



SAPIENZA
UNIVERSITÀ DI ROMA

FACOLTÀ DI
INGEGNERIA DELL'INFORMAZIONE, INFORMATICA E STATISTICA

DIPARTIMENTO
DI INGEGNERIA INFORMATICA, AUTOMATICA E GESTIONALE

Neural Networks course
Final project

A.Y. 2020/21

Report of C3AE Project

Professors:

Aurelio Uncini
Simone Scardapane

Students:

Francesco Petri 1797147
Giovanni Pecorelli 1799865

Contents

1	Introduction	2
1.1	Related works	2
1.2	Report organization section	2
2	Datasets	3
2.1	Preprocessing	3
2.2	Data generation	4
2.2.1	Training data augmentation	4
3	Theory Stuff	5
3.1	Spiegazione modello C3AE	5
3.2	Spiegazione problema	5
3.3	Implementation	5
4	Experiments	6
4.1	??first experiment??	6
4.2	??other full experiments??	6
4.3	Ablation Study	7
5	Project structure	8
6	Conclusions	9

Chapter 1

Introduction

In the past decade ((years?)) soft biometrics has emerged out to be a new area of interest for the researchers due to its growing real-world applications. This includes a classic learning problem in computer vision: the estimation of demographic traits, such as the age. Researchers are trying to develop models which can accurately estimate the age or the age group of a person using different biometric traits. Currently, neural networks give the best classification results for age estimation using human faces. Many CNNs (convolutional neural networks) such as AlexNet, VggNet, GoogLeNet and ResNet ((citations needed?)) are able to accomplish this task with promising performance.

However, to obtain more precise accuracy these networks have grown deeper and larger. This trend has resulted in increasingly higher computational costs in either training or deploying. In particular, deploying the previously mentioned models on mobile phones, cars and robots is next to impossible due to the model size and computational cost.

Recently other models have been proposed with the aim to reduce the number of parameters, thus yielding lightweight models without weakening their efficiency. In this work, we investigate the limits of compact models for small-scale images and focus on one the most compact models for age classification, implementing it in practice to evaluate its performance.

1.1 Related works

The following report presents the development of the final project for the Neural Networks course at Università degli studi di Roma "La Sapienza", A.Y. 2020/21.

Our work is based on the study made by ((citation)). In the paper ((citation)) they propose a **Compact** basic model, **C**ascaded training and multi-scale **C**ontext, aiming to tackle small-scale image **A**ge **E**stimation. The model is called **C3AE**.

The proposed model is able to achieve a state-of-the-art performance compared with alternative compact models and even outperforms many bulky models. With an extremely compact model of 0.25 MB for the full model, which is possibly the smallest model that has been obtained so far on the facial recognition, C3AE is suitable to be deployed even on low-end mobiles and embedded platforms. A discussion on which techniques have been used to attain the desired results are discussed in a later chapter. [1]

1.2 Report organization section

((Report organization section?))

Chapter 2

Datasets

Our C3AE implementation was trained and tested on the following datasets:

- **Wiki**: a large dataset containing 62,328 labelled images¹ collected from Wikipedia [2]. Despite the size, it lacks samples of very young or very old people, and it is quite noisy. The cropped and aligned version of the dataset was used in order to ensure each picture has a single face in it. This dataset was used as a pre-training set for ((one experiment)) and as the training set for the whole ablation study.
- **UTKFace**: a dataset containing over 20,000 labelled images². It covers a wider range of ages compared to Wiki, therefore this dataset was used as the main training set for ((that same experiment)) and as a fine-tuning set after the pre-training with Wiki.
- **FG-NET**: a dataset containing 1000 labelled images³ [3]. It is significantly higher-quality and better-curated than the other datasets listed here, but it is extremely small as well. This dataset was used as the test set for all experiments.

2.1 Preprocessing

Each dataset initially came in the form of a set of image files each with the corresponding ground-truth age encoded in the file name in some way. In order to make the datasets usable by our model, the following preprocessing procedure was applied to each of them.

First of all, the age information is extracted from the file name of each image through a dataset-dependent regular expression.

Second, the face in each image is detected with the MTCNN face-detection network⁴ [4], and the resulting information is used as a base to position the three bounding boxes to be used later to generate the different crops that C3AE uses as multi-scale context (see chapter 3: Theory Stuff).

Then, faces with associated age outside the $[0, 120]$ range are filtered out of the dataset. Images where less or more than one face is detected are also discarded.

Finally, images, age labels and bounding boxes are all organized into a **pandas** table which is saved to disk in the **pickle** format, in order to be loaded later by the other parts of the code.

¹Collected and available at <https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/>

²Collected and available at <https://susanqq.github.io/UTKFace/>

³Collected and available at https://yanweifu.github.io/FG_NET_data/

⁴Imported as external code from <https://github.com/ipazc/mtcnn>

2.2 Data generation

Training, validation and test data are generated in real time during the respective phase from the preprocessed dataset.

This requires a small amount of additional processing, which normally consists of applying a reflect-type padding to handle bounding boxes that are partially outside the border of the image, then generating three cropped images, one per bounding box, and finally resizing each of those crops to 64×64 pixels, the accepted input size of the C3AE model.

2.2.1 Training data augmentation

Training data undergo additional transformations during this generation process in order to augment the dataset and make the trained model more robust:

- Random erasing: before adding the reflect padding, an arbitrary portion of the image is deleted and replaced with random noise [5].
- Random shift: Before cropping, each bounding box is independently shifted by a random amount. If any part of a box would move past the border of the padded image, it stops at the border instead.
- Random contrast, brightness and color temperature change⁵
- Random rotation
- Vertical flipping

Each operation has a random probability to be applied and is independent of the others, and the intensity of each transformation (except flipping) is also random, but limited.

⁵The temperature change code was imported as external from <https://www.askaswiss.com/2016/02/how-to-manipulate-color-temperature-opencv-python.html> and adapted to randomize the intensity of the change.

Chapter 3

Theory Stuff

Morbi maximus dui vitae lorem pretium cursus. Cras risus nisl, viverra a libero ut, vehicula iaculis dui. Suspendisse potenti. Fusce mi odio, maximus ut arcu blandit, condimentum bibendum lorem. Sed vel tempor est. Nulla a orci lobortis, fermentum nisi non, vulputate sapien. Etiam efficitur placerat enim, a malesuada nisl cursus quis. Cras et elit ut lorem imperdiet posuere sed sit amet dolor. Aliquam id velit ac lorem consequat ullamcorper.

3.1 Spiegazione modello C3AE

È un modello molto bello

3.2 Spiegazione problema

qqqqqqqq

3.3 Implementation

aaaaaa

Chapter 4

Experiments

In this chapter we detail the experiments performed on the model described in chapter 3: Theory Stuff and discuss the respective results.

4.1 ??first experiment??

An initial experiment was performed by training the C3AE model for 400 epochs on the Wiki dataset as a means to test the `tensorflow` environment and the functioning of the implemented C3AE model within it.

The evolution of training and validation loss is shown in Figure 4.1. The experiment ran to completion with no runtime errors, and it can be observed that the model was able to decrease its loss, and that such loss reached an asymptote around epoch 50.

For this reason, we concluded that a training time of 400 epochs is too much for this model and decided to set the epoch limit to 100 for all following experiments, assuming that they would have a similar evolution to this one and therefore all significant improvement would happen much before the 100th epoch, in order to reduce the experiments' computation time.

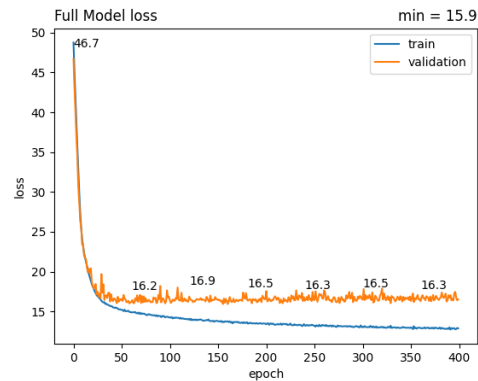


Figure 4.1: Initial experiment loss (400 epochs on Wiki)

4.2 ??other full experiments??

????????????????????

(One on Wiki, one on UTK, one on Wiki+UTK, all tested w/ FGNET)

????????????????

4.3 Ablation Study

A separate set of experiments was performed to study the impact on performance of the following components of the model and the training process: the context module and the cascade module of the C3AE model, and the training data augmentation.

We trained the following variants of the full C3AE model:

- *Full model*: the standard model with no changes, to serve as a benchmark against the other variants.
- *No augmentation*: the data augmentation transformations on the training data are disabled.
- *No context*: the context module of C3AE is excluded. Therefore, only one crop, the outermost one, is given as input to the model, and obviously there is no concatenation phase.
- *No cascade*: the cascade module of C3AE is excluded. Consequently, the intermediate layer between the concatenated feature vector and the age output loses its meaning of two-point representation of age and becomes a plain hidden layer. This also means that this variant does not compute any KL Divergence and outputs only the final age estimation; the loss depends only on the age MAE as well.
- *No context and no cascade*: both modules are disabled.

Each variant was trained for 100 epochs on the Wiki dataset. ((Use of test set needs some debugging))

((plot and discuss results))

Chapter 5

Project structure

Probabilmente non ci serve, ma non voglio scordarmi che esiste

Chapter 6

Conclusions

Maecenas quis magna id lorem auctor consequat ut consectetur justo. Suspendisse dui lectus, cursus id nisl sit amet, dapibus sodales ipsum. Duis urna urna, sollicitudin non molestie a, cursus pellentesque ante. Donec laoreet dolor ac ligula mollis lacinia. Aenean at tristique leo. Vestibulum in ultricies est, ac ullamcorper lectus. Phasellus placerat nulla ante, vitae pharetra risus scelerisque ut. Sed consequat sem rhoncus, congue erat in, dignissim nisi. Aenean accumsan quis diam ornare ornare. Vivamus quis vehicula urna, ut rhoncus ligula.

Bibliography

- [1] Pippo, Pluto e Paperino. Guida di Paperopoli. ICLR, 1050:11, 2021
- [2] Rasmus Rothe and Radu Timofte and Luc Van Gool. *Deep expectation of real and apparent age from a single image without facial landmarks*. IJCV 126(2-4):144–157, 2016.
- [3] Yun Fu, Guodong Guo, and Thomas S Huang. *Age synthesis and estimation via faces: A survey*. IEEE Trans. on Pattern Analysis and Machine Intelligence 32(11):1955–1976, 2010.
- [4] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, Yu Qiao. *Joint Face Detection and Alignment using Multi-task Cascaded Convolutional Networks*. IEEE Signal Processing Letters 23(10):1499–1503, 2016.
- [5] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. *Random erasing data augmentation*. arXiv preprint arXiv:1708.04896, 2017.