



AGH

AGH UNIVERSITY OF KRAKOW

Department of Computer Science

Project nr 2
Visualization of structured and unstructured data

"PCA kernelPCA"

Author:
Degree programme:
Supervisor:

Jacek Tyszkiewicz
Computer Science

Krakow, 2025

Table of Contents

1. PCA kernelPCA 3

1.1. Perform PCA for Breast Cancer Dataset (3 points)..... 3

1.2. KernelPCA (3 points) 8

1.3. Classic PCA Analysis on MNIST (4 points) 13

1.3.1. Important:..... 13

1. PCA kernelPCA

1.1. Perform PCA for Breast Cancer Dataset (3 points)

- You can find this dataset in the scikit-learn library, import it and convert to pandas dataframe. Original labels are '0' and '1'; for better readability change these names to: 'benign' and 'malignant'.
- Visualize correlations between pairs of features (due to the greater number of features use pandas corr() function instead of pairplot and seaborn heatmap()).
- Examine (zbada) explained variance, draw a plot showing relation between total explained variance and number of principal components used.
- Use recursive feature elimination (available in scikit-learn module) or another feature ranking algorithm to split 30 features into 15 “more important” and 15 “less important” features.
Then repeat the last step from the full data set – draw a plot showing relation between total explained variance and number of principal components used for all 3 cases.

Explain the result briefly.

- You can find this dataset in the scikit-learn library, import it and convert to pandas dataframe. Original labels are '0' and '1'; for better readability change these names to: 'benign' and 'malignant'.

```
1 from IPython.display import display
2 pd.set_option('display.max_columns', None)
3 pd.set_option('display.width', None)
4 pd.set_option('display.max_colwidth', None)

1 from sklearn.datasets import load_breast_cancer
2 import pandas as pd
3
4 data = load_breast_cancer()
5 df = pd.DataFrame(data.data, columns=data.feature_names)
6
7
8 df['target'] = data.target
9 df['target'] = df['target'].map({0: 'malignant', 1: 'benign'})
10
11 display(df.head())

1 corr_matrix = df.drop(columns='target').corr()
2
```

```

3 plt.figure(figsize=(16, 12))
4 sns.heatmap(corr_matrix, annot=False, cmap='coolwarm', center=0)
5 plt.title('Feature_Correlation_Matrix_(Breast_Cancer_Dataset)')
6 plt.show()

```

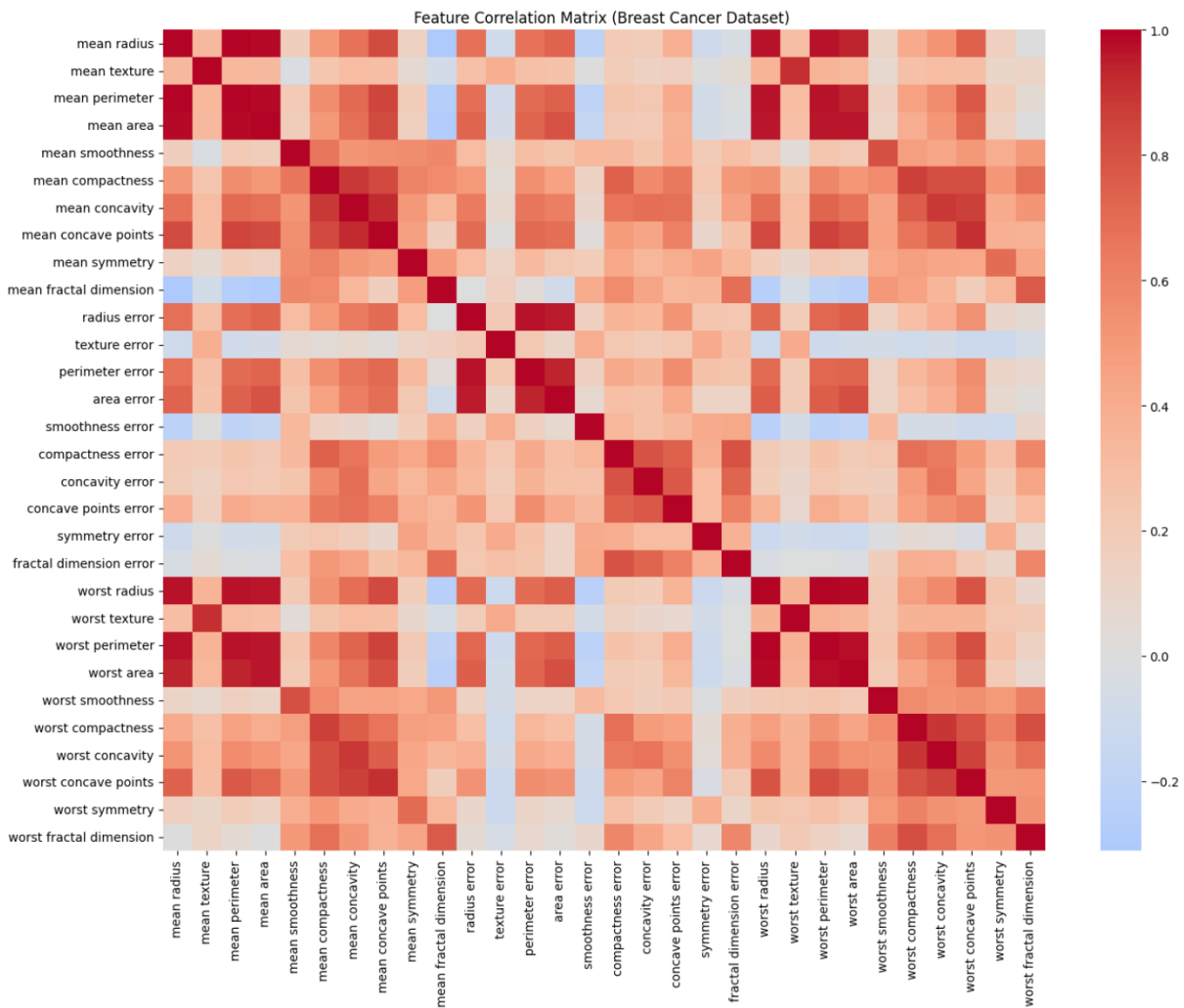


Fig. 1.1. Feature Heatmap

```

1 X = data.data
2
3 X_scaled = StandardScaler().fit_transform(X)
4
5 pca = PCA()
6 X_pca = pca.fit_transform(X_scaled)
7
8 explained_variance_ratio = pca.explained_variance_ratio_
9 cumulative_variance = np.cumsum(explained_variance_ratio)
10
11 plt.figure(figsize=(10, 6))
12 plt.plot(range(1, len(cumulative_variance)+1), cumulative_variance, marker='o',
13         linestyle='--')
14 plt.title('Explained variance for number of principal components used')
15 plt.xlabel('number of principal components')
16 plt.ylabel('cumulative explained variance')
17 plt.grid(True)
18 plt.show()

```

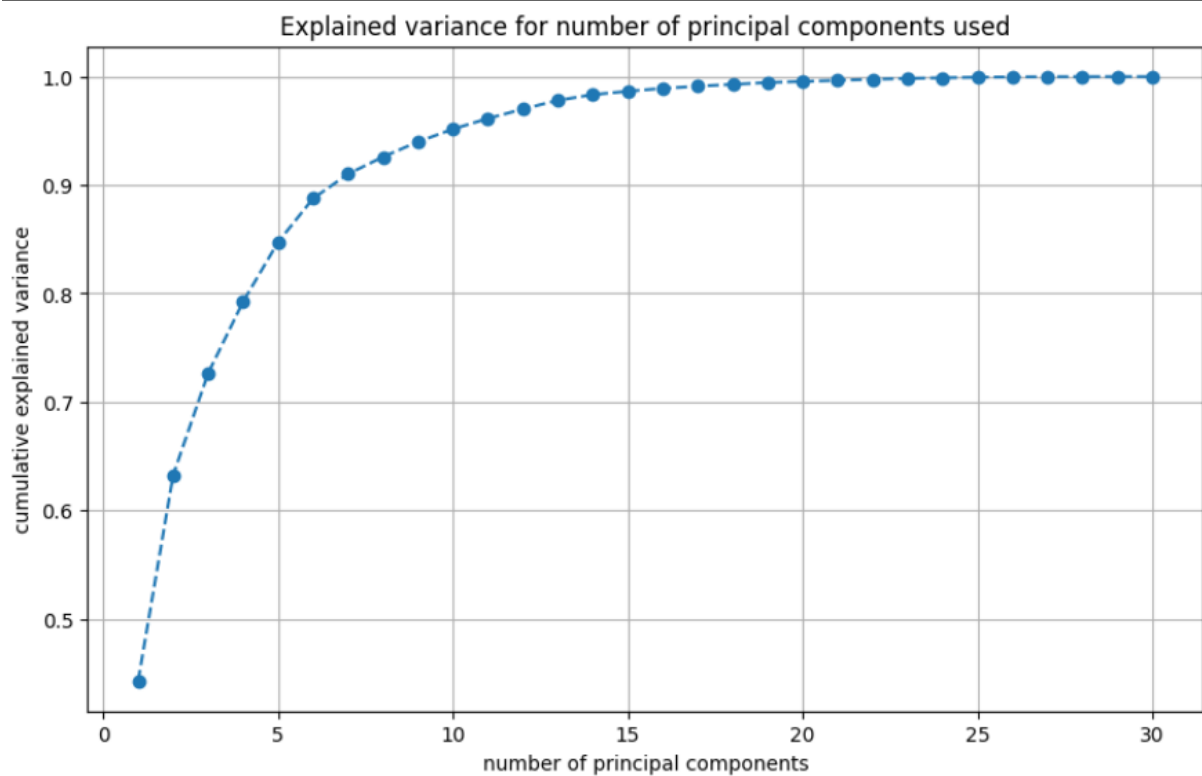
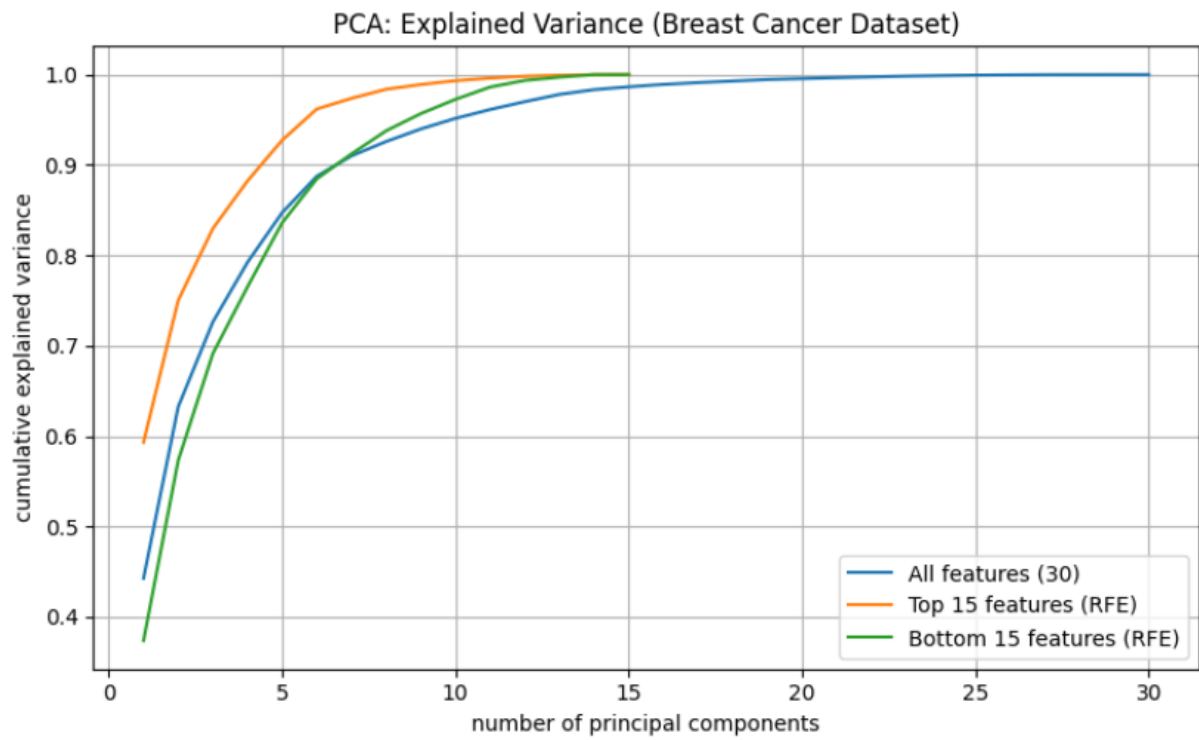


Fig. 1.2. Feature Heatmap

```
1 X = data.data
2 y = data.target
3
4 X = StandardScaler().fit_transform(X)
5
6 estimator = LogisticRegression(max_iter=1000)
7 selector = RFE(estimator, n_features_to_select=15)
8 selector.fit(X, y)
9
10 X_all = X
11 X_top15 = X[:, selector.support_]
12 X_bottom15 = X[:, ~selector.support_]
13
14
15 def explained_variance(X_data):
16     pca = PCA().fit(X_data)
17     return np.cumsum(pca.explained_variance_ratio_)
18
19 var_all = explained_variance(X_all)
20 var_top15 = explained_variance(X_top15)
21 var_bottom15 = explained_variance(X_bottom15)
22
23
24 plt.figure(figsize=(8, 5))
25 plt.plot(range(1, len(var_all) + 1), var_all, label="All_features_(30)")
26 plt.plot(range(1, len(var_top15) + 1), var_top15, label="Top_15_features_(RFE)"
27 )
28 plt.plot(range(1, len(var_bottom15) + 1), var_bottom15, label="Bottom_15_
29 features_(RFE)")
30
31 plt.xlabel("number_of_principal_components")
32 plt.ylabel("cumulative_explained_variance")
33 plt.title("PCA: Explained Variance (Breast Cancer Dataset)")
34 plt.legend()
35 plt.grid(True)
36 plt.tight_layout()
37 plt.show()
```

**Fig. 1.3.** Feature Heatmap

1.2. KernelPCA (3 points)

- Visualize in 2D datasets used in this lab, experiment with the parameters of the KernelPCA method. Change kernel and gamma params.

Documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA.html>

```

1 data = load_breast_cancer()
2 X = data.data
3 y = data.target
4 X_scaled = StandardScaler().fit_transform(X)
5
6 gammas = [0.001, 0.01, 0.1, 1, 10]
7
8 fig, axes = plt.subplots(1, len(gammas), figsize=(18, 4))
9 fig.suptitle("KernelPCA_(RBF)_-effect_of_the_gamma_parameter", fontsize=16)
10
11 # Pętla po gamma
12 for i, gamma in enumerate(gammas):
13     kpca = KernelPCA(n_components=2, kernel='rbf', gamma=gamma)
14     X_kpca = kpca.fit_transform(X_scaled)
15
16     ax = axes[i]
17     ax.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y, cmap='coolwarm', s=10)
18     ax.set_title(f"gamma_{gamma}")
19     ax.set_xticks([])
20     ax.set_yticks([])
21
22 plt.tight_layout(rect=[0, 0, 1, 0.93])
23 plt.show()

```

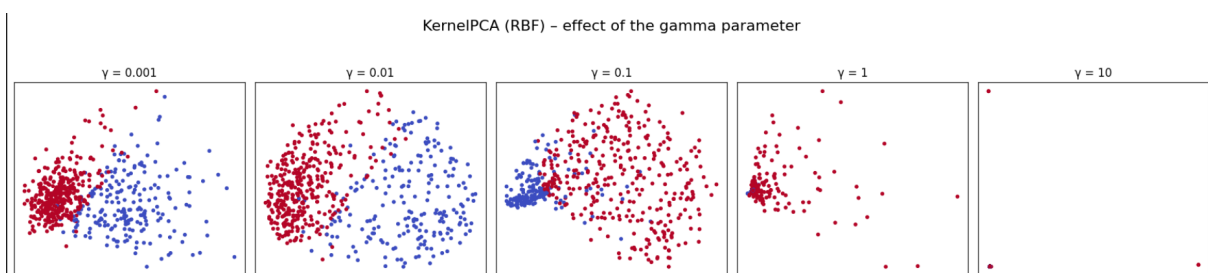


Fig. 1.4. Feature Heatmap

```

1 gammas = [0.00001, 0.0001, 0.001, 0.01, 0.1]
2 degrees = [2, 3, 5, 6, 7]
3
4
5 fig, axes = plt.subplots(len(degrees), len(gammas), figsize=(16, 9))
6 fig.suptitle("KernelPCA (poly) - influence of gamma and polynomial degree",
7             fontsize=16)
8
9 for i, degree in enumerate(degrees):
10     for j, gamma in enumerate(gammas):
11         kpca = KernelPCA(n_components=2, kernel='poly', gamma=gamma, degree=
12                          degree, coef0=1)
13         X_kpca = kpca.fit_transform(X_scaled)
14
15         ax = axes[i, j]
16         ax.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y, cmap='coolwarm', s=10)
17         ax.set_title(f"gamma={gamma}, deg={degree}")
18         ax.set_xticks([])
19         ax.set_yticks([])
20
21 plt.tight_layout(rect=[0, 0, 1, 0.95])
22 plt.show()

```

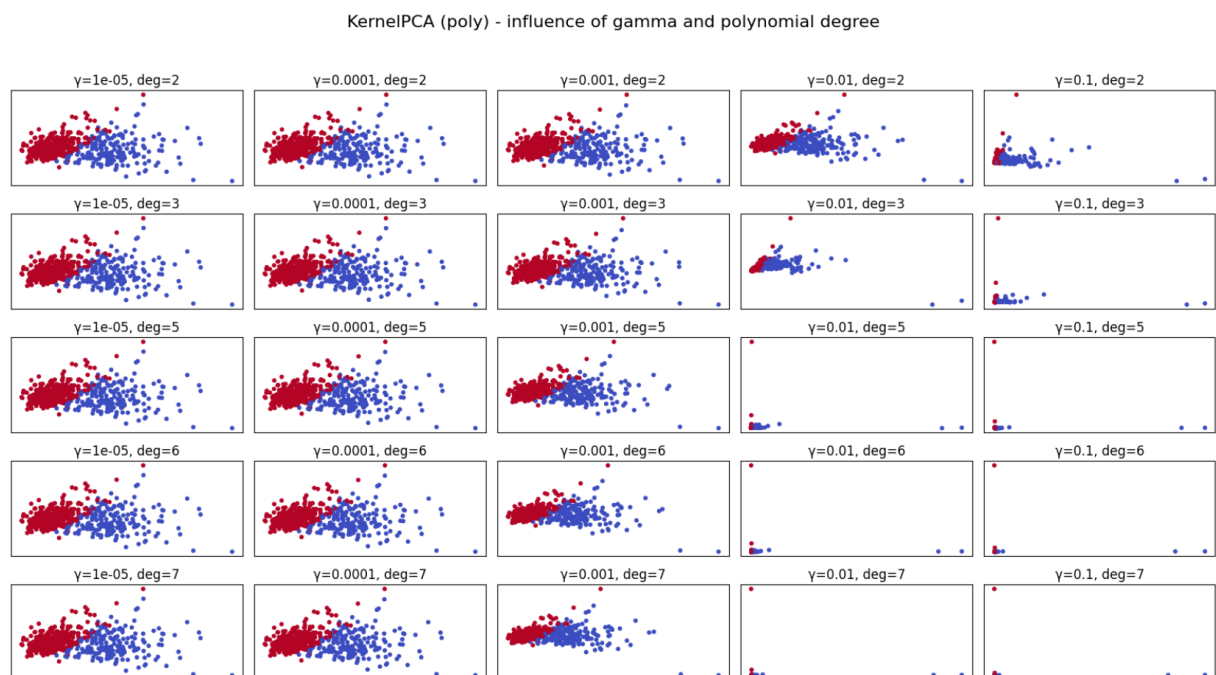


Fig. 1.5. Feature Heatmap

```

1 kpca = KernelPCA(n_components=2, kernel='cosine')
2 X_kpca = kpca.fit_transform(X_scaled)
3
4
5 plt.figure(figsize=(6, 5))
6 plt.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y, cmap='coolwarm', s=15)
7 plt.title("KernelPCA - kernel='cosine'")
8 plt.xlabel("PC1")
9 plt.ylabel("PC2")
10 plt.grid(True)
11 plt.tight_layout()
12 plt.show()

```

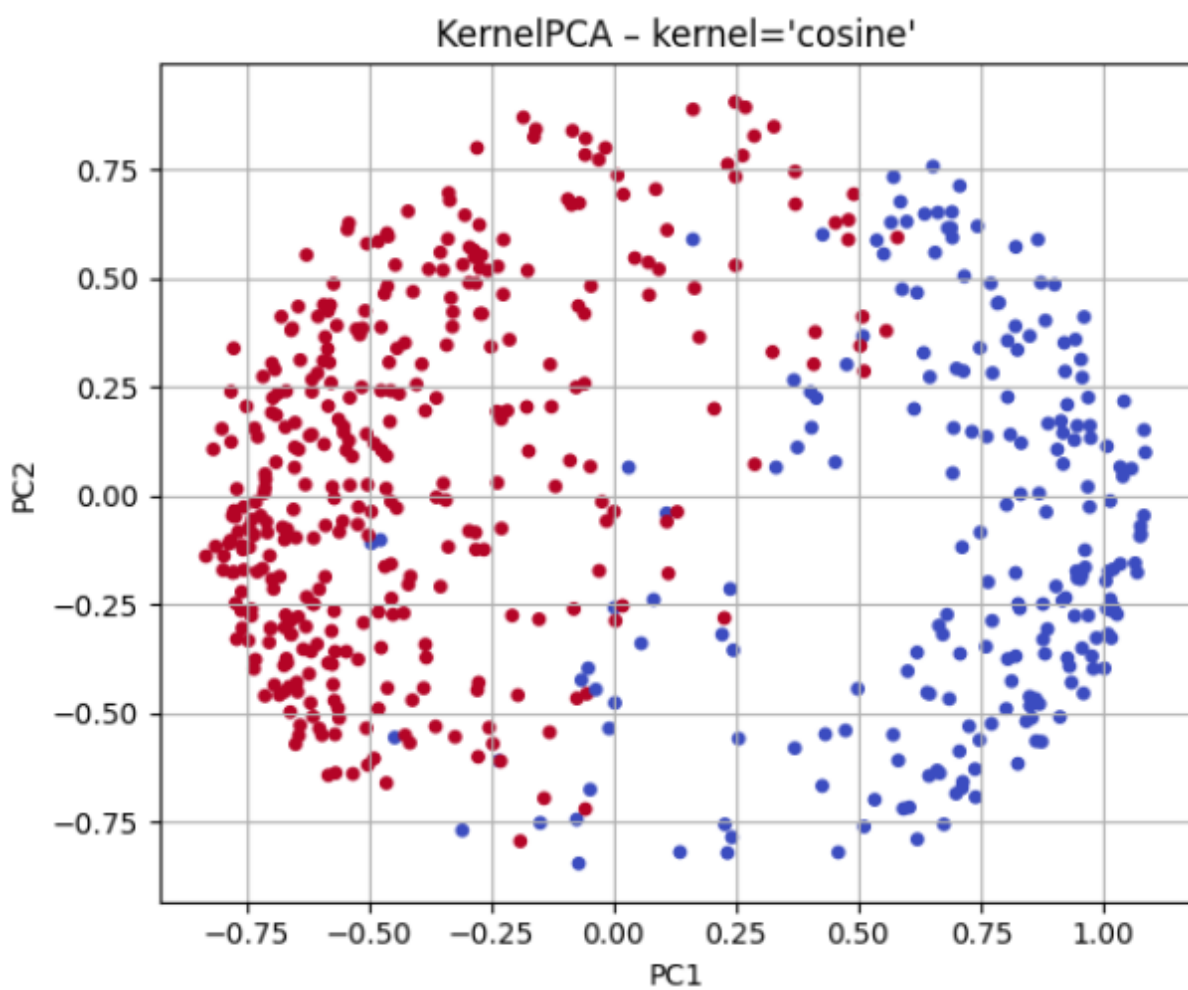


Fig. 1.6. Feature Heatmap

```
1 kpca = KernelPCA(n_components=2, kernel='linear')
2 X_kpca = kpca.fit_transform(X_scaled)
3
4
5 plt.figure(figsize=(6, 5))
6 plt.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y, cmap='coolwarm', s=15)
7 plt.title("KernelPCA - kernel='linear'")
8 plt.xlabel("PC1")
9 plt.ylabel("PC2")
10 plt.grid(True)
11 plt.tight_layout()
12 plt.show()
```

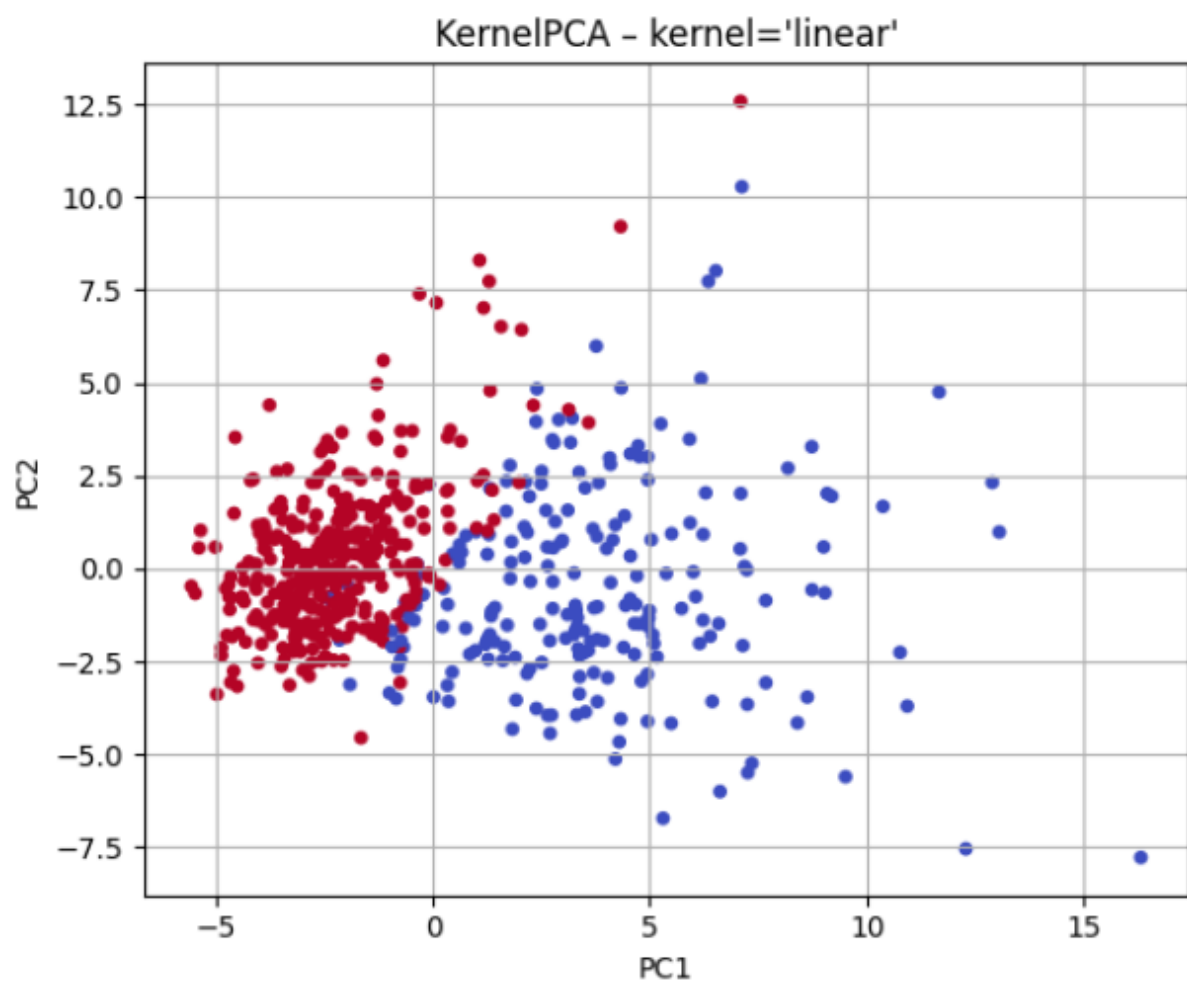


Fig. 1.7. Feature Heatmap

```

1  gammas = [0.001, 0.01, 0.1, 1]
2  coef0s = [-1, 0, 1]
3
4
5  fig, axes = plt.subplots(len(coef0s), len(gammas), figsize=(16, 9))
6  fig.suptitle("Effect_of_gamma_and_coef0_in_KernelPCA_(sigmoid)", fontsize=16)
7
8  for i, coef0 in enumerate(coef0s):
9      for j, gamma in enumerate(gammas):
10         kpca = KernelPCA(n_components=2, kernel='sigmoid', gamma=gamma, coef0=
            coef0)
11         X_kpca = kpca.fit_transform(X_scaled)
12
13         ax = axes[i, j]
14         ax.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y, cmap='coolwarm', s=10)
15         ax.set_title(f"gamma={gamma},_coef0={coef0}")
16         ax.set_xticks([])
17         ax.set_yticks([])
18
19 plt.tight_layout(rect=[0, 0, 1, 0.95])
20 plt.show()

```

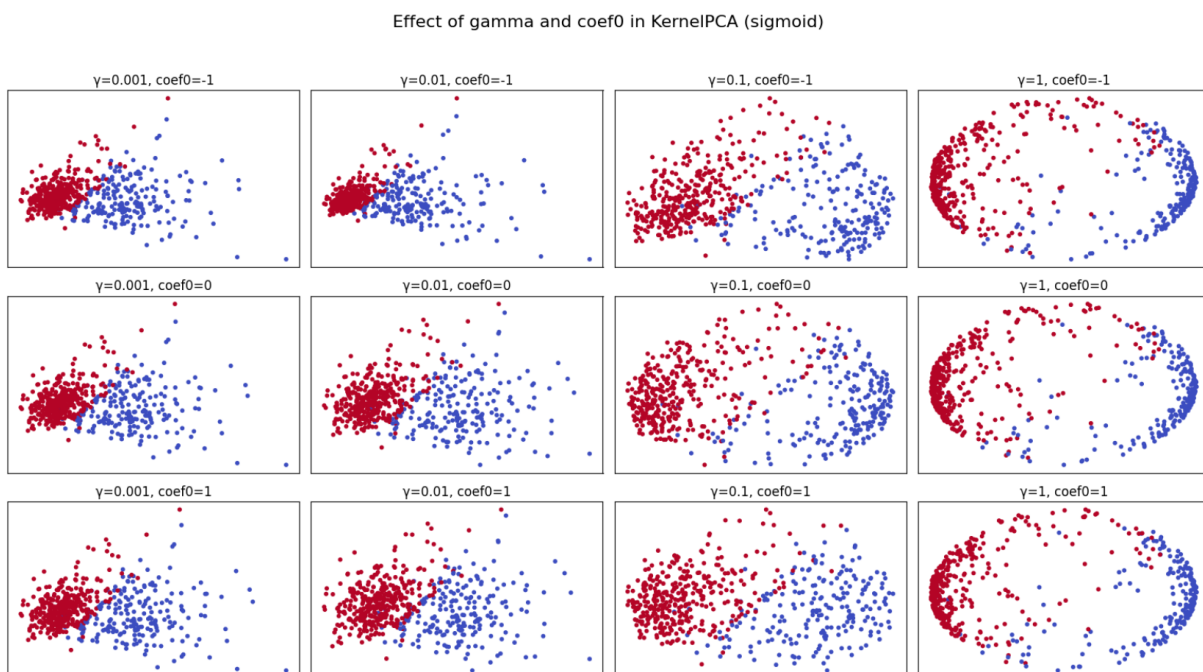


Fig. 1.8. Feature Heatmap

1.3. Classic PCA Analysis on MNIST (4 points)

- Download the MNIST data set (there is a function to load it in libraries such as scikit-learn, keras). It is a collection of black and white photos of handwritten digits with a resolution of 28x28 pixels, which together gives 784 dimensions.
- Try to visualize this dataset using PCA and KernelPCA, don't expect full separation of the data.
- Similar to the exercises, examine explained variance. Draw explained variance vs number of principal components plot.
- Find number of principal components for 99%, 95%, 90%, and 85% of explained variance.
- Draw some sample MNIST digits and from PCA do its original space. Make an inverse (odwrotną) transformation for number of components (komponentów) corresponding (odpowiadających) with explained variance shown above and draw the reconstructed images.
The idea of this exercise is to see visually how depending on the number of components some information is lost.
- Perform the same reconstruction using KernelPCA (make comparisons for the same components number).

1.3.1. Important:

Write a brief analysis (200–500 words) addressing:

1. What are the key strengths and limitations of the PCA variant you explored?
2. In what real-world scenarios would this approach be particularly useful?
3. How does the computational complexity scale with dataset size?
4. What improvements or modifications would you suggest to the method?

```

1 from sklearn.datasets import fetch_openml
2 mnist = fetch_openml('mnist_784', version=1, as_frame=False)
3
4
5 X, y = mnist.data, mnist.target.astype(int)
6
7
8 X_scaled = StandardScaler().fit_transform(X)
9
10 X_sample = X_scaled[:2000]
11 y_sample = y[:2000]
12
13 pca = PCA(n_components=2)
14 X_pca = pca.fit_transform(X_sample)
15
16 kpca = KernelPCA(n_components=2, kernel='rbf', gamma=0.01)
17 X_kpca = kpca.fit_transform(X_sample)
18
19 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(14, 6))
20
21 scatter1 = ax1.scatter(X_pca[:, 0], X_pca[:, 1], c=y_sample, cmap='tab10', s=3)
22 ax1.set_title("PCA - MNIST projection to 2D")
23
24 scatter2 = ax2.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y_sample, cmap='tab10', s=
    =3)
25 ax2.set_title("KernelPCA (rbf) - MNIST projection to 2D")
26
27 plt.tight_layout()
28 plt.show()

```

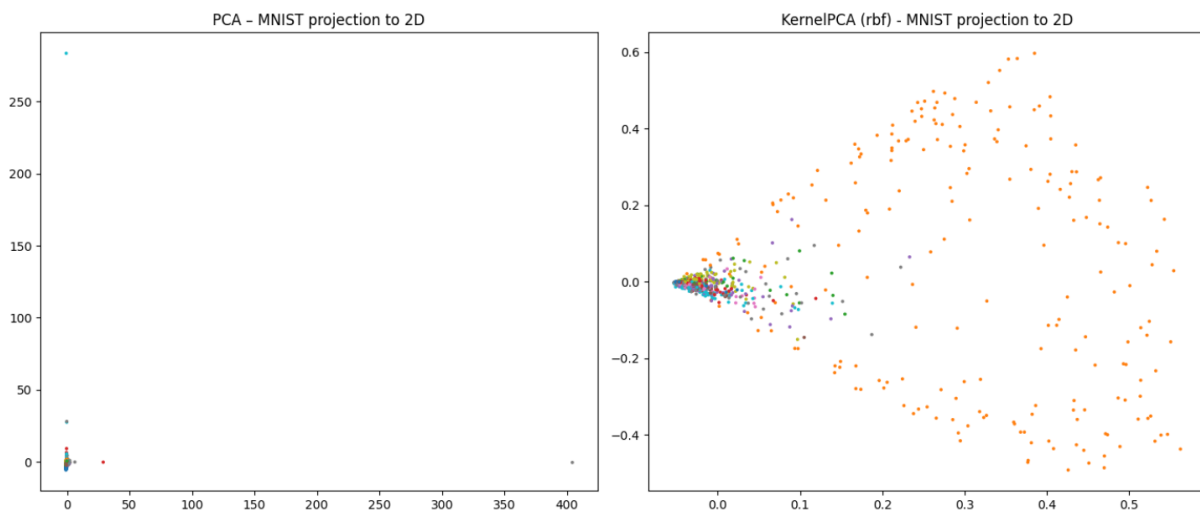


Fig. 1.9. Feature Heatmap

```

1 explained_variance = np.cumsum(pca.explained_variance_ratio_)
2
3
4 progi = [0.85, 0.90, 0.95, 0.99]
5 for prog in progi:
6     liczba = np.argmax(explained_variance >= prog) + 1
7     print(f'Dla {int(prog*100)}% wariancji potrzeba: {liczba} komponentów')
8
9 plt.figure(figsize=(10, 5))
10 plt.plot(range(1, len(explained_variance) + 1), explained_variance, marker='o')
11 plt.title('Skumulowana wariancja wyjaśniona przez PCA')
12 plt.xlabel('Liczba głównych składowych')
13 plt.ylabel('Skumulowana wyjaśniona wariancja')
14 plt.grid(True)
15 plt.tight_layout()
16 plt.show()

```

```

To explain ≥ 85% of variance, you need: 1 components
To explain ≥ 90% of variance, you need: 1 components
To explain ≥ 95% of variance, you need: 1 components
To explain ≥ 99% of variance, you need: 1 components

```

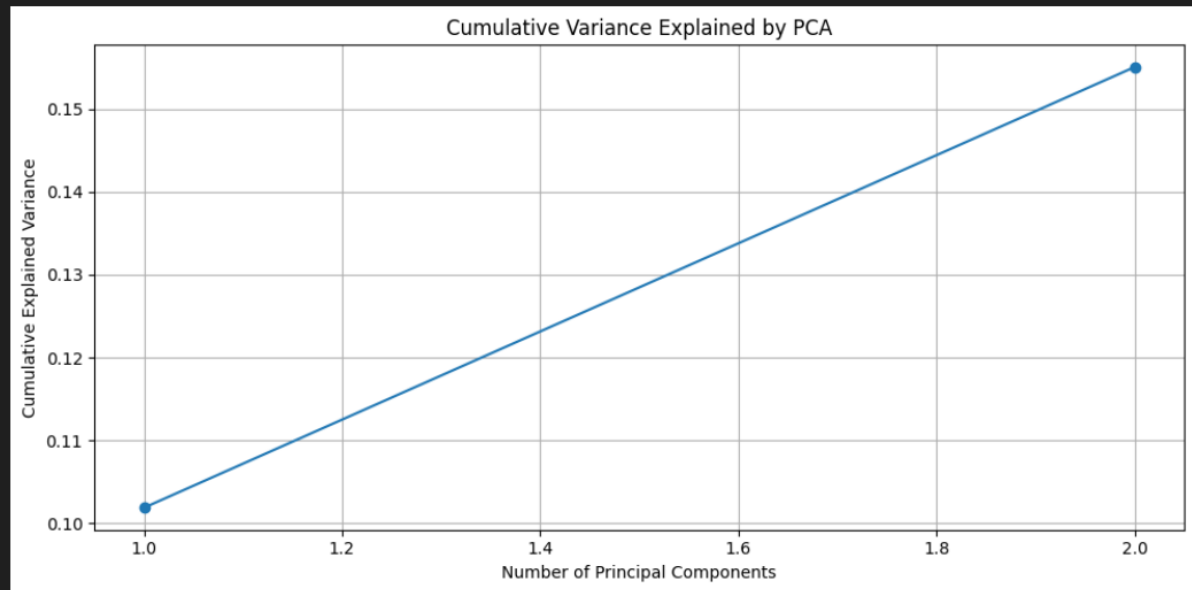


Fig. 1.10. Feature Heatmap


```

1 X = mnist.data[:10]
2 scaler = StandardScaler()
3 X_scaled = scaler.fit_transform(X)
4
5 n_components = 10
6 pca = PCA(n_components=n_components)
7 X_reduced = pca.fit_transform(X_scaled)
8 X_reconstructed = pca.inverse_transform(X_reduced)
9
10 X_original = scaler.inverse_transform(X_scaled)
11 X_reconstructed = scaler.inverse_transform(X_reconstructed)
12
13 fig, axes = plt.subplots(2, 10, figsize=(15, 4))
14 for i in range(10):
15     axes[0, i].imshow(X_original[i].reshape(28, 28), cmap="gray")
16     axes[0, i].axis("off")
17     axes[0, i].set_title("Original")
18
19     axes[1, i].imshow(X_reconstructed[i].reshape(28, 28), cmap="gray")
20     axes[1, i].axis("off")
21     axes[1, i].set_title(f"{n_components}_PCs")
22
23 plt.suptitle("MNIST_image_reconstruction_using_PCA")
24 plt.tight_layout()
25 plt.show()

```

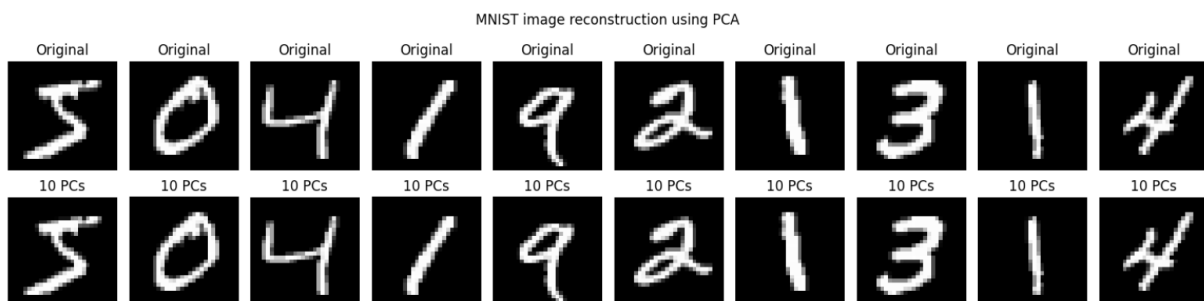


Fig. 1.11. Feature Heatmap

```

1 X = mnist.data[:10]
2 y = mnist.target[:10].astype(int)
3
4 scaler = StandardScaler()
5 X_scaled = scaler.fit_transform(X)
6
7 kpca = KernelPCA(n_components=10, kernel='rbf', gamma=0.01,
8                 fit_inverse_transform=True)
9 X_reduced = kpca.fit_transform(X_scaled)
10 X_reconstructed = scaler.inverse_transform(kpca.inverse_transform(X_reduced))
11
12 fig, axes = plt.subplots(2, 10, figsize=(15, 3))
13 for i in range(10):
14     axes[0, i].imshow(X[i].reshape(28, 28), cmap='gray')
15     axes[0, i].axis('off')
16     axes[0, i].set_title(f"Label: {y[i]}")
17
18     axes[1, i].imshow(X_reconstructed[i].reshape(28, 28), cmap='gray')
19     axes[1, i].axis('off')
20
21 plt.suptitle("KernelPCA: Original and Reconstruction (RBF, gamma=0.01)",
22             fontsize=14)
23 plt.tight_layout()
24 plt.show()

```

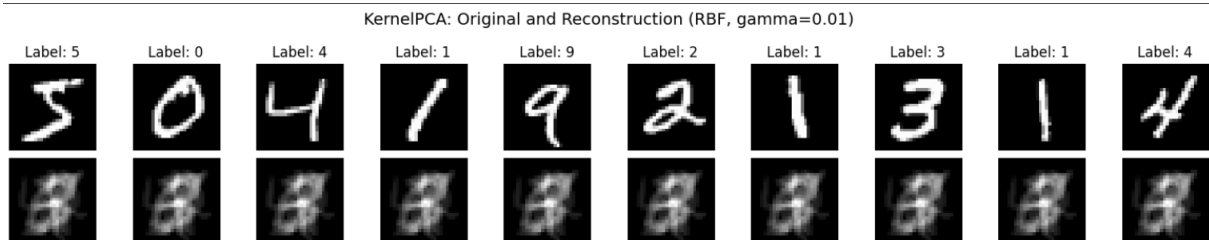


Fig. 1.12. Feature Heatmap

Important:

Write a brief analysis (200–500 words) addressing:

1. What are the key strengths and limitations of the PCA variant you explored?
2. In what real-world scenarios would this approach be particularly useful?
3. How does the computational complexity scale with dataset size?
4. What improvements or modifications would you suggest to the method?

Key strengths of the PCA variant I explored is that it extends traditional PCA to handle non-linear dimensionality reduction by implicitly mapping data into a higher-dimensional feature space. This is especially useful when the data is not linearly separable. The main limitation of this method is its high computational complexity. For example, when I tried to use the full MNIST dataset, the kernel matrix (K) had a size of $70,000 \times 70,000$, which is extremely memory-intensive.

This approach is especially useful in fields like image processing or finance, where datasets often contain many features and exhibit non-linear relationships. It is also valuable for visualizing high-dimensional data in two or three dimensions.

The computational complexity of PCA is determined by the algorithm used to compute the eigenvalues — in `scikit-learn` implementations, this is based on Singular Value Decomposition (SVD). For a data matrix $X \in \mathbb{R}^{n \times p}$, the complexity of full SVD is:

- $\mathcal{O}(np^2)$ when $n \geq p$
- $\mathcal{O}(n^2p)$ when $p > n$

In the case of **Kernel PCA**, the method first computes a kernel (similarity) matrix $K \in \mathbb{R}^{n \times n}$, where n is the number of samples. This leads to:

- **Memory complexity:** $\mathcal{O}(n^2)$
- **Time complexity** for eigen-decomposition: up to $\mathcal{O}(n^3)$

I think we should reduce the time complexity, but it's not easy. One possible solution is to use a lower number of samples from the original data, but they should be representative. Another modification I would suggest is automatic hyperparameter optimization, which could be performed using methods like Bayesian Optimization (BOGP), Tree-structured Parzen Estimator (TPE), SMAC, or genetic algorithms. However, this requires defining a loss function.

(*) Optional: 4. Modern PCA Applications (2 points)

Choose **ONE** of the following advanced tasks:

Option A: Robust PCA for Noisy Image Reconstruction

1. Take 10 sample images from MNIST
2. Add random noise and some structured outliers (e.g., salt and pepper noise)
3. Apply both standard PCA and Robust PCA for denoising (odszumienia)
4. Compare the results quantitatively (ilościowo) using metrics like PSNR, SSIM and qualitatively (jakościowo)

Option B: Incremental PCA for Large Datasets

1. Load the Fashion-MNIST dataset (a more complex alternative to MNIST)
2. Implement a streaming data scenario by processing the dataset in small batches
3. Apply Incremental PCA and compare its performance (both in terms of computation time and accuracy) with standard PCA
4. Visualize how the principal components evolve as more data is processed

Option C: PCA in Convolutional Autoencoder

1. Implement a simple convolutional autoencoder for MNIST digit reconstruction
2. Extract the bottleneck layer representations (latent space)
3. Compare these representations with traditional PCA by:
 - Visualizing the 2D projections of both
 - Measuring reconstruction quality at comparable compression ratios
 - Analyzing the feature importance in both approaches

```

1  X = mnist.data[:10].astype("float32") / 255.0
2  original_images = X.reshape(-1, 28, 28)
3  X_clean = X.reshape(10, -1)
4
5
6  noisy_images = np.array([random_noise(img, mode='s&p', amount=0.3) for img in
7      original_images])
8  X_noisy = noisy_images.reshape(10, -1)
9
10 pca = PCA(n_components=5)
11 X_pca_recon = pca.inverse_transform(pca.fit_transform(X_noisy))
12
13 # Robust PCA denoising
14 _, X_rpca_recon, _, _ = robust_pca(X_noisy, n_components=5)
15
16 # PSNR & SSIM
17 pca_psnr = [psnr(orig, recon) for orig, recon in zip(X_clean.reshape(-1, 28,
18     28), X_pca_recon.reshape(-1, 28, 28))]
19 rpca_psnr = [psnr(orig, recon) for orig, recon in zip(X_clean.reshape(-1, 28,
20     28), X_rpca_recon.reshape(-1, 28, 28))]
21
22 pca_ssim = [ssim(orig, recon, data_range=1.0) for orig, recon in zip(X_clean.
23     reshape(-1, 28, 28), X_pca_recon.reshape(-1, 28, 28))]
24 rpca_ssim = [ssim(orig, recon, data_range=1.0) for orig, recon in zip(X_clean.
25     reshape(-1, 28, 28), X_rpca_recon.reshape(-1, 28, 28))]
26
27 # Table
28 import pandas as pd
29 df = pd.DataFrame({
30     "Obraz": list(range(1, 11)),
31     "PSNR_PCA": pca_psnr,
32     "PSNR_RobustPCA": rpca_psnr,
33     "SSIM_PCA": pca_ssim,
34     "SSIM_RobustPCA": rpca_ssim
35 })
36 print(df.round(3))
37
38 # Wizualizacja
39 for i in range(3):
40     fig, axs = plt.subplots(1, 4, figsize=(10, 3))
41     axs[0].imshow(original_images[i], cmap="gray")
42     axs[0].set_title("Original")
43
44     axs[1].imshow(noisy_images[i], cmap="gray")
45     axs[1].set_title("Noisy")
46
47     axs[2].imshow(X_pca_recon[i].reshape(28, 28), cmap="gray")
48     axs[2].set_title("PCA")

```

```

45     axs[3].imshow(X_rpca_recon[i].reshape(28, 28), cmap="gray")
46     axs[3].set_title("Robust_PCA")
47
48     for ax in axs:
49         ax.axis("off")
50     plt.suptitle(f"Image_{i+1}")
51     plt.tight_layout()
52     plt.show()

```

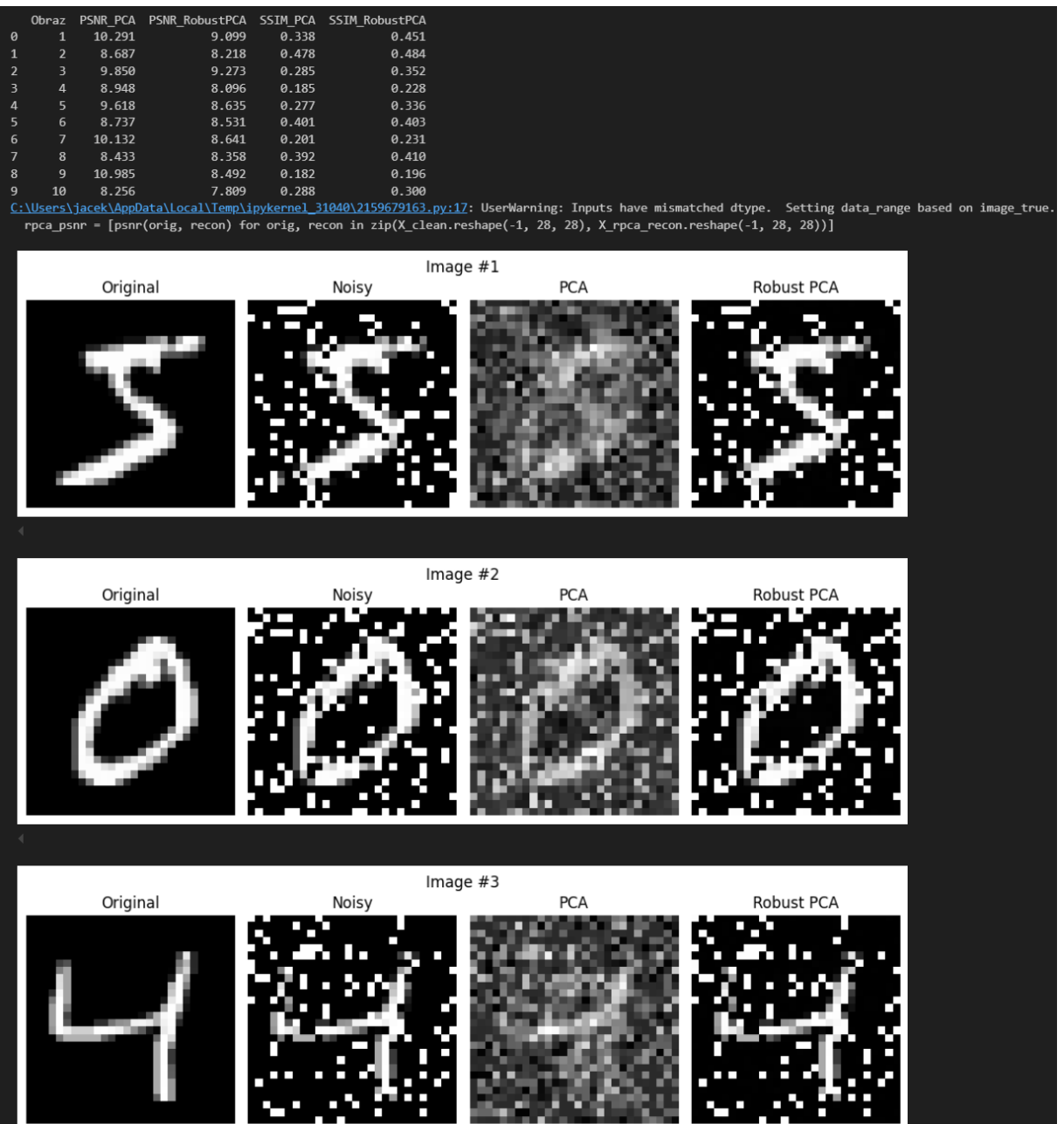


Fig. 1.13. Feature Heatmap

Analysis

PSNR is a metric that measures the pixel-wise difference between the original and reconstructed images. A higher PSNR indicates lower error and better quality, but it does not necessarily correspond to visual perception. On the other hand, SSIM is a perceptual metric that evaluates image similarity based on structure, luminance, and contrast. It better reflects human visual perception—higher SSIM values imply higher visual similarity.

We can see that PPCA achieves better PSNR than Robust PCA, but this doesn't reflect perceived visual quality. However, in SSIM, Robust PCA performs better than PCA. Since SSIM is more aligned with human visual perception, it is considered a more reliable metric. In my opinion, the image transformed by Robust PCA is clearer than the one produced by standard PCA.

A key strength of Robust PCA is that it preserves the important structure in data even under corruption, such as salt-and-pepper noise. It also achieves better visual quality in image reconstruction than standard PCA. However, the method comes with certain limitations. It typically requires more memory and longer convergence time due to iterative optimization. In real-world applications, Robust PCA proves useful in fields like medical imaging—where maintaining structural fidelity is essential—or in anomaly detection tasks.

Regarding computational complexity, Robust PCA tends to scale poorly with increasing dataset size, mainly due to its reliance on repeated SVD computations. This makes it more computationally intensive than traditional PCA. Reducing this complexity is challenging, but one potential direction is to work with a smaller subset of representative samples.

To improve the method further, automatic hyperparameter optimization could be introduced. Techniques such as Bayesian Optimization (BOGP), Tree-structured Parzen Estimator (TPE), SMAC, or genetic algorithms could be leveraged to fine-tune parameters like the regularization weights or convergence criteria. However, implementing such optimization requires a well-defined and effective loss function tailored to the reconstruction quality or task-specific requirements.