

**Міністерство освіти і науки України
Національний технічний університет України «КПІ» імені Ігоря Сікорського
Кафедра обчислювальної техніки ФІОТ**

**ЗВІТ
з лабораторної роботи №5
з навчальної дисципліни «Вступ до технології Data Science»**

Тема:

**РЕАЛІЗАЦІЯ МЕТОДІВ МАШИННОГО НАВЧАННЯ
(MACHINE LEARNING (ML))**

Виконав:

Студент 3 курсу кафедри ІІІ ФІОТ,
Навчальної групи ІІІ-11
Лошак В.І.

Перевірив:

Професор кафедри ОТ ФІОТ
Писарчук О.О.

Київ 2023

I. Мета роботи:

Виявити дослідити та узагальнити особливості аналізу даних з використанням методів та технологій машинного навчання (Machine Learning (ML))

II. Завдання:

Завдання III рівня складності 9 балів: реалізувати на вибір ТРИ з п'яти сформованих груп технічних вимог.

Група технічних вимог_1:

Реалізувати кластеризацію вхідних даних, отриманих Вами у ході виконання Дз_1, модельних та (або) реальних – на власний вибір. Методи Machine Learning з переліку: *k-means* (*k*-середніх); *Support Vector Machine* (машина опорних векторів); *k-nearest neighbors* (найближчих сусідів); ієрархічна кластеризація – для кластеризації обраних даних обрати самотійно. Провести аналіз та пояснення отриманих результатів, сформулювати висновки.

Група технічних вимог_2:

Реалізувати кластеризацію за кольоровою ознакою об'єктів на самотійно обраному цифровому зображенні. Методи Machine Learning з переліку: *k-means* (*k*-середніх); *Support Vector Machine* (машина опорних векторів); *k-nearest neighbors* (найближчих сусідів) – для кластеризації обраного зображення обрати самотійно.

За необхідності провести покращення якості зображення: зміна кольору; підвищення контрасту; фільтрація, тощо. Етапи покращення якості та кластеризації повинні забезпечувати виділення геометричних або кольорових ознак обраного на цифровому зображенні об'єкту для його подальшої ідентифікації. Провести аналіз отриманих результатів, сформулювати висновки.

Група технічних вимог_3:

Підрахувати кількість об'єктів на обраному цифровому зображенні. Об'єкти, що підлягають обрахунку обрати самотійно. Зміст етапів попередньої обробки зображень (корекція кольору, фільтрація, векторизація, кластеризація) має бути результатом R&D процесів, що конкретизується обраним зображенням і об'єктами для підрахунку. Провести аналіз отриманих результатів, сформулювати висновки.

III. Результати виконання лабораторної роботи.

1. Реалізувати кластеризацію вхідних даних

В лабораторній здійснюється обробка даних з електричного мікроскопу наведених за посиланням: <https://www.kaggle.com/competitions/data-science-bowl-2018/data>

Функції для зчитування Id та власне самих зображень.

```
1 def image_ids_in(root_dir, ignore=[]):
2     ids = []
3     for id in os.listdir(root_dir):
4         if id in ignore:
5             print('Skipping ID:', id)
6         else:
7             ids.append(id)
8     return ids

1 def read_image(image_id, space="rgb"):
2     image_file = STAGE1_TRAIN_IMAGE_PATTERN.format(image_id, image_id)
3     image = skimage.io.imread(image_file)
4     # Drop alpha which is not used
5     image = image[:, :, :3]
6     #hsv is hue saturation value format of the image
7     if space == "hsv":
8         image = skimage.color.rgb2hsv(image)
9     return image
```

Рис. 1 – Функції для читання вхідних даних в заданому форматі(rgb, hsv)

Для успішного виконання поставленого завдання необхідно кластеризувати дані за найактивнішим кольором що присутній у зображенні. Для кластеризації обрано алгоритм K-means що є простим і добре пристосований до задач такого типу. Щоб знайти оптимальні параметри цього алгоритму використовуємо метод Ліктя. Вибрано декілька зразків даних і знайдено оптимальний параметр К для них використовуючи вище наведений алгоритм.

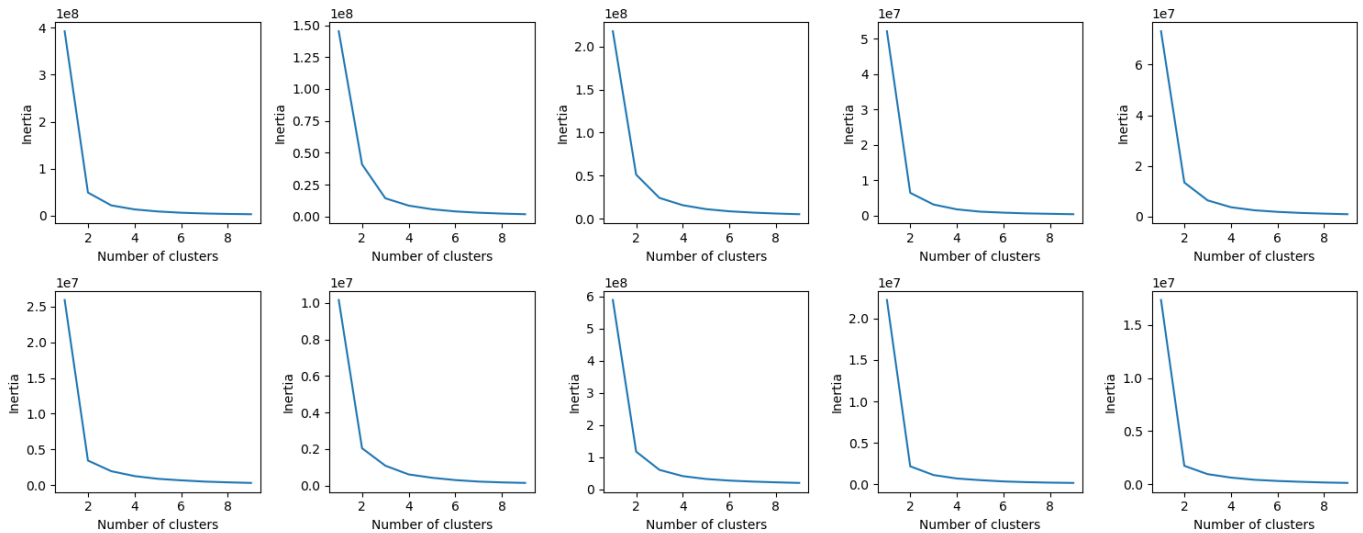


Рис. 2— Оптимальні значення параметра кластеризації.

```

1 def get_dominant_colors(img, top_colors=2):
2     img_l = img.reshape((img.shape[0] * img.shape[1], img.shape[2]))
3     clt = KMeans(n_clusters = top_colors)
4     clt.fit(img_l)
5     # grab the number of different clusters and create a histogram
6     # based on the number of pixels assigned to each cluster
7     numLabels = np.arange(0, len(np.unique(clt.labels_)) + 1)
8     (hist, _) = np.histogram(clt.labels_, bins = numLabels)
9     # normalize the histogram, such that it sums to one
10    hist = hist.astype("float")
11    hist /= hist.sum()
12    return clt.cluster_centers_, hist

```

Рис. 3— Функція що використовує оптимальний параметр кластеризації за замовчуванням. Використовується для екстракції доміантного кольору зображення.

Повертає гістограму топ кольорів у зображенні .

Для отримання залишкових даних про всі зображення у вигляді зображення таких як ширина та висота використовуємо:

```

1 def get_images_details(image_ids):
2     details = []
3     for image_id in image_ids:
4         image_hsv = read_image(image_id, space="hsv")
5         height, width, l = image_hsv.shape
6         dominant_colors_hsv, dominant_rates_hsv = get_dominant_colors(image_hsv, top_colors=1)
7         dominant_colors_hsv = dominant_colors_hsv.reshape(1, dominant_colors_hsv.shape[0] * dominant_colors_hsv.sh
8         info = [image_id, width, height, dominant_colors_hsv.squeeze()]
9         details.append(info)
10    return details

```

Рис. 4— Функція для завантаження даних в програму(все крім самого зображення)

```

1 def plot_elbow_kmeans(data, k_max):
2     scores = []
3     for k in range(1, k_max):
4         kmeans = KMeans(n_clusters=k)
5         kmeans.fit(data)
6         scores.append(kmeans.inertia_)
7
8
9     plt.plot(range(1, k_max), scores)
10    plt.xlabel('Number of clusters')
11    plt.ylabel('Inertia')

```

[159]

Рис. 5— Функція для побудови графіку «ліктя».

Щоб визначити який параметр кластеризації використовувати для безпосередньо кластеризації зображень в датасеті по групах побудовано графік.

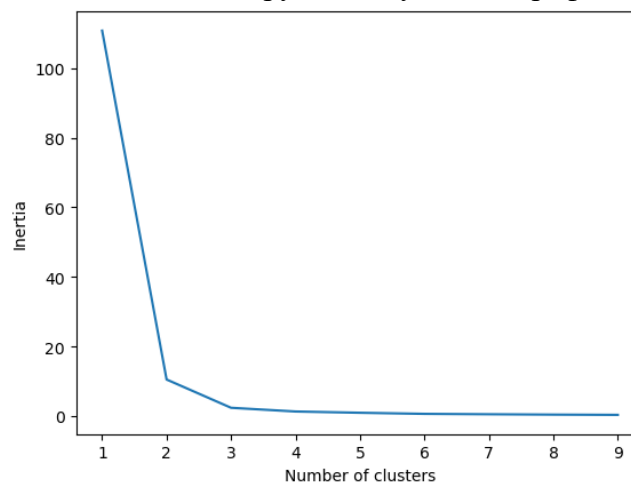


Рис. 6— Доцільно використовувати параметр 3 щоб отримати три групи зображень

```

1 kmeans = KMeans(n_clusters=3).fit(x)
2 clusters = kmeans.predict(X)
3 train_df[HSV_CLUSTER] = clusters
4 train_df

```

[31] Python

... c:\Users\vikto\miniconda3\envs\venv1\Lib\site-packages\sklearn\cluster_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=3.
warnings.warn(
...

	image_id	width	height	hsv_dominant	hsv_cluster
0	00071198d059ba7f5914a526d124d28e6d010c92466da2...	256	256	[0.0, 0.0, 0.02408716538373055]	0
1	003cee89357d9fe13516167fd67b609a164651b2193458...	256	256	[0.0, 0.0, 0.035367060642616806]	0
2	00ae65c1c6631ae6f2be1a449902976e6eb8483bf6b074...	320	256	[0.6969922721241343, 0.2417388976327306, 0.760...	1
3	0121d6759c5adb290c8e828fc882f37dfaf3663ec885c6...	320	256	[0.6439642758207026, 0.21017601724122487, 0.80...	1
4	01d44a26f6680c42ba94c9bc6339228579a95d0e2695b1...	320	256	[0.7435220601544179, 0.2324891863323699, 0.786...	1
...
665	fec226e45f49ab81ab71e0eaa1248ba09b56a328338dce...	256	256	[0.0, 0.0, 0.06294316310508558]	0
666	feffce59a1a3eb0a6a05992bb7423c39c7d52865846da3...	256	256	[0.0, 0.0, 0.0762871237362146]	0
667	ff3407842ada5bc18be79ae453e5bdaa1b68afc842fc22...	696	520	[0.0, 0.0, 0.06788620819681142]	0
668	ff3e512b5fb860e5855d0c05b6cf5a6bcc7792e4be1f0b...	256	256	[0.0, 0.0, 0.015257891486672656]	0
669	ff599c7301daa1f783924ac8be3ce7b42878f15a39c2d...	360	360	[0.0, 0.0, 0.0475662672476387]	0

Рис. 7— Датафрейм з інформацією про зображення

Для кожного кластеру виведено нну кількість зображень що належать цьому кластеру на екран.

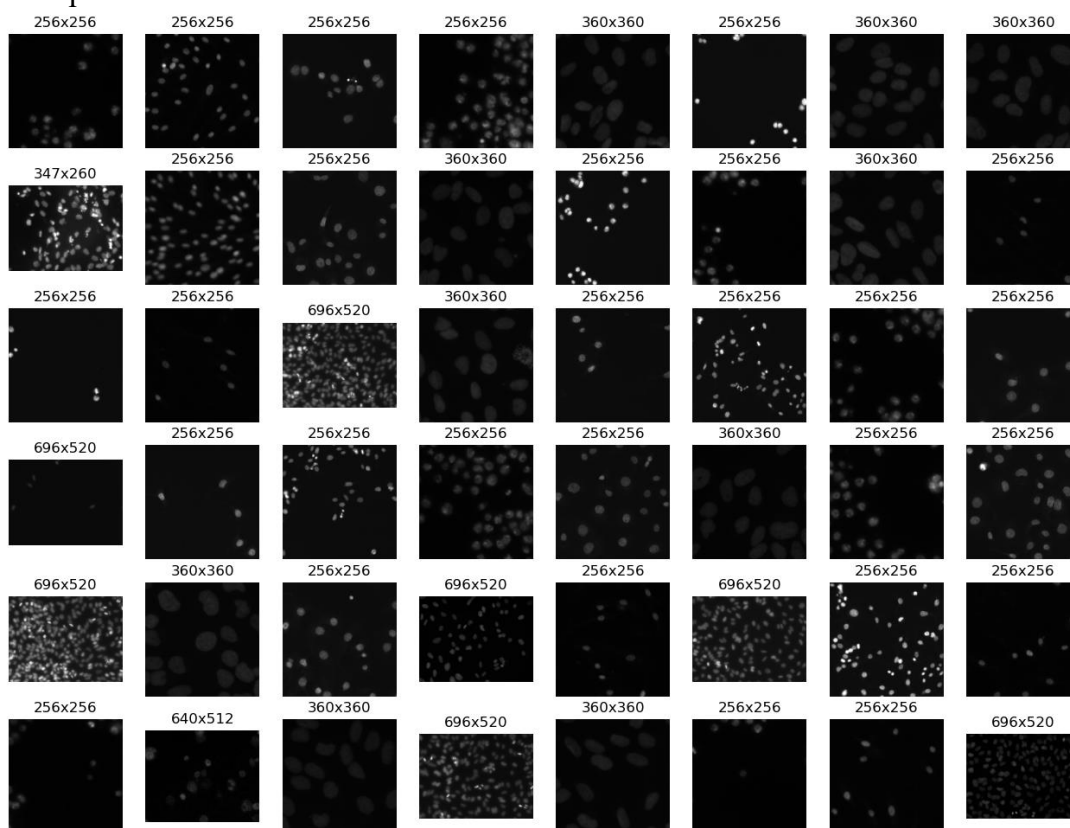


Рис. 8—Кластер 1.

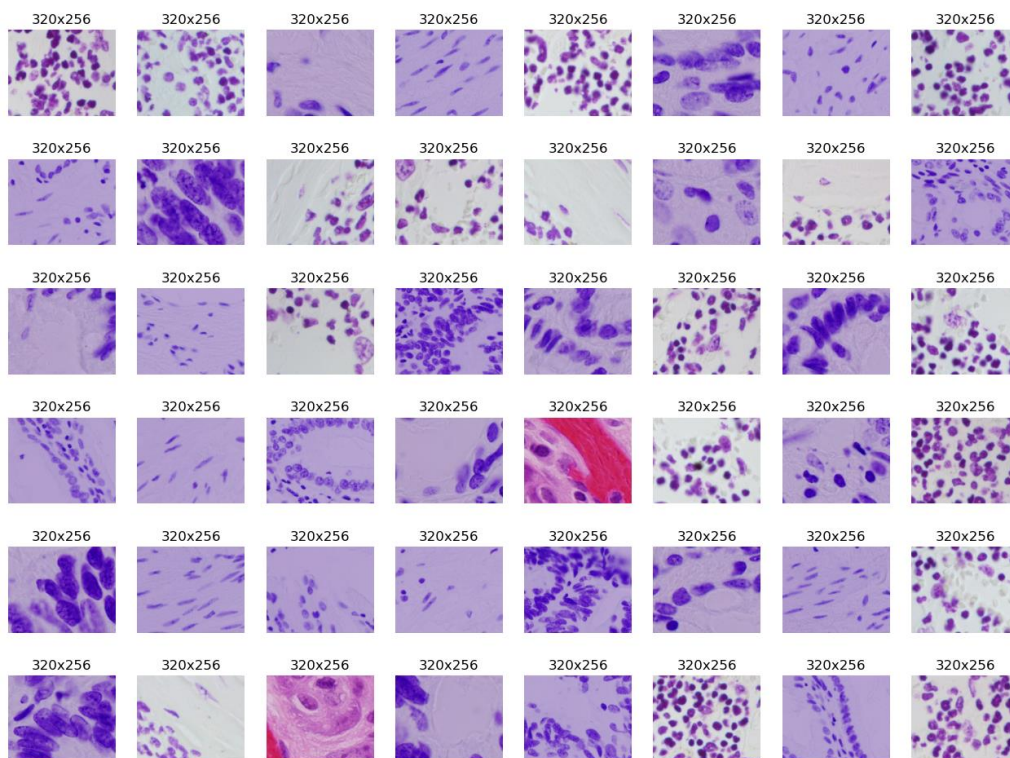


Рис. 9—Кластер 2.

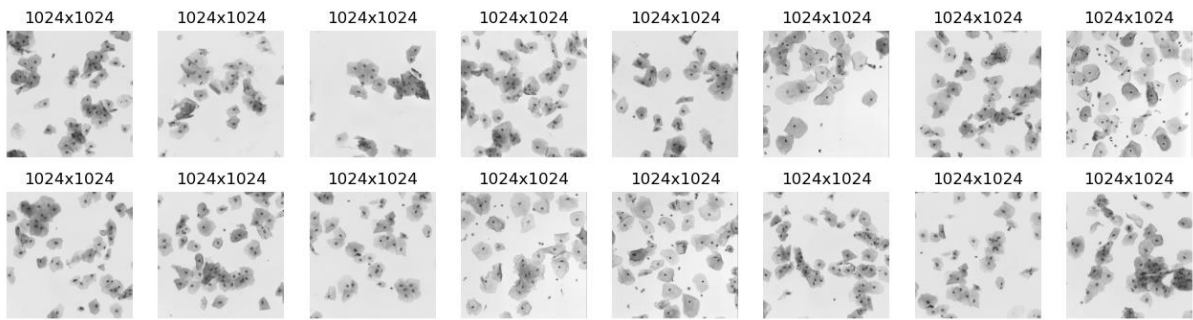


Рис. 10—Кластер 3

Отже як видно з результату кластеризація пройшла успішно. Присутній дизбаланс в класах але загалом зображення розподілені за домінуючим кольором зображення.

2. *Реалізувати кластеризацію за кольоровою ознакою об'єктів на самостійно обраному цифровому зображенні.*

Для аналізу використано перше зображення з другого кластеру в датасеті.

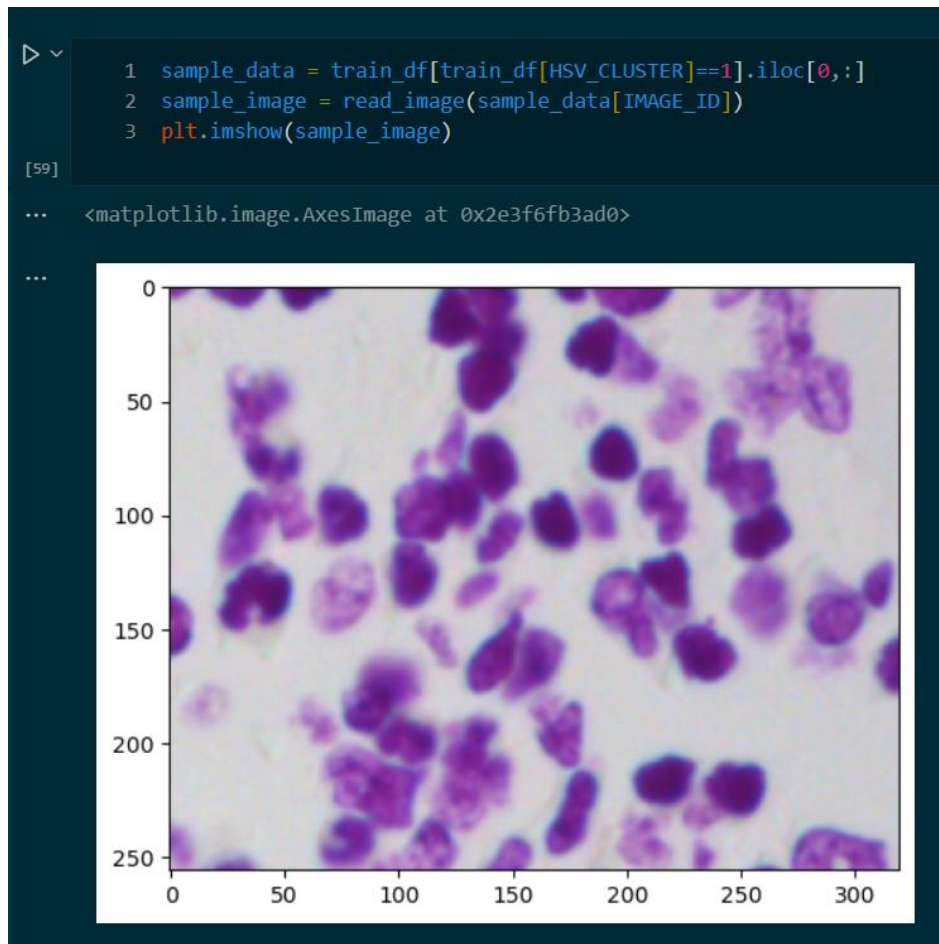


Рис. 11—Зображення для кластеризації кольорів.

Щоб краще розуміти що саме відбувається зображено всі пікселі присутні в зображенні на 3D графіку шкалами якого є відповідно R G та B показники зображення

```

1 def plot_pixels(image_resized, title, colors=None, N=10000):
2     if colors is None:
3         colors = image_resized
4
5     # choose a random subset
6     rng = np.random.RandomState(0)
7     i = rng.permutation(image_resized.shape[0])[:N]
8     colors = colors[i]
9     R, G, B = image_resized[i].T
10
11     fig = plt.figure(figsize=(10, 8))
12     ax = fig.add_subplot(111, projection='3d')
13     ax.scatter(R, G, B, c=colors, marker='.')
14     ax.set(xlabel='Red', ylabel='Green', zlabel='Blue')
15     ax.set_title(title)

```

Рис. 12— Функція для побудови графіку пікселів за кольорами.

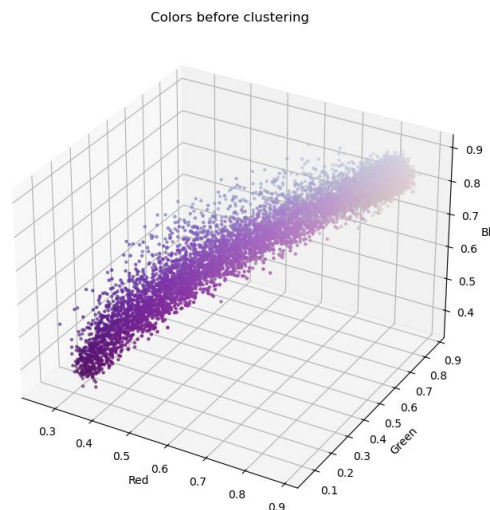


Рис. 13— перед застосуванням кластеризації пікселі зображення мають вигляд.

Для того щоб обрати оптимальний параметр кластеризації знову використано лікоть.

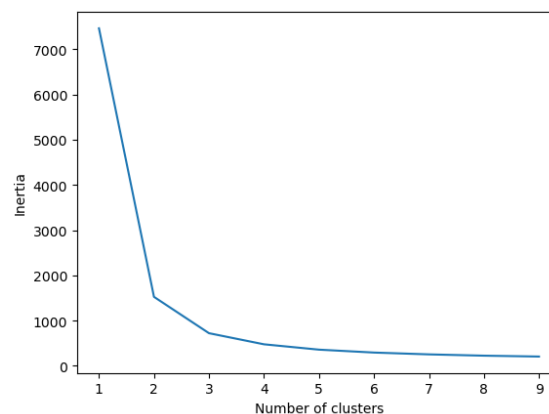


Рис. 14— оптимальний параметр для кластеризації 3

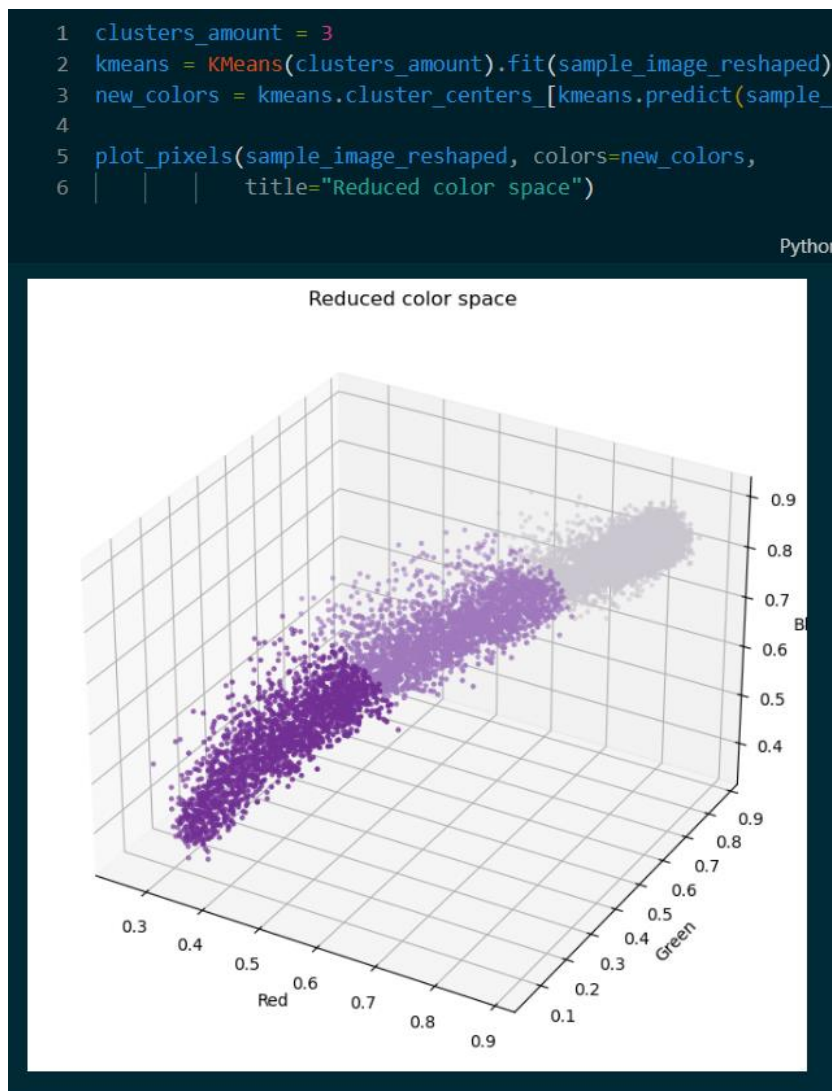


Рис. 15 —Пікселі зображення забарвлені кольором свого кластера.



Рис. 16 — Реконструйоване зображення в зжатому просторі кольорів.

Залежно від кількості кластерів зображення буде використовувати різну кількість базових кольорів. В ході цієї лабораторної роботи використано алгоритмом підбору параметрів K-means для зменшення втрат під час стиснення зображення.

3. Підрахувати кількість об'єктів на обраному цифровому зображенні.

Об'єкти, що підлягають обрахунку обрати самостійно.

Об'єктом підрахунку є клітини що відображені на індивідуальному зображенні. Для обробки використано те саме зображення що і в попередньому завданні. Для підрахунку обрано використовувати алгоритм K-means наступним чином. До dataframe використаного в попередньому завданні для кластеризації за кольорами додано також координати кожного пікселя в зображенні. Після цього локалізовано всі пікселі зображення що не належать до заднього плану в окремий датафрейм. На цьому датафреймі що містить лише пікселі клітин проведено пошук за відстанню. Таким чином знайдено кластери для пікселів клітин базуючись на їх близькості. Після підбору оптимального параметра кластеризації цей параметр за значенням дорівнює кількості клітин що видно на фото.



Рис. 17 — три базові кольори кластеризованого зображення.

```
1 pixels_df = pd.DataFrame(new_colors)
2 pixels_df.columns = ['R', 'G', 'B']
3 pixels_df['Position'] = [[i//sample_data
4
5 pixels_df
```

[153]

	R	G	B	Position
0	0.446170	0.186919	0.58271	[0, 0]
1	0.446170	0.186919	0.58271	[0, 1]
2	0.446170	0.186919	0.58271	[0, 2]
3	0.446170	0.186919	0.58271	[0, 3]
4	0.446170	0.186919	0.58271	[0, 4]
...
81915	0.622792	0.469767	0.74089	[255, 315]
81916	0.797796	0.781329	0.81663	[255, 316]
81917	0.797796	0.781329	0.81663	[255, 317]
81918	0.797796	0.781329	0.81663	[255, 318]
81919	0.797796	0.781329	0.81663	[255, 319]

Рис. 18 — dataframe з доданою інформацією про положення пікселя.

	R	G	B	Position
0	0.446170	0.186919	0.58271	[0, 0]
1	0.446170	0.186919	0.58271	[0, 1]
2	0.446170	0.186919	0.58271	[0, 2]
3	0.446170	0.186919	0.58271	[0, 3]
4	0.446170	0.186919	0.58271	[0, 4]
...
81911	0.622792	0.469767	0.74089	[255, 311]
81912	0.622792	0.469767	0.74089	[255, 312]
81913	0.622792	0.469767	0.74089	[255, 313]
81914	0.622792	0.469767	0.74089	[255, 314]
81915	0.622792	0.469767	0.74089	[255, 315]

32470 rows × 4 columns

Рис. 19 — Пікселі що є частиною клітини.

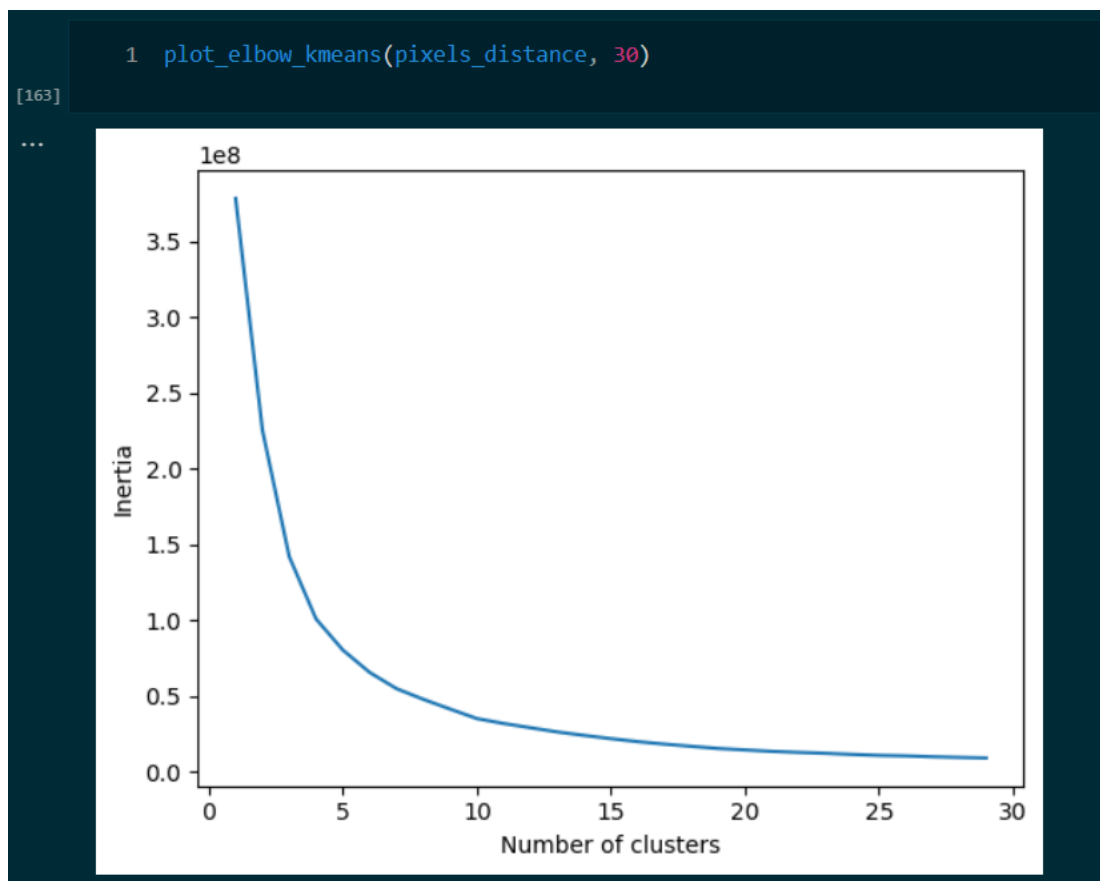


Рис. 20 — Пошук оптимального параметра кластеризації з використанням ліктя.
Оптимальним є значення параметра 30

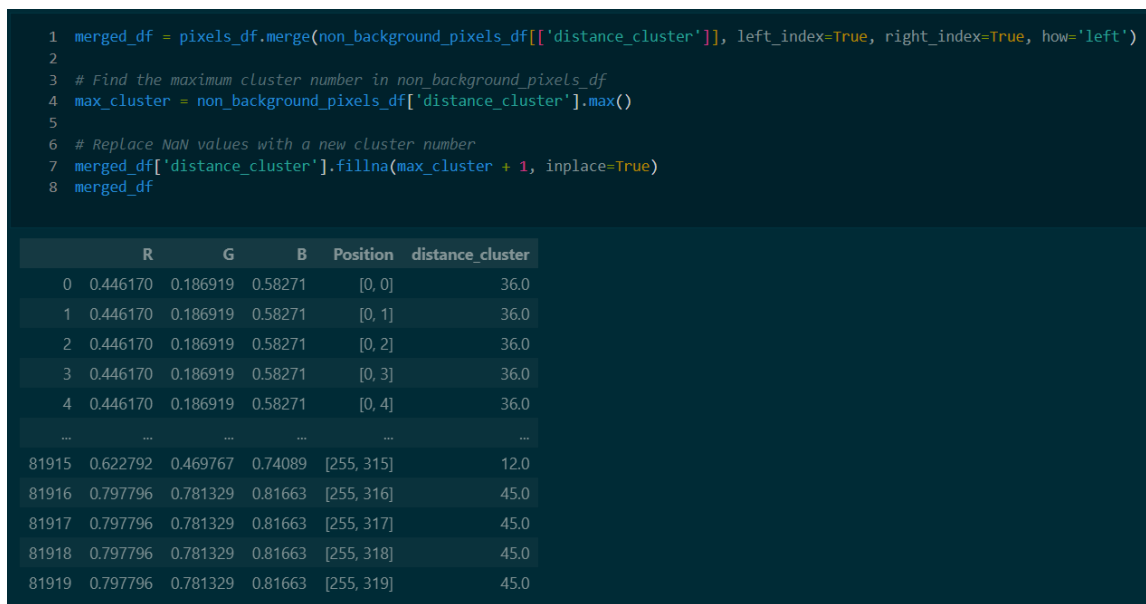


Рис. 20 — Комбінація dataframe без та з кластеризацією за відстанню в один. Присвоєння пікселям фона окремого кластеру.



Рис. 21 — Візуалізація кластерів за відстанню у вигляді зображення з випадкових кольорів для кожного окремо взятого кластера

IV. Висновки.

Результати лабораторної роботи показують потенціал і багатофункціональність застосування базового алгоритму K-means для досить складних процесів обробки зображень. Як для кластиризації зображень, так і для виділення features конкретно взятого зображення. Також його можна використовувати для побудови алгоритмів компресії, для стилізації зображень і в безлічі інших галузей.

Виконав: студент ФІОТ Лошак В.І. ПІ-11