
6 img1="./img/regular_patttern_with_perturb.png"
7
8
9 im1 = pil.Image.open(img1)
10
11
12 plt.figure(figsize=(8,8))
13 plt.imshow(im1)
14 plt.show()
```
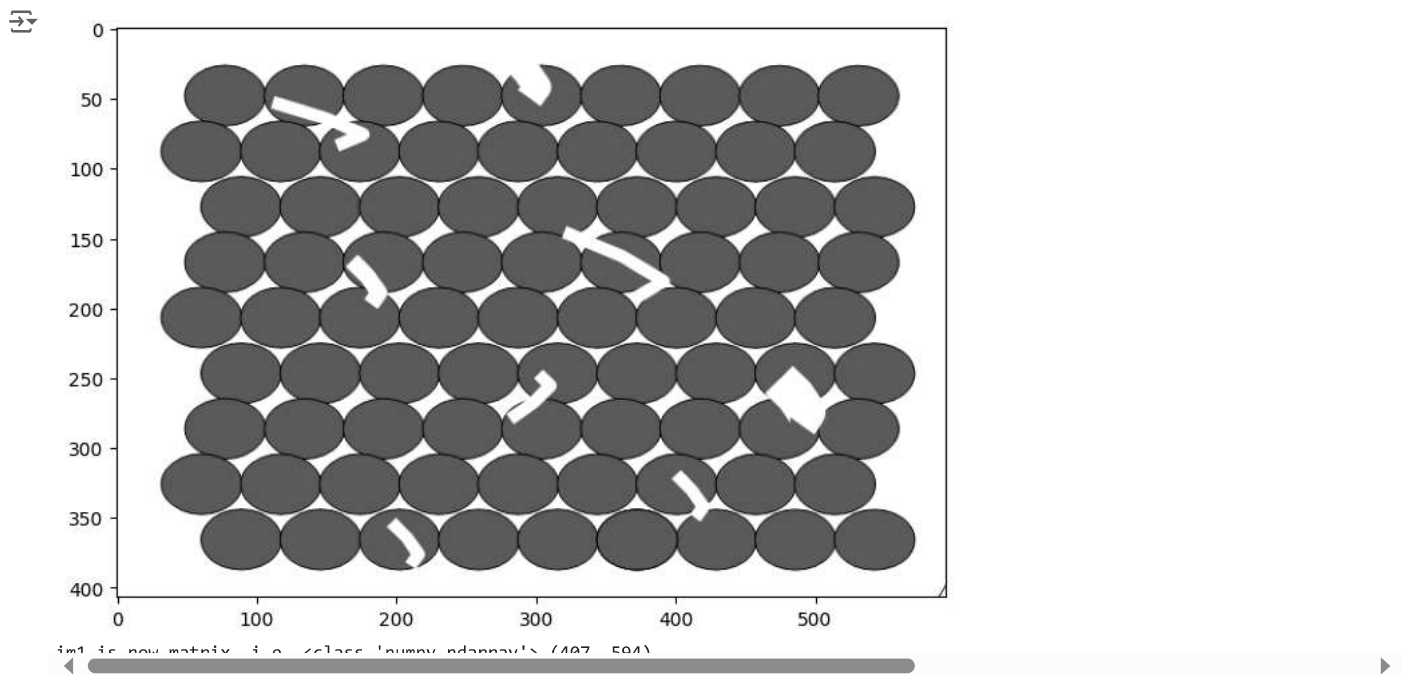


```
1 from scipy import signal
2 im1=np.mean(im1,2)  # gray img, 2-D matrix,...X
```

```
1 plt.figure(figsize=(8,8))
2 plt.imshow(im1,cmap='gray')
3 plt.show()
4
5 print("im1 is now matrix, i.e.", type(im1), im1.shape)
6
```
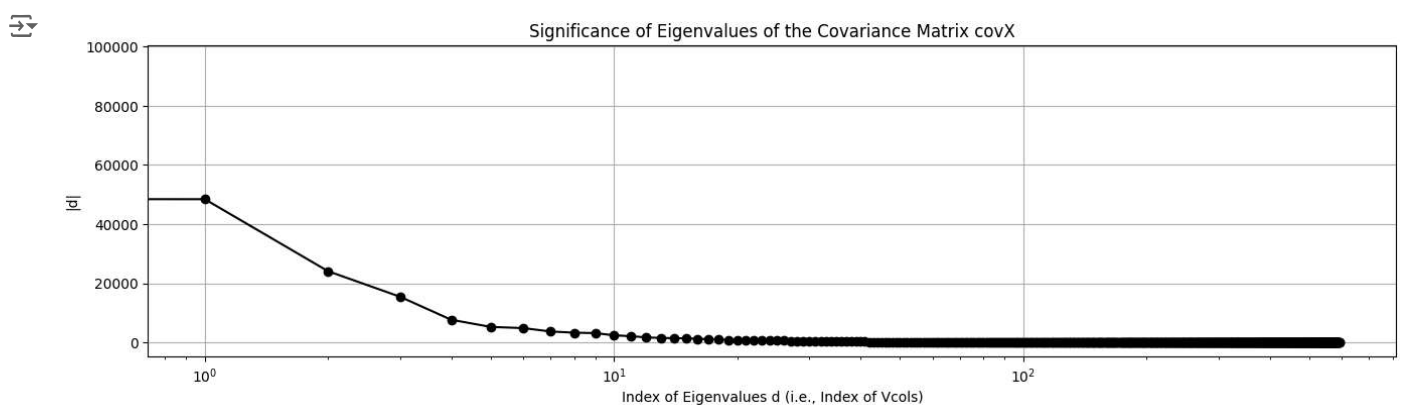
im1 is now matrix, i.e. <class 'numpy.ndarray'> (407, 594)

## Localization of uncorrelated events by filtering correlated events using PCA (X-Xpca => occurrence of anomalies in the image)

i.e. filtering correlated events => occurrence of anomalies

```
 1 X = im1.copy()
 2 X = (X - np.mean(X)) / np.std(X)
 3
 4 covX = np.dot(X.T, X)  # Covariance matrix (M x M)
 5
 6 d, Vcols = np.linalg.eig(covX)  # Eigenvalues and eigenvectors
 7
 8 plt.figure(figsize=(16, 4))
 9 plt.title("Significance of Eigenvalues of the Covariance Matrix covX")
10 plt.semilogx(abs(d), '-ok')
11 plt.ylabel('|d|')
12 plt.xlabel('Index of Eigenvalues d (i.e., Index of Vcols)')
13 plt.grid()
14 plt.show()
15
16
```



## Compression using PCA

(And its use for the localization of anomalous events by filtering uncorrelated events using PCA)

i.e., residuals = original - image without uncorrelated events => occurrence of anomalies

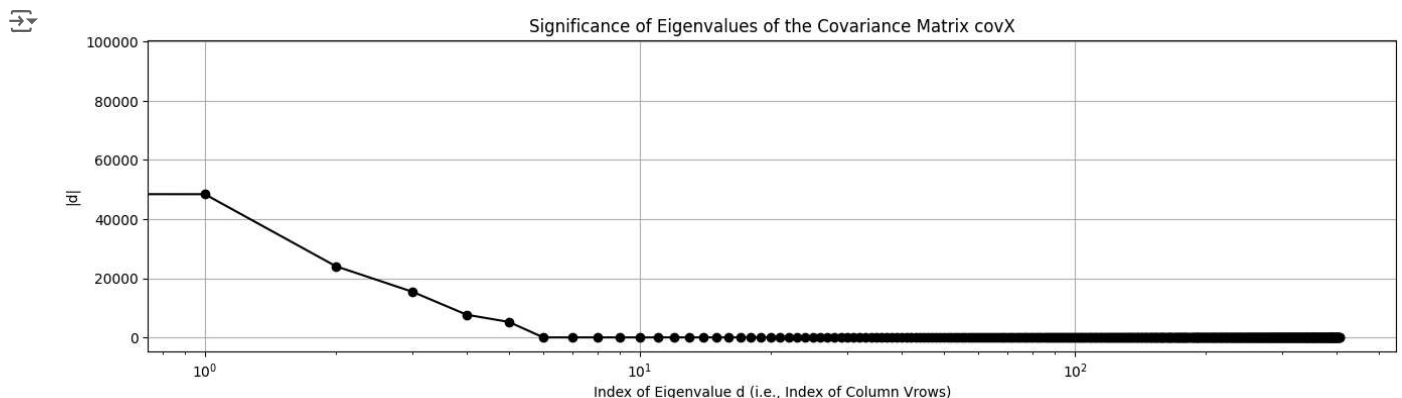## Compression through Columns and Subsequently through Rows

Through Columns: I choose $m << M$ most significant eigenvectors, i.e., $n$ columns $\mathbf{V}_{cols}$ for which the eigenvalues $d$ are the largest (see the previous graph) (the compressed image will have $N$ rows and $m$ columns).

```
1 print("X ... N x M =",X.shape)
2 m=6
3 print("choosing m=",m)
4
5 Xcols=np.dot(X,Vcols[:,:m])  # compressed na m sloupcu
6
7 print("Xcols ... N x m =",Xcols.shape)
8
```

```
X ... N x M = (407, 594)
choosing m= 6
Xcols ... N x m = (407, 6)
```

Along the rows: I select $n << N$ most significant eigenvectors, i.e., $n$ columns $\mathbf{V}_{rows}$ for which the eigenvalues $d$ are the largest (see the following graph).
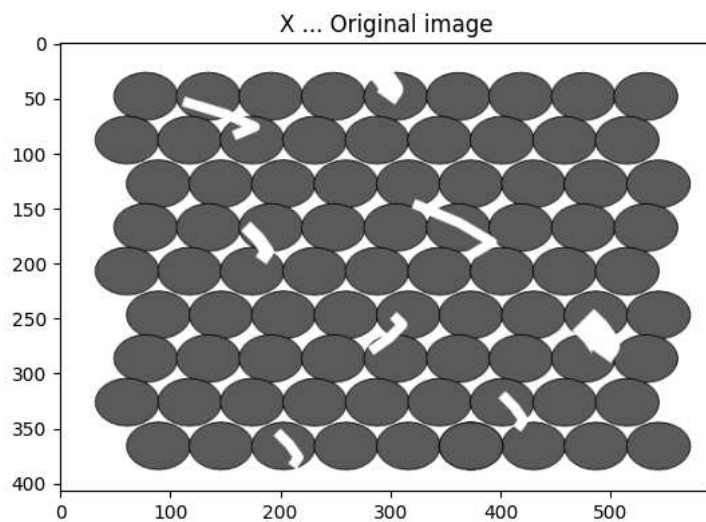
```
1 # Compression along the rows
2 Xrows = Xcols.T
3 covX = np.dot(Xrows.T, Xrows)  # Covariance matrix (N x N)
4 d, Vrows = np.linalg.eig(covX)
5
6 plt.figure(figsize=(16, 4))
7 plt.title("Significance of Eigenvalues of the Covariance Matrix covX")
8 plt.semilogx(abs(d), '-ok')
9 plt.ylabel('|d|')
10 plt.xlabel('Index of Eigenvalue d (i.e., Index of Column Vrows)')
11 plt.grid()
12 plt.show()
13
```
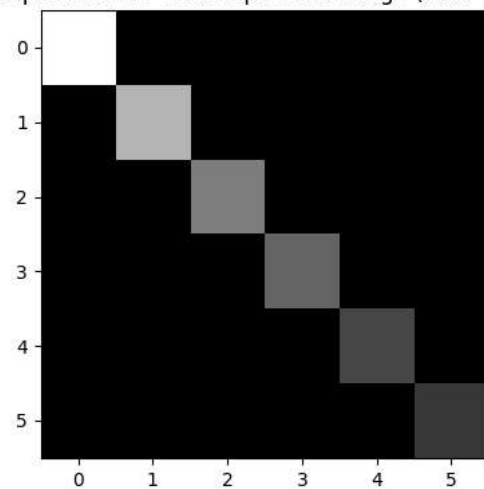


```
1 print("Xrows ... n x M =", Xrows.shape)
2 n = 6
3 print("Choice of n =", n)
4 Xcompressed = np.dot(Xrows, Vrows[:, :n])  # Compressed image with n rows and m columns
5 print("Xcompressed ... n x m =", Xcompressed.shape)
6
7 ### Decompression - retaining only correlated events (the degree of correlation depends on the choice of m and n)
8
9 Xpca_rows = np.dot(Xcompressed, Vrows[:, :n].T)
10 Xpca_cols = Xpca_rows.T
11 Xpca = np.dot(Xpca_cols, Vcols[:, :m].T)
12
13 Xpca = np.real(Xpca)
14
15 plt.figure(figsize=(16, 20))
16 plt.subplot(411)
17 plt.imshow(X, cmap='gray')
18 plt.title("X ... Original image")
19 plt.subplot(412)
20 plt.imshow(abs(Xcompressed), cmap='gray')
21 plt.title("Xcompressed ... PCA compressed image (NxM -> mxn)")
22 plt.subplot(413)
23 plt.imshow(Xpca, cmap='gray')
24 plt.title("Xpca ... PCA compressed and decompressed image (NxM -> mxn -> NxM)")
25 plt.subplot(414)
26 dX = (Xpca - X) ** 4  # Thresholding (enhancing the difference between white and black), using powers of 2, 4, 6, 8
27 plt.imshow(-dX, cmap='gray')
```

```
28 plt.title("Detected fault locations (biased) as biased residuals dX=(Xpca-X)^4")
29 plt.show()
30
```
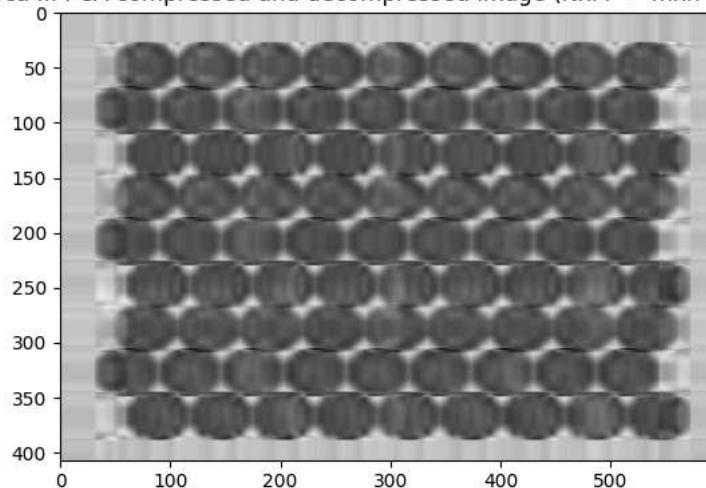
```
Xrows ... n x M = (6, 407)
Choice of n = 6
Xcompressed ... n x m = (6, 6)
```

**X ... Original image**



**Xcompressed ... PCA compressed image (NxM -> mxn)**



**Xpca ... PCA compressed and decompressed image (NxM -> mxn -> NxM)**



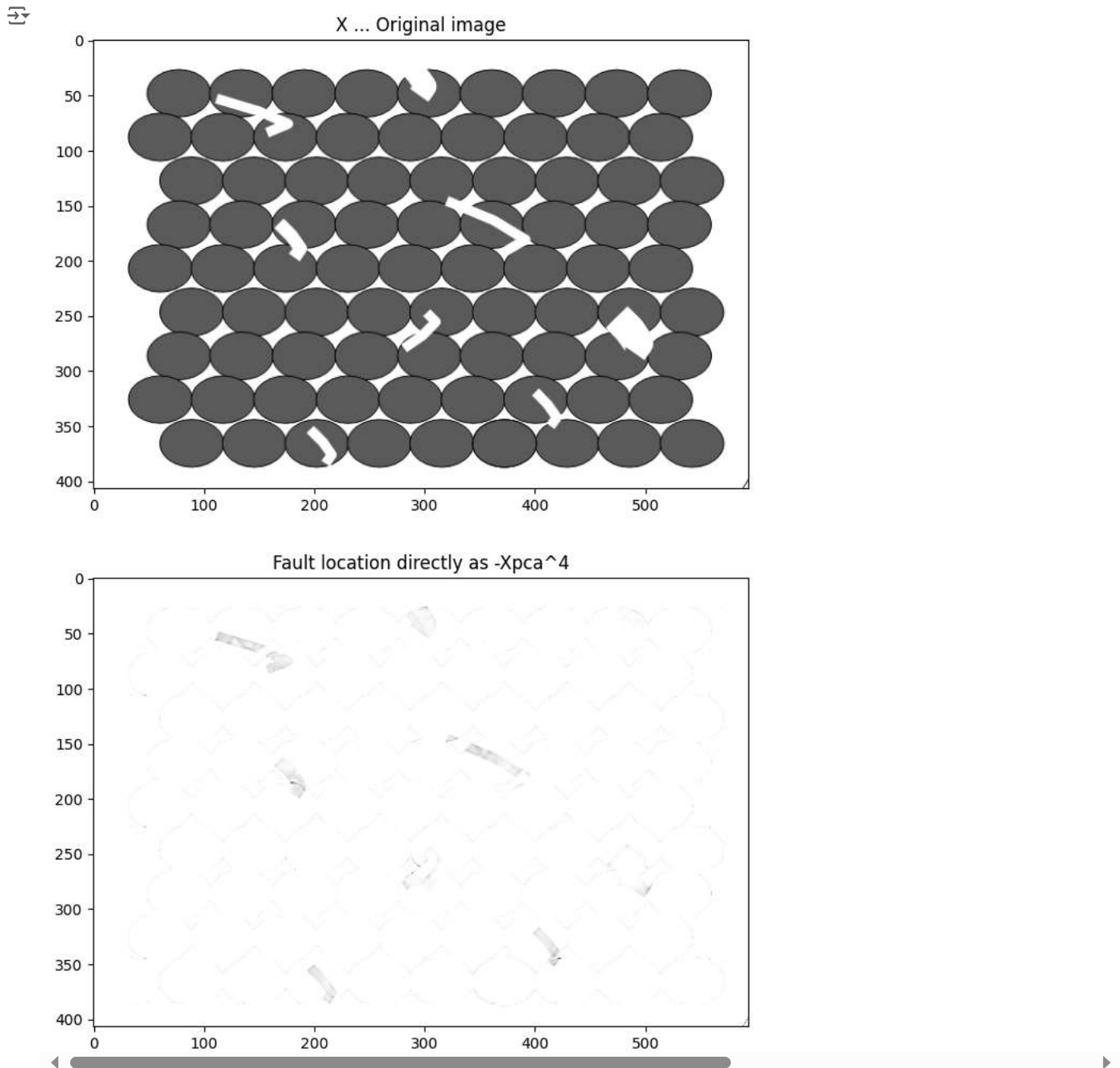**Detected fault locations (biased) as biased residuals dX=(Xpca-X)^4**

```
1 Xpca.shape
```

    (407, 594)

## Localization of uncorrelated events by filtering correlated events using PCA (Xpca => occurrence of anomalies in the image)

i.e. filtering correlated events => occurrence of anomalies

```
1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib.image as mpimg
5 import PIL as pil
6 from scipy import signal
7
8 # Path to the image
9 img1="./img/regular_patttern_with_perturb.png"
10 im1 = pil.Image.open(img1)
11 im1 = np.mean(im1, 2)  # Image as a single matrix
12
13 # Copy and normalize the image
14 X = im1.copy()
15 X = (X - np.mean(X)) / np.std(X)
16
17 # Calculate the covariance matrix (M x M)
18 covX = np.dot(X.T, X)
19
20 # Eigenvalues and eigenvectors
21 d, Vcols = np.linalg.eig(covX)
22
23 # Define the value 'm'
24 m = 6
25
26 # Project the data onto the first 'm' columns of Vcols
27 Xcols = np.dot(X, Vcols[:, m:])
28
29 # Project the data along the columns to get 'Xrows'
30 Xrows = Xcols.T
31
32 # Calculate the covariance matrix (N x N) for 'Xrows'
33 covX = np.dot(Xrows.T, Xrows)
34
35 # Eigenvalues and eigenvectors for 'Xrows'
36 d, Vrows = np.linalg.eig(covX)
37
38 # Define the value 'n'
39 n = 6
40
41 # Compress the image to 'n' rows and 'm' columns
42 Xcompressed = np.dot(Xrows, Vrows[:, n:])
43
44 # Perform the decompression
45 Xpca_rows = np.dot(Xcompressed, Vrows[:, n:].T)
46 Xpca_cols = Xpca_rows.T
47 Xpca = np.dot(Xpca_cols, Vcols[:, m:].T)
48 Xpca = np.real(Xpca)
49
50 # Create a figure for visualization
51 plt.figure(figsize=(10, 12))
52 plt.subplot(211)
53 plt.title("X ... Original image")
54 plt.imshow(X, cmap='gray')
55 plt.subplot(212)
56 plt.imshow(-Xpca**4, cmap='gray')
57 plt.title("Fault location directly as -Xpca^4")
58 plt.show()
59
```

X ... Original image


Fault location directly as -Xpca^4

## Tasks:

**T3.1** Try PCA on your own image for the detection of uncorrelated anomalies in the image (or noise). **[0.3 Points]**

**T3.2** Explore the influence of the choice of the number of eigenvectors on the quality of detection. **[0.3 Points]**

**T3.3** Create concise documentation for the solution, code, and the achieved results (3-5 pages of A4). Properly cite the used sources! **[0.6 Points]**

```
1 # 3.1) Try PCA on your own image for the detection of uncorrelated anomalies in the image (or noise)
2 # 3.2) Explore the influence of the choice of the number of eigenvectors on the quality of detection.
3
4 %matplotlib inline
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import matplotlib.image as mpimg
8 import PIL as pil
9 img1="/content/image1.webp"
10
11
12 im1 = pil.Image.open(img1)
13
14 from scipy import signal
15 im1=np.mean(im1,2)  # gray img, 2-D matrix,...X
16
17 X = im1.copy()
```
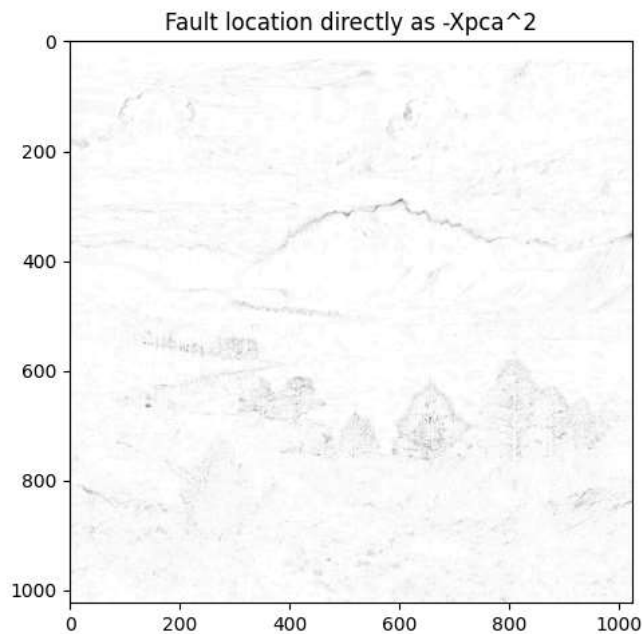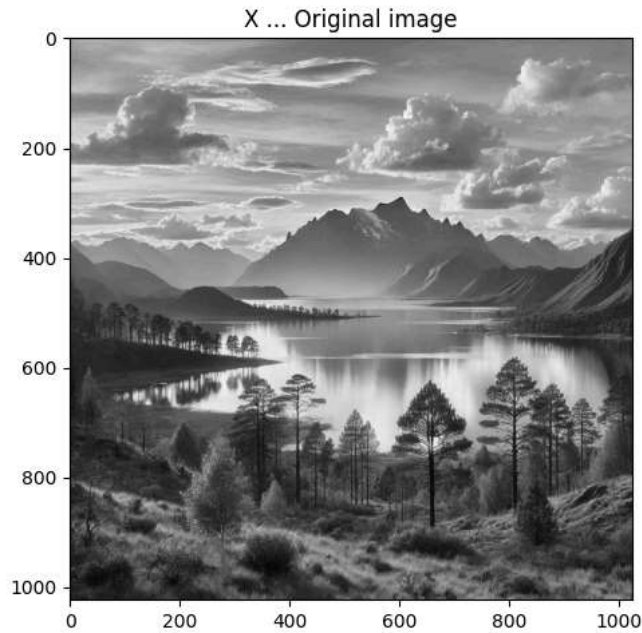
```python
18 X = (X - np.mean(X)) / np.std(X)
19
20 # Calculate the covariance matrix (M x M)
21 covX = np.dot(X.T, X)
22
23 # Eigenvalues and eigenvectors
24 d, Vcols = np.linalg.eig(covX)
25
26 print("X ... N x M =",X.shape)
27
28 m = 8
29
30 # Project the data onto the first 'm' columns of Vcols
31 Xcols = np.dot(X, Vcols[:, m:])
32
33 # Project the data along the columns to get 'Xrows'
34 Xrows = Xcols.T
35
36 # Calculate the covariance matrix (N x N) for 'Xrows'
37 covX = np.dot(Xrows.T, Xrows)
38
39 # Eigenvalues and eigenvectors for 'Xrows'
40 d, Vrows = np.linalg.eig(covX)
41
42 # Define the value 'n'
43 n = 8
44
45 # Compress the image to 'n' rows and 'm' columns
46 Xcompressed = np.dot(Xrows, Vrows[:, n:])
47
48 # Perform the decompression
49 Xpca_rows = np.dot(Xcompressed, Vrows[:, n:].T)
50 Xpca_cols = Xpca_rows.T
51 Xpca = np.dot(Xpca_cols, Vcols[:, m:].T)
52 Xpca = np.real(Xpca)
53
54 # Create a figure for visualization
55 plt.figure(figsize=(10, 12))
56 plt.subplot(211)
57 plt.title("X ... Original image")
58 plt.imshow(X, cmap='gray')
59 plt.subplot(212)
60 plt.imshow(-Xpca**2, cmap='gray')
61 plt.title("Fault location directly as -Xpca^2")
62 plt.show()
```

X ... Original image
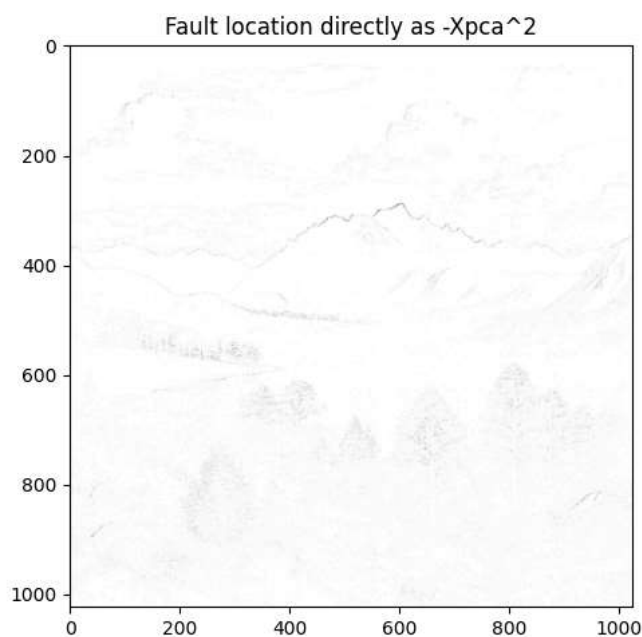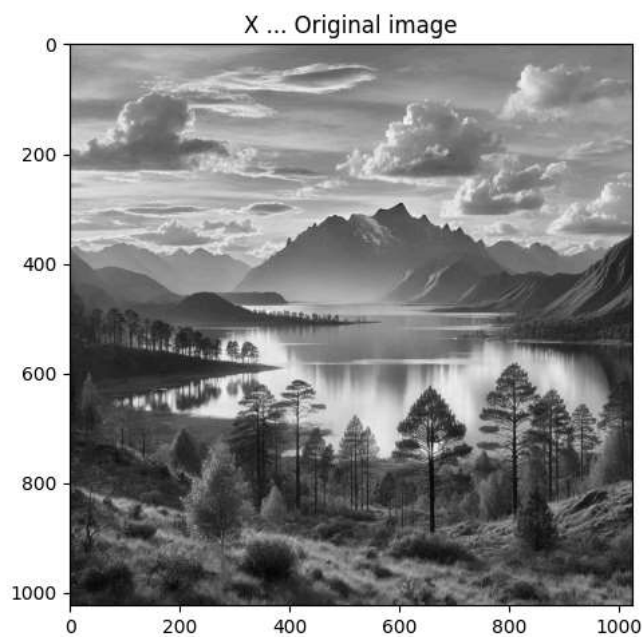


Fault location directly as -Xpca^2



```
1
2 X = im1.copy()
3 X = (X - np.mean(X)) / np.std(X)
4
5 # Calculate the covariance matrix (M x M)
6 covX = np.dot(X.T, X)
7
8 # Eigenvalues and eigenvectors
9 d, Vcols = np.linalg.eig(covX)
10
11 print("X ... N x M =",X.shape)
12
13 m = 50
14
15 # Project the data onto the first 'm' columns of Vcols
16 Xcols = np.dot(X, Vcols[:, m:])
17
18 # Project the data along the columns to get 'Xrows'
19 Xrows = Xcols.T
20
21 # Calculate the covariance matrix (N x N) for 'Xrows'
22 covX = np.dot(Xrows.T, Xrows)
23
24 # Eigenvalues and eigenvectors for 'Xrows'
25 d, Vrows = np.linalg.eig(covX)
26
27 # Define the value 'n'
28 n = 50
```

```python
28  n = 50
29
30  # Compress the image to 'n' rows and 'm' columns
31  Xcompressed = np.dot(Xrows, Vrows[:, n:])
32
33  # Perform the decompression
34  Xpca_rows = np.dot(Xcompressed, Vrows[:, n:].T)
35  Xpca_cols = Xpca_rows.T
36  Xpca = np.dot(Xpca_cols, Vcols[:, m:].T)
37  Xpca = np.real(Xpca)
38
39  # Create a figure for visualization
40  plt.figure(figsize=(10, 12))
41  plt.subplot(211)
42  plt.title("X ... Original image")
43  plt.imshow(X, cmap='gray')
44  plt.subplot(212)
45  plt.imshow(-Xpca**2, cmap='gray')
46  plt.title("Fault location directly as -Xpca^2")
47  plt.show()
```

X ... N x M = (1024, 1024)



X ... Original image



Fault location directly as -Xpca^2

```python
1  X = im1.copy()
2  X = (X - np.mean(X)) / np.std(X)
3
4  # Calculate the covariance matrix (M x M)
5  covX = np.dot(X.T, X)
6
7  # Eigenvalues and eigenvectors
```

```
 8 d, Vcols = np.linalg.eig(covX)
 9
10 print("X ... N x M =",X.shape)
11
12 m = 1
13
```