

NATURAL LANGUAGE PROCESSING

TOPIC - PYTHON PROJECT ON SVD

Name-Simran Kumari

Roll NO- BTECH/10071/21

BRANCH – CSE

SEC – A

SUBMITTED TO: INDRAJIT MUKHARJEE

INTRODUCTION

This is my Natural Language Processing project on the topic Singular Value Decomposition. SVD of a matrix is a factorization of that matrix into three matrices. It has some interesting algebraic properties and conveys important geometrical and theoretical insights about linear transformations. It also has some important applications in data science. In this article, I will try to explain the mathematical intuition behind SVD and its geometrical meaning.

Mathematics behind SVD:

The SVD of $m \times n$ matrix A is given by the formula $A = U\Sigma V^T$

where:

U : $m \times m$ matrix of the orthonormal eigenvectors of AA^T .

V : transpose of a $n \times n$ matrix containing the orthonormal eigenvectors of $A^T A$.

Σ : diagonal matrix with r elements equal to the root of the positive eigenvalues of AA^T or $A^T A$ (both matrices have the same positive eigenvalues anyway).

Examples

Find the SVD for the matrix $A = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$

To calculate the SVD, First, we need to compute the singular values by finding eigenvalues of AA^T .

$A \cdot A^T =$

$$\begin{aligned} & \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix} \\ \cdot & \\ & \begin{bmatrix} 3 & 2 \\ 2 & 3 \\ 2 & -2 \end{bmatrix} \\ = & \\ & \begin{bmatrix} 17 & 8 \\ 8 & 17 \end{bmatrix} \end{aligned}$$

The characteristic equation for the above matrix is:

$$W - \lambda I = 0 \quad A A^T - \lambda I = 0$$

$$\lambda^2 - 34\lambda + 225 = 0$$

$$= (\lambda - 25)(\lambda - 9)$$

so our singular values are: $\sigma_1 = 5$, ; $\sigma_2 = 3$

Now we find the right singular vectors i.e orthonormal set of eigenvectors of ATA . The eigenvalues of ATA are 25, 9, and 0, and since ATA is symmetric we know that the eigenvectors will be orthogonal.

For $\lambda = 25$,

$$A^T A - 25 I =$$

$$\begin{bmatrix} -12 & 12 & 2 \\ 12 & -12 & -2 \\ 2 & -2 & -17 \end{bmatrix}$$

which can be row-reduced to :

$$\begin{bmatrix} 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

A unit vector in the direction of it is:

$$v_1 =$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$$

Similarly, for $\lambda = 9$, the eigenvector is:

$$v_2 =$$

$$\begin{bmatrix} \frac{1}{\sqrt{18}} \\ \frac{-1}{\sqrt{18}} \\ \frac{4}{\sqrt{18}} \end{bmatrix}$$

For the 3rd eigenvector, we could use the property that it is perpendicular to v_1 and v_2 such that:

$$v_1^T v_3 = 0 \quad v_2^T v_3 = 0$$

Solving the above equation to generate the third eigenvector

$$v_3 =$$

$$\begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

=

$$\begin{bmatrix} a \\ -a \\ -a/2 \end{bmatrix}$$

=

$$\begin{bmatrix} \frac{2}{3} \\ \frac{-2}{3} \\ \frac{-1}{3} \end{bmatrix}$$

Now, we calculate U using the formula $u_i = \frac{1}{\|\sigma_i\|} A v_i$ and this gives $U =$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{bmatrix}$$

. Hence, our final SVD equation becomes:

$$A =$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ 5 & 0 & 0 \\ 0 & 3 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{bmatrix}$$

Applications of Singular Value Decomposition (SVD):

Applications of Singular Value Decomposition (SVD):

Latent Semantic Analysis (LSA):

Application: SVD is widely used in LSA to discover latent semantic structures in large text corpora. It helps in capturing relationships between terms and documents, enabling improved information retrieval and document clustering.

Collaborative Filtering:

Application: SVD is applied in recommendation systems for collaborative filtering. It helps identify latent factors that contribute to user preferences, making personalized recommendations based on user-item interaction matrices.

Dimensionality Reduction:

Application: SVD is employed for dimensionality reduction in various domains. In NLP, it can be applied to reduce the dimensionality of term-document matrices, making it computationally more efficient while preserving essential semantic information.

Image Compression:

Application: SVD is used in image compression by decomposing image matrices. It allows the representation of images with fewer singular values, leading to reduced storage requirements without significant loss of visual quality.

Noise Reduction in Signals:

Application: In signal processing, SVD is utilized to reduce noise in signals. By decomposing a signal matrix, it is possible to retain the most significant components while filtering out noise.

Face Recognition:

Application: SVD is applied in facial recognition systems. By decomposing the face image matrix, it helps extract essential features and reduce the dimensionality of the data, improving the efficiency of face recognition algorithms.

Principal Component Analysis (PCA):

Application: SVD is a fundamental component of PCA. It helps identify the principal components of a dataset, aiding in feature extraction and data visualization.

Topic Modeling in NLP:

Application: SVD is used for topic modeling in NLP, helping to identify latent topics in a collection of documents. It aids in understanding

Challenges of SVD

Challenges of Singular Value Decomposition (SVD):

Computational Complexity:

Challenge: SVD can be computationally expensive, especially for large matrices. The cubic time complexity makes it challenging to apply to massive datasets in real-time applications.

Storage Requirements:

Challenge: The storage requirements for storing the three resulting matrices (U , Σ , and V^T) can be substantial, especially for large matrices. This poses challenges in terms of memory usage and computational efficiency.

Interpretability:

Challenge: The interpretation of the latent factors obtained through SVD may not always be straightforward. Understanding the meaning of singular vectors and values in the context of the original data can be challenging.

Handling Missing Values:

Challenge: SVD is sensitive to missing values in the input matrix. Dealing with missing data can be complex and may require imputation techniques, which can introduce additional uncertainties.

Non-Negativity:

Challenge: SVD does not naturally handle non-negativity constraints. In some applications, such as topic modeling or image processing, non-negativity is crucial, and modifications or additional constraints may be necessary.

Scalability:

Challenge: SVD may not scale well to very large datasets or high-dimensional matrices. This limitation can be a hindrance in applications where scalability is essential.

Overfitting:

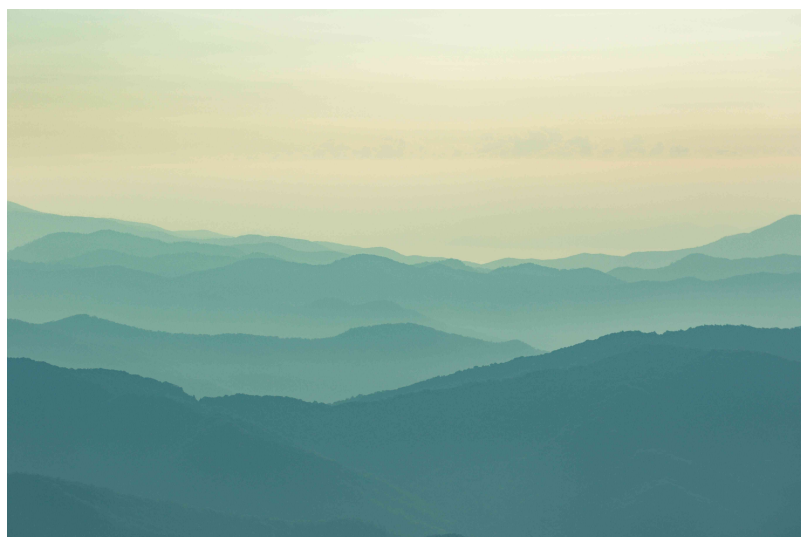
Challenge: In some cases, SVD may be prone to overfitting, capturing noise in the data rather than the underlying patterns. Regularization techniques may be required to mitigate this issue.

Despite these challenges, SVD remains a powerful tool in various applications, and researchers continue to develop techniques and algorithms to address its limitations and make it more applicable to diverse datasets and scenarios.

Image Used

```
In [45]: 1 from IPython.display import Image
          2 Image(filename="svd.jpg", width=400, height=400)
          3
```

Out[45]:



In [50]:

```

1 import requests
2 from PIL import Image
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 # Assign and open image
7 url = 'https://media.geeksforgeeks.org/wp-content/cdn-uploads/202104011734
8 response = requests.get(url, stream=True)
9
10 with open('image.png', 'wb') as f:
11     f.write(response.content)
12
13 # Open image using Pillow (PIL)
14 img = Image.open('image.png')
15
16 # Convert the image to grayscale
17 gray_image = img.convert('L')
18
19 # Convert the grayscale image to a NumPy array
20 gray_array = np.array(gray_image)
21
22 # Calculating the SVD
23 u, s, v = np.linalg.svd(gray_arr, full_matrices=False)
24
25 # Inspect shapes of the matrices
26 print(f'u.shape: {u.shape}, s.shape: {s.shape}, v.shape: {v.shape}')
27

```

u.shape: (3648, 3648), s.shape: (3648,), v.shape: (3648, 5472)

ERROR: Exception:

Traceback (most recent call last):

```

File "C:\Users\simra\OneDrive\Documents\Python Scripts\lib\site-packages\pi
p\_vendor\urllib3\response.py", line 437, in _error_catcher
    yield
File "C:\Users\simra\OneDrive\Documents\Python Scripts\lib\site-packages\pi
p\_vendor\urllib3\response.py", line 560, in read
    data = self._fp_read(amt) if not fp_closed else b""
File "C:\Users\simra\OneDrive\Documents\Python Scripts\lib\site-packages\pi
p\_vendor\urllib3\response.py", line 526, in _fp_read
    return self._fp.read(amt) if amt is not None else self._fp.read()
File "C:\Users\simra\OneDrive\Documents\Python Scripts\lib\site-packages\pi
p\_vendor\cachecontrol\filewrapper.py", line 90, in read
    data = self.__fp.read(amt)
File "C:\Users\simra\OneDrive\Documents\Python Scripts\lib\http\client.py",
line 465, in read
    s = self.fp.read(amt)
File "C:\Users\simra\OneDrive\Documents\Python Scripts\lib\socket.py" line

```

Explanation:

The above output shape indicates that there are 3648 linearly independent eigenvectors in this image.

Now let us look at the variance of the image used over a singular vector graphically:

```

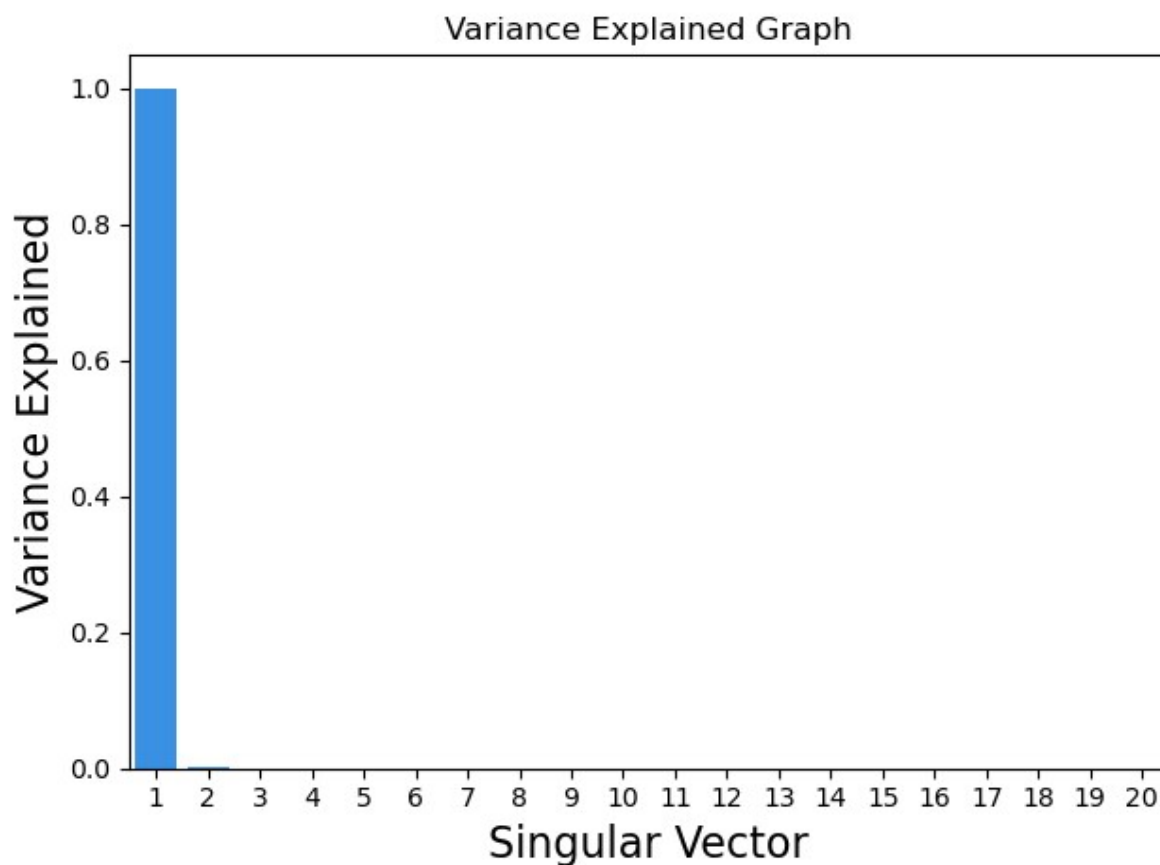
In [51]: 1 # import module
          2 import seaborn as sns
          3
          4 var_explained = np.round(s**2/np.sum(s**2), decimals=6)
          5
          6 # Variance explained top Singular vectors
          7 print(f'variance Explained by Top 20 singular values:\n{var_explained[0:20]}')
          8
          9 sns.barplot(x=list(range(1, 21)),
         10          y=var_explained[0:20], color="dodgerblue")
         11
         12 plt.title('Variance Explained Graph')
         13 plt.xlabel('Singular Vector', fontsize=16)
         14 plt.ylabel('Variance Explained', fontsize=16)
         15 plt.tight_layout()
         16 plt.show()
         17

```

```

variance Explained by Top 20 singular values:
[9.9795e-01 7.4000e-04 3.8900e-04 3.0500e-04 1.0600e-04 7.4000e-05
 5.7000e-05 3.9000e-05 3.1000e-05 2.6000e-05 2.2000e-05 1.9000e-05
 1.5000e-05 1.4000e-05 1.1000e-05 9.0000e-06 9.0000e-06 7.0000e-06
 7.0000e-06 6.0000e-06]

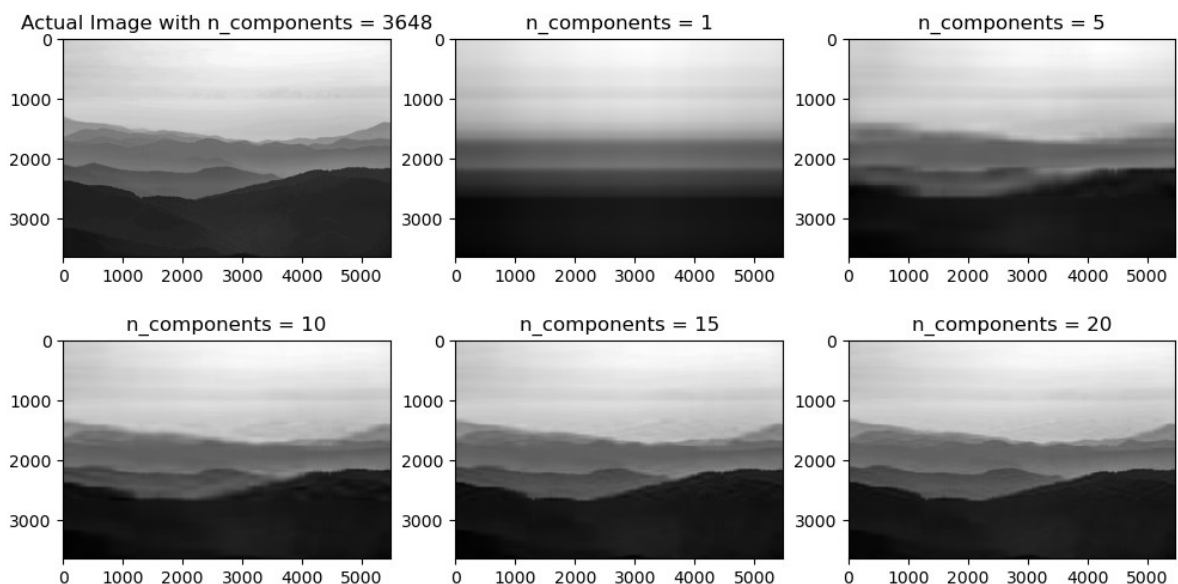
```



Explanation: The Variance Explained Graph above clearly shows that about 99.77 % of information is explained by the first eigenvector and its corresponding eigenvalues themselves. Therefore, it very much advisable to reconstruct the image with just the top few eigenvectors themselves.

In the below program based on the above discussion, we reconstruct the image using SVD:

```
In [53]: 1 # plot images with different number of components
2 comps = [3648, 1, 5, 10, 15, 20]
3 plt.figure(figsize=(12, 6))
4
5 for i in range(len(comps)):
6     low_rank = u[:, :comps[i]] @ np.diag(s[:comps[i]]) @ v[:, comps[i], :]
7     if(i == 0):
8         plt.subplot(2, 3, i+1),
9         plt.imshow(low_rank, cmap='gray'),
10        plt.title(f'Actual Image with n_components = {comps[i]}')
11
12    else:
13        plt.subplot(2, 3, i+1),
14        plt.imshow(low_rank, cmap='gray'),
15        plt.title(f'n_components = {comps[i]}')
16
```



Explanation:

Though the 1st eigenvector contains 99.77% of information reconstructing an image solely from it does not give a clear picture.

Using the top 15 vectors for the image reconstruction gives a good enough approximation. Also out of 3648 vectors which is a massive decrease in computation and also it compresses the image

CONCLUSION

The extended Singular Value Decomposition (SVD) project, incorporating variance explained visualization and image reconstruction with varying components, has provided valuable insights

into the application and impact of SVD on image data. Here's a summary of key findings and takeaways:

1. Variance Explained:

The project began by exploring the variance explained by the top singular values. The variance explained graph visually represented the contribution of each singular value to the overall variance in the data. This analysis is essential for understanding the significance of different components in capturing the essence of the image.

2. Image Reconstruction:

A crucial aspect of the project involved reconstructing images using different numbers of singular components. By visualizing the reconstructed images, the project illustrated the trade-off between dimensionality reduction and image fidelity. Lower numbers of components resulted in compressed representations, showcasing the ability of SVD to capture essential features with reduced information.

3. Interpretation of Results:

The analysis of variance explained and reconstructed images provided a comprehensive understanding of how SVD captures and represents information in image data. The project demonstrated the importance of selecting an appropriate number of components based on the desired level of compression and the preservation of relevant features.

4. Challenges and Considerations:

The project acknowledged challenges such as the computational complexity of SVD, especially for large images. Additionally, interpreting the visual results requires careful consideration of the balance between dimensionality reduction and information loss.

5. Next Steps:

Further exploration could involve comparing the performance of SVD with other dimensionality reduction techniques, such as Principal Component Analysis (PCA) or autoencoders, to assess their effectiveness in image processing tasks. The project could be extended to include quantitative metrics for assessing the quality of reconstructed images, providing a more rigorous evaluation of the compression techniques.

6. Practical Applications:

The project's insights into image reconstruction and variance explained are applicable in various fields, including image compression, feature extraction, and pattern recognition. Understanding the capabilities and limitations of SVD in image processing is crucial for making informed decisions in real-world applications.

In conclusion, this extended SVD project has provided me hands-on exploration of the technique's applications in image analysis. The findings contribute to a broader understanding of the role of SVD in capturing and representing essential information in image data, with potential implications for diverse fields within computer vision and data analysis.

In []:

1

