

covid-19 search engine

October 5, 2022

Marco Ghezzi Matteo Limoncini

1 Introduction

This is the report of the project for the Information Retrieval course. This project aims to create a search engine being able to respond to some queries.

This project tries to solve the problem of searching and making queries over a large dataset of images. 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

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.

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

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 developed different Convolutional Neural Networks.

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 our 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 ct scans taken from real patients in São Paulo, Brazil: every image is paired with a label to specify if the person associated to 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

In order to read and process the images, it's necessary to find the optimal method in term of space and time efficiency.

Since every image is associated with a label, they've been divided in two directories: the first containing all the images labeled with covid and the second containing all the images labeled with non covid. For the classification task the labeled images need to be divided in three different datasets: train, validation and test sets. We used a split ratio of 0.8, 0.1 and 0.1 respectively. Starting from the images in these directories, we created three different datasets: one for training data, one for validation data and the last one for test data. While creating these datasets, we specified the size of the images 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 we normalized all pixel values between 0 and 1 since the computation of high numeric values may become more complex.

In order to create a system that is able to answer the first query with the most precision and less time consuming (identifying the ct scan images with covid) we evaluated different approaches. Initially we focused on Machine Learning techniques.

2.3 Classification

Since we're dealing with an image binary classification problem, we decided to choose a Convolutional Neural Network. A Convolutional Neural Network (CNN) is an artificial neural network suited for images, with a sequence of layers, and every layer transforms one volume of activations to another through

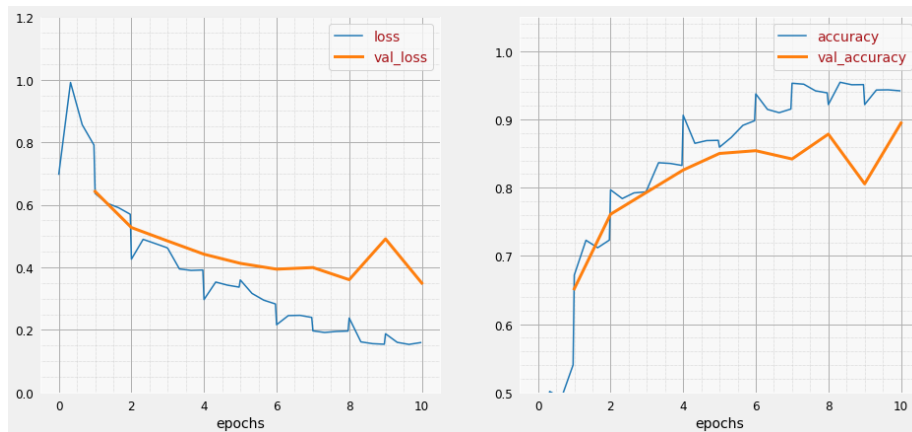


Figure 1: Performance of CNN for binary classification

a differentiable function. We used three main types of layers to build our neural network: Convolutional Layer ¹, Pooling Layer ² and Fully-Connected Layer.

Using this CNN, we're able to obtain interesting results: this network is able to correctly classify a new image with an accuracy of 89% with less than a minute of training time over 1986 images.

2.4 Clustering

The second method we took into consideration was clustering. First of all, since we're dealing with images thus with multiple dimensions, we need to choose the correct feature extraction method.

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 current number of features, always preserving the maximum amount of information. Another application of PCA is to compress the data and hence reduce the computational time. In this project, we will use this technique to tackle both problems.

- VGG16

Another idea is to use a Convolutional Neural Networks as feature extraction methods. 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

¹https://www.tensorflow.org/api_docs/python/tf/keras/layers/Conv2D

²https://www.tensorflow.org/api_docs/python/tf/keras/layers/MaxPool2D

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. We decided to use this CNN to perform the feature extraction stage, so we removed the last layer of the network and we use the last dense layer 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 that we will use to cluster the images. Since this layer has over 4000 dimensions we decided to reduce its dimensions using PCA. PCA allow us to reduce the dimensionality keeping as much information as possible, after this step we have a smaller feature vector thus we can apply a classic cluster algorithm.

2.4.2 Clustering algorithm

After feature extraction we need 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, according to a certain distance measure, to each other than objects in different groups.

There are different approaches and algorithms to perform clustering tasks 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, is 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 run.

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, each cluster that is defined this way, substitutes the corresponding items in the dataset. Then is evaluated the similarity between different clusters in order to define their similarity. It is an easy to implement algorithm and always generate same clusters but it's slow for large datasets, 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

	Computational time
Classification with CNN	0:00:51.235565
Clustering with PCA - KMeans	0:03:40.649903
Clustering with VGG - KMeans	0:02:44.163987
Image similarity - Euclidean	0:02:18.823460
Image similarity - Cosine	0:05:50.546088

Figure 2: Comparison of different execution times

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 term of in-cluster densities.

We decided to use K-Means because in our case it's important to use an algorithm able to scale to large dataset and able to produce a result in not so much time and it's a not problem fixing a number of clusters predetermined, number of clusters is 2 (images with covid and images without covid).

3 Conclusion and future work

In this work was used different approaches to solve the same queries. We've summarised the different approaches with different execution times and precision as we can see in figure 2.

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] Ibtihail 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, <https://doi.org/10.1016/j.jbi.2017.01.002>. (<https://www.sciencedirect.com/science/article/pii/S1532046417300023>)
- [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. <https://doi.org/10.1155/2019/9658350>
- [6] Very Deep Convolutional Networks for Large-Scale Image Recognition Karen Simonyan, Andrew Zisserman