# AIN 433 - Introduction to Computer Vision Laboratory

**HACETTEPE UNIVERSITY**

Department of Computer Engineering

# Dimensionality Reduction, Image Retrieval, and Classification

**Due Date: 08.11.2023 - Wednesday (23:00:00)**

## Overview

In this assignment, goal is to get familiar with the representation of images and processing them. This assignment consists of three parts. For the first part you will implement and analyze dimension reduction algorithm. You will use covariance matrix based PCA to represent images in lower dimensions. For the second part you will implement a simple image retrieval algorithm with different representations. For the third part you will implement a classification algorithm based on logistic regression.

**PART 1: Dimensionality Reduction with PCA**

In this part of the assignment you are expected to reduce the dimension of the given images' representations and plot them. You will use PCA algorithm for dimension reduction.

- You are given 256x256 images in "Dataset1" folder. You have 11 images.

- Use intensity values as feature of corresponding image. Represent your dataset with a matrix, M , that has a shape 65536x11.

- Normalize your representation via subtracting mean of your dataset from images therefore you need to calculate the mean of each image and generate a mean vector that has a shape of 1x11. Normalized 65536x11 shaped matrix is called $D$.

- Then, calculate a 11x11 covariance matrix via $Cov = D^T D$.

- The next step is to calculate eigenvalues and eigenvectors of the covariance matrix. You can use a built-in function for that purpose.

- Sort eigenvectors($V$) according to their corresponding eigenvalues (in descending order). Choose the first n eigenvectors.

- Project all images to the lower dimensional space via $V^T D$.

- Plot the results and comment.

**You will implement your PCA algorithm using Python. Your first aim is to represent each image with 3 points**, which means your new feature vector will be projected to 1x3 instead of 1x65536. You have to plot your result and comment about result. You can also reduce the dimension to 2 and 1.

**PART 2: Image Retrieval**

In this part of the assignment you will implement a simple image retrieval system. You have 300 images belongs to 10 different classes (in Dataset2 folder). Your aim is to rank images according to the relevance with the given query image (from QUERY_IMAGES folder). To calculate the similarity between images, you have represent them with feature vectors. You will use two different representation for this part of the assignment;

- PCA coefficients

- Color histogram of image

For each image:

- Calculate the feature vectors of the image,

- Calculate distance between selected image (image's features) and all the other images (all other images' features),

- Generate a ranked list based on the distance between them.

- Use k-means algorithm to cluster the images.

- Evaluate results via mean average precision (MAP) metric for each feature you used. You have to report MAP for each class and dataset.

You will use Euclidean distance to calculate the distance between images. Although the histograms are not in the Euclidean space, Euclidean Distance will still measure the "similarity" with an acceptable rate. Please explain what similarity means in this context. Explain what is the distance between two images, and why is it related to similarity.

Think about what the distance function outputs and show 10 most similar images for an image in each class. Comment about the results and features according to retrieved images. You can use different features, different color spaces for color histogram etc. Please remember; you have to explain "why". No external code is allowed to calculate histogram or distance. You can use built-in filtering functions.

# Metric Calculation

- Precision: fraction of retrieved images that are relevant

$$Precision = \frac{\#(relevant\ images\ retrieved)}{\#(retrieved\ images)}$$

- Recall: fraction of relevant image that are retrieved

$$Recall = \frac{\#(relevant\ images\ retrieved)}{\#(relevant\ images)}$$

- MAP calculation for one query:

    - start with rank (K) 1 and go with one-rank-at-time
    - if the image at rank K is relevant, calculate precision for corresponding K
    - take the average of all precision values at all relevant images

- MAP: average precision averaged across a set of queries

- MAP for each class: use query images for each class and calculate MAP for them separately.

- To understand calculation of the metric you can look at this presentation.

**PART 3: Classification**

In this part of the assignment you will implement a simple classification system. You have 300 images belongs to 10 different classes (in Dataset2 folder). Your aim is to implement and train a logistic regression model (you must do your very own, it is forbidden to use any libraries that do logistic regression directly) and show the results by testing it over query images (from QUERY_IMAGES folder). Then you have to classify these by using logistic regression. It is enough to classify two classes (selection is up to you) for getting 10 points from this part, but if you classify all of the classes, you will also get 15 points of bonus from this part. Comment about the stages and your results at the report, also comment about your approach to apply logistic regression over all the classes if you did the bonus part.

# The Report

You are expected to implement these three parts using Python 3, apply your implementation to a set of images, evaluate the outputs and report your observations.

## What should you write in the report?

- Explain dimension reduction, importance etc.
- Explain the logic of the given algorithm, why you sorted the eigenvalues and choose eigenvectors accordingly.
- Plot 3 dimensional data (result of your implementation) and comment about the results. Analyze your observations.
- Give a brief description of histogram and PCA components.
- Explain the logic behind the given image retrieval algorithm.
- Comment about the results considering features, why the result is good/bad.
- Observe retrieved images and comment about the advantages/disadvantages of the given algorithm and representation methods.
- You can use different color spaces for color histogram and observe the result.
- Calculate MAP metric for each experiment and comment about reasons the performance of the method and features that you used.
- Explain the logic behind the logistic regression algorithm.
- Comment about the results for classification, why the result is good/bad.
- Observe classified images and comment about advantages/disadvantages of the classification method.
- Comment about your approach for the bonus part if you did.
- You are encouraged to write your report in LaTeX
- You should give visual results by using a table structure.

# What to Hand In

Your submission must contain following:

- README.txt *(Text file containing the details about your implementation, how to run your code, which libraries have to installed, the organization of your code, functions etc. The template must be followed will be shared.)*
- src/ *(directory containing all your code)*
- Report.pdf

File hierarchy is as follows and must be zipped before submitted (Not .rar, only not compressed .zip files because the submit system just supports .zip files).

```
- b<StudentID>.zip
    - README.txt
    - Report.pdf
    - <src>
      - *.*
```

**Note that you MUST exactly score ONE from the submit system, otherwise you will have 20% of point deduction even if your hierarchy is correct. (For MacOS users, you can check the Piazza post for the script that zips the folder without any extra content which causes getting zero from submit)**

# Grading

The assignment will be graded out of 100:

- **45 % (part 1):** CODE: 33 REPORT: 12

- **45 % (part 2):** CODE: 33 REPORT: 12

- **10 % (part 3):** CODE: 7 REPORT: 3

- **15 % (part 3 - bonus):** CODE: 11 REPORT: 4

**The score for the report will be multiplied by your code score for that part which is divided by maximum score that can be taken from code part. Each part will be evaluated individually. For example, for first part say that some scored 22 for code part and 9 for report part, his/her final score for that part will be calculated as follows: 22+(22/33)\*9=28**

**You MUST use justify -iki yana yasla in Turkish- page alignment and passive voice at your report, otherwise 20% of your score will be deducted FOR EACH VIOLATION for the relevant part of your report. Note that also your report's alignment must be well designed, otherwise you may also face with some point deductions for bad alignments/designs**

**You MUST write comments to your code as necessary and also your code MUST be readable.**

# Academic Integrity

All work on assignments must be done individually unless stated otherwise. You are encouraged to discuss with your classmates about the given assignments, but these discussions should be carried out in an abstract way. That is, discussions related to a particular solution to a specific problem (either in actual code or in the pseudocode) will not be tolerated. In short, turning in someone else's work, in whole or in part, as your own will be considered as a violation of academic integrity. Please note that the former condition also holds for the material found on the web as everything on the web has been written by someone else.

# References

[1] `https://en.wikipedia.org/wiki/Principal_component_analysis`
[2] Ian T. Jolliffe and Jorge Cadima, "Principal Component Analysis: A Review and Recent Developments", Adaptive data analysis: theory and applications, volume 374, issue 2065, `https://doi.org/10.1098/rsta.2015.0202`, 2016.
[3] `https://en.wikipedia.org/wiki/Color_histogram`
[4] `https://en.wikipedia.org/wiki/K-means_clustering`
[5] J. A. Hartigan and M. A. Wong, "A K-Means Clustering Algorithm", Journal of the Royal Statistical Society. Series C (Applied Statistics), Vol. 28, No. 1, pp. 100-108, 1979.