

# Covid-19 search engine

Marco Ghezzi, Matteo Limoncini

October 21, 2022

## Abstract

The aim of this project is to create an information retrieval system that make it possible to search and execute queries over a large dataset of medical images; the main purpose of the system is finding the images most similar to a given one. We applied two different machine learning techniques in order to build our retrieval system: supervised learning, using an artificial neural network to build a classification system, and unsupervised learning, in particular clustering. These different techniques were compared by analyzing different performance metrics and computational time. As a final result, we noted that employing a Convolutional Neural Network for classification (supervised learning) produces the best results in terms of precision and time efficiency.

## 1 Introduction

Content-based image retrieval (CBIR) is an image search technique designed to find images that are most similar to a given query: it measures the similarity of two images based on the similarity of the properties of their visual components. Even though this method is fully automated, it is affected by the “semantic gap”, that is the gap between the features extracted from an image and the high-level concepts contained in the image. Thus, feature extraction is the foundation for CBIR.

In this project we tackle the problem of creating a system able to search and make queries over a large dataset of images, retrieving the most similar images of a given query. The main goal is to build a content-based medical image retrieval system using different approaches and comparing their results based on classification metrics and computational time.

## 2 Research question and methodology

Two different Machine Learning approaches were used: the first method is performing supervised training with binary classification, the second one is using unsupervised learning with clustering.

The most important stage before applying one of the clustering algorithm is the feature extraction stage in which a visual concept is converted to a numerical

form: these features could be in the form of global features (i.e. color, shape, texture ...) or local features (like corners blobs or edges) that are invariant against scale, translation and rotation changes.

After the feature extraction stage, clustering is performed. Clustering is an unsupervised learning algorithm that gathers image descriptors into a single group that semantically differs from other groups: the most used clustering algorithm in CBIR is  $k$ -means, even though it fails in handling outliers and noisy data.

In this work we also decided to address the problem of classification since we had prior knowledge of the labels of the images: in order to classify the images in a supervised learning way, we decided to develop a Convolutional Neural Network, an Artificial Neural Network usually employed to analyze visual objects.

The final step is the similarity measurement between the extracted features from the query image and all other images in the dataset to retrieve the most relevant images. There are different metrics that could be used: in this work we considered the Euclidean distance [1] and the cosine distance [2].

$$d(a, b) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2} \quad (1)$$

$$d(a, b) = \frac{A \cdot B}{||A|| \cdot ||B||} \quad (2)$$

## 2.1 Data

The data collected in this dataset are computed tomography scans taken from real patients in São Paulo, Brazil: every image is paired with a label to specify if the person associated with the image was tested positive for Covid-19.

The dataset consists of 2482 images of different pixel sizes; the labels are approximately equally divided.

## 2.2 Image Preprocessing

Every image is associated with a label and they've been divided in two directories: the first containing all the images labeled as Covid and the second containing all the images labeled as non Covid. Starting from the images in these directories, three different datasets were created, since we're dealing with a classification task: one for training data, one for validation data and the last one for test data, using a split ratio of 0.8, 0.1 and 0.1 respectively. The size of the images was specified during the creation of these datasets in order to find the best trade-off between result accuracy and training time: in fact this resize operation is necessary because all the images in the dataset have different sizes. Then all the pixel values were normalized between 0 and 1 since the computation of high numeric values may become more complex.

Different Machine Learning approaches were evaluated in order to create a system that is able to answer user's queries with the most precision and less time consuming.

## 2.3 Classification

Since this is an image binary classification problem, we have considered using an Artificial Neural Network, in particular a Convolutional Neural Network.

A Convolutional Neural Network (CNN) is an Artificial Neural Network suited for image classification, image recognition and medical image analysis. It consists of an input layer, hidden layers and an output layer. Three main types of hidden layers were used to build our neural network: Convolutional Layer <sup>1</sup>, Pooling Layer <sup>2</sup> and Fully-Connected Dense Layer <sup>3</sup>.

Using this CNN, we obtained interesting results: this network is able to correctly classify a new image with an accuracy of 82% with less than a minute of training time over 1986 images.

In Figure 1 we can see that the training of the CNN produces overfitting. This is due to the fact that the network learns details, like image brightness and lung shape, that are not useful for correctly classifying whether an image contains Covid or not.

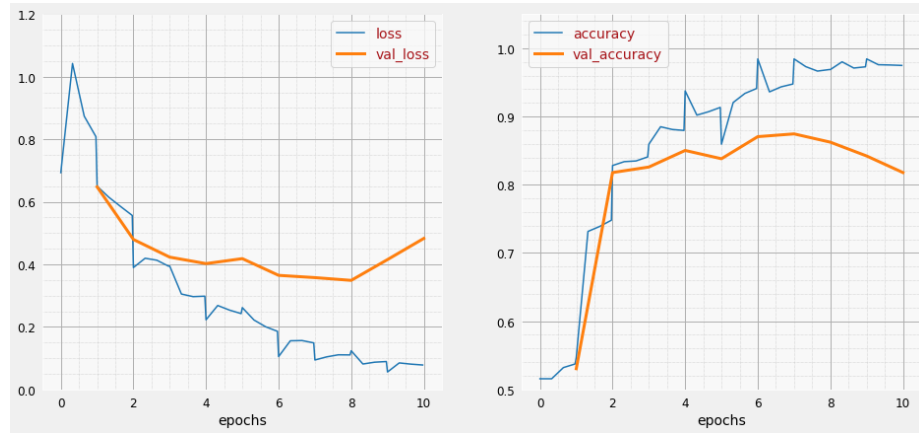


Figure 1: Performance of a CNN for binary classification

## 2.4 Clustering

The second method used in this project was clustering. First of all, since the dataset contains images with multiple dimensions, it was important to choose the correct feature extraction method.

<sup>1</sup>[https://www.tensorflow.org/api\\_docs/python/tf/keras/layers/Conv2D](https://www.tensorflow.org/api_docs/python/tf/keras/layers/Conv2D)

<sup>2</sup>[https://www.tensorflow.org/api\\_docs/python/tf/keras/layers/MaxPool2D](https://www.tensorflow.org/api_docs/python/tf/keras/layers/MaxPool2D)

<sup>3</sup>[https://www.tensorflow.org/api\\_docs/python/tf/keras/layers/Dense](https://www.tensorflow.org/api_docs/python/tf/keras/layers/Dense)

### 2.4.1 Feature extraction methods

- Principal component analysis(PCA)

It allows us to extract features and at the same time reduce the dimensions of our dataset to any number less than the current number of features, always preserving the maximum amount of information which is measured through the variance of data along the various dimensions. Another application of PCA is to compress the data, hence reducing the computational time. In this project, this technique was used to tackle both problems.

- VGG16

Another idea is to use a Convolutional Neural Network as a feature extraction method. VGG16 is a Convolutional Neural Network (CNN) model proposed by Karen Simonyan and Andrew Zisserman at the University of Oxford. VGG16 has two targets. The first is to detect objects within an image coming from 200 classes, which is called object localization. The second is to classify images, each labeled with one of 1000 categories, which is called image classification. This CNN was chosen to perform the feature extraction stage, so the last layer of the network was removed and the last dense layer was used as a feature vector. This means that the new final layer is a fully-connected layer with 4,096 output nodes. This vector of 4,096 numeric values is the feature vector used to cluster the images.

### 2.4.2 Clustering algorithm

After feature extraction we needed to apply the cluster algorithm. Clustering is an unsupervised learning technique (no labels needed) that has the goal of grouping objects in such a way that objects in the same group (*cluster*) are more similar to each other according to a certain distance measure, than objects in different groups.

There are different approaches and algorithms to perform clustering which can be divided into three sub-categories:

- Partition-based clustering: E.g.  $k$ -means
- Hierarchical clustering: E.g. agglomerative
- Density-based clustering: E.g. DBSCAN

#### **$k$ -means**

$k$ -means is a simple, relatively fast and easy scalable partition based algorithm, its objective is to minimize the average squared Euclidean distance of images from their cluster centers.  $k$ -means needs to know in advance the number of clusters, it's sensitive to outliers and can consider only linear boundaries, so it's not optimal for clusters with non linear structure. Another problem is that  $k$ -means has a random initialization so can generate different clusters in different runs.

### Agglomerative clustering

Agglomerative clustering, also called bottom-up clustering, treats each image as a singleton cluster and after that it merges clusters until a single cluster will contain all the images. It aggregates items starting with the most similar, creating a new cluster that substitutes the corresponding items in the dataset. Then the similarity between different clusters is evaluated. The implementation of this algorithm is straightforward; it always generates identical clusters in different executions but it's slow for large datasets. Its complexity is at least  $O(n^2)$ .

### DBSCAN

DBSCAN views clusters as areas of high density separated by areas of low density. It has good performance with arbitrary shapes clusters and it is robust to outliers; on the other hand it has some parameters difficult to determine (like the threshold under which two points are considered neighbors) and it doesn't behave very well if clusters are very different in terms of in-cluster densities.

In this project  $k$ -means was used because it is able to scale to a large dataset and to produce a result in less time. It's not a problem fixing a number of predetermined clusters because the number of clusters is 2 (images with Covid and images without Covid).

## 3 Experimental results

Different approaches were used to respond to the same queries.

The different approaches with different execution times and precision are summarized in figure 2.

	Computational time
Classification with CNN	0:00:53.347152
Clustering with PCA - KMeans	0:03:01.008863
Clustering with VGG - KMeans	0:02:39.943830
Image similarity - Euclidean	0:02:28.449627
Image similarity - Cosine	0:06:21.365513

Figure 2: Comparison of different execution times

In the figure [3] we can see the comparison between the confusion matrix

of the clustering with PCA and the confusion matrix of the clustering with VGG16.

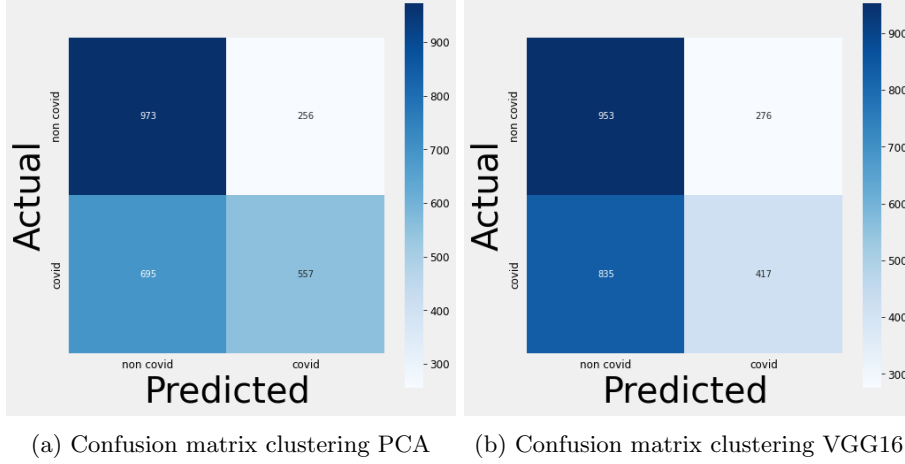


Figure 3: Comparison between PCA and VGG16

From a more generic point of view with clustering there was a great number of false negative, but the system was able to recognize rather well the true negative with a small number of false positive. There’s not a great difference between the two approaches of clustering with different feature extraction methods.

## 4 Conclusion and future work

In conclusion the classification task, using a Convolutional Neural Network, produces better results using less training time than all the other methods; this is due to the fact that this supervised learning technique makes use of labeled images, so it is easier to classify two different images. It is also quite useful to note that even the clustering technique, an unsupervised learning method, produces good results, correctly classifying more than 60% of the given images; it can be an alternative way to analyze images when their respective labels are unknown.

One of the first things to do in order to improve the quality of the final results is to enhance the image preprocessing; it is possible to apply data augmentation in order to find the correct pattern that properly classifies an image.

Another option is to improve the explainability of the process, in particular how a method achieves its corresponding result. For example some explainable AI tools can be used, a set of processes and methods that allows human users to comprehend and trust the results and output created by the machine learning algorithms. Explainability helps ensure that the system is working as expected and in this case makes it easier to understand why a certain image is chosen instead of another.

## References

- [1] Kumar, A., Kim, J., Cai, W., Fulham, M., & Feng, D. (2013). Content-based medical image retrieval: a survey of applications to multidimensional and multimodality data. *Journal of digital imaging*, 26(6), 1025-1039.
- [2] Yan, K., Wang, X., Lu, L., & Summers, R. M. (2018). DeepLesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning. *Journal of medical imaging*, 5(3), 036501
- [3] Ibtihaal M. Hameed, Sadiq H. Abdulhussain & Basheera M. Mahmmod — D T Pham (Reviewing editor) (2021) Content-based image retrieval: A review of recent trends, *Cogent Engineering*, 8:1, DOI: 10.1080/23311916.2021.1927469
- [4] Ling Ma, Xiabi Liu, Yan Gao, Yanfeng Zhao, Xinming Zhao, Chunwu Zhou, A new method of content based medical image retrieval and its applications to CT imaging sign retrieval, *Journal of Biomedical Informatics*, Volume 66, 2017, Pages 148-158, ISSN 1532-0464
- [5] Afshan Latif, Aqsa Rasheed, Umer Sajid, Jameel Ahmed, Nouman Ali, Naeem Iqbal Ratyal, Bushra Zafar, Saadat Hanif Dar, Muhammad Sajid, Tehmina Khalil, "Content-Based Image Retrieval and Feature Extraction: A Comprehensive Review", *Mathematical Problems in Engineering*, vol. 2019, Article ID 9658350, 21 pages, 2019.
- [6] Simonyan, K., & Zisserman, A. (2015). Very Deep Convolutional Networks for Large-Scale Image Recognition. *CoRR*, abs/1409.1556.