



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
2	Datasets	3
2.1	Preprocessing	3
2.2	Data generation	4
2.2.1	Training data augmentation	4
3	The C3AE Model	5
3.1	Model details	5
3.2	Implementation	6
4	Experiments	8
4.1	First experiment	8
4.2	Performance evaluation over multiple datasets	8
4.2.1	Wiki	9
4.2.2	UTK	9
4.2.3	Wiki + UTK	9
4.3	Ablation Study	10
4.3.1	No augmentation	10
4.3.2	No context	10
4.3.3	No cascade	10
4.3.4	Full ablation	11
5	Conclusions	12

Chapter 1

Introduction

In the past decade 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 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 want to 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.

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.

1.1 Related works

Our work is based on the study made by Chao Zhang, Shuaicheng Liu, Xun Xu and Ce Zhu in the paper "*C3AE: Exploring the Limits of Compact Model for Age Estimation*" [1] in which they propose a **C**ompact 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 size of 0.25 MB for the full model, which is possibly the smallest that has been obtained so far on the facial recognition, C3AE is suitable to be deployed even on low-end mobiles and embedded platforms.

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: The C3AE Model).

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

The C3AE Model

3.1 Model details

The C3AE plain model takes as input a source image and tries to output a plausible estimation for the age of the portrayed person.

Since the objective is to obtain a lightweight model the total number of parameters needs to be as low as possible, as long as it doesn't affect the overall performance. For this reason the input image size is limited ($64 \times 64 \times 3$) and the other channels sizes are also small.

Standard convolutional layers are adequate for the trade-off between accuracy and compactness (as opposed to the separable convolution block used in the bigger models like MobileNets and ShuffleNets), followed by batch normalization, Relu and average pooling (BRA).

The model is composed of five of these standard convolutions and two fully connected layers as shown in Table 3.1.

Layer	Kernel	Stride	Output size	Parameters
Image	-	1	$64 \times 64 \times 3$	-
Conv. 1	$3 \times 3 \times 32$	1	$62 \times 62 \times 32$	896
BRA	-	1	$31 \times 31 \times 32$	128
Conv. 2	$3 \times 3 \times 32$	1	$29 \times 29 \times 32$	9248
BRA	-	1	$14 \times 14 \times 32$	128
Conv. 3	$3 \times 3 \times 32$	1	$12 \times 12 \times 32$	9248
BRA	-	1	$6 \times 6 \times 32$	128
Conv. 4	$3 \times 3 \times 32$	1	$4 \times 4 \times 32$	9248
BN + ReLu	-	1	$4 \times 4 \times 32$	128
Conv5	$1 \times 1 \times 32$	1	$4 \times 4 \times 32$	1056
Feat.	$1 \times 1 \times 32$	1	12	6156
Predict.	$1 \times 1 \times 1$	1	1	13

Table 3.1: Architecture of the model

To estimate people's age C3AE considers two objectives simultaneously: the first one minimizes the Kullback-Leibler loss between distributions, and the second one optimizes the squared loss between discrete ages.

In the following sections are mentioned this and other techniques used in C3AE.

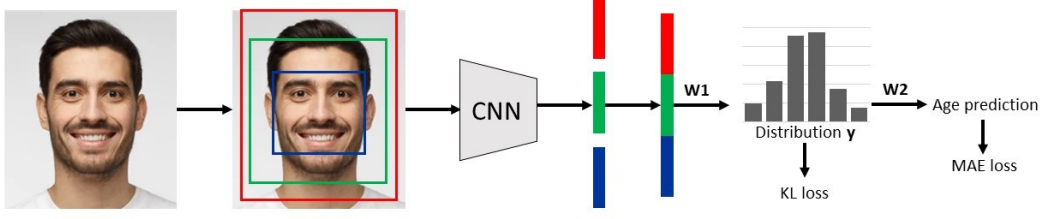


Figure 3.1: Overview of the model age estimation process

Context-based Regression

The resolution and the size of small-scale image is limited, so the idea is to exploit facial information at different granularity levels. After the face recognition task we identify three different crops of different sizes of the subject’s face, like in Figure 3.1. Each cropped image has a special view on the face. The smaller image contains rich local information; in return the bigger one may contain global and scene information. The three crops are then fed into the shared CNN network, and finally the bottlenecks of the three-scale facial images are aggregated by concatenation.

Two-point Age Representation

With C3AE we don’t predict directly the age as a number. Instead, this model uses a two-point age representation through a distribution over two discrete adjacent bins.

For example let’s consider the corresponding representation of an age of 22 with 10 bins. In this case the set of bins is [10, 20, 30, 40, 50, 60, 70, 80, 90, 100] and the corresponding vector representation of the age is [0, 0.8, 0.2, 0, 0, 0, 0, 0, 0, 0].

Cascade Training

From the above section, the age value can be represented as a distribution vector. The mapping from this vector to age value is decomposed into two steps, where we define two different losses for the two cascade tasks.

The first one (L_{kl}) measures the discrepancy between ground-truth label and predicted age distribution. We adopt KL-Divergence as the measurement.

The second loss (L_{reg}) controls the prediction of the final age and is implemented as an L1 distance or mean absolute error (MAE loss).

In the training process the two loss functions are considered in the cascade style as shown in Figure 3.1 but they are still trained jointly, and the total loss is given as $L_{total} = \alpha L_{kl} + L_{reg}$, where α is the hyperparameter to balance the two losses.

3.2 Implementation

Our implementation of the C3AE model can be found in [this GitHub repository](#) ((change href color?)). The code for this part of the project is based on the Tensorflow library for Python and was written in Visual Studio Code text editor, collaborating remotely using the Live Share extension feature, which let us write the code together.

The model

The code of the model implementation can be found in the directory `master/model.py`. We followed as strictly as possible the structure indicated in the reference paper, both in the number, composition and ordering of the layers and also in the choice of parameters such as kernel dimensions, bottleneck dimension and squeeze factors. The Keras package integrated in Tensorflow provided all the basic functions needed to implement the building blocks of the C3AE plain model, shared by all three of the image cuts.

In the last part of the file we organized the code so that we could be able to replicate the ablation studies conducted in the original C3AE paper. Ablation studies are procedure where certain parts of the network are removed, in order to gain a better understanding of the network's behaviour. In our case, the two ablations affect the presence or the absence of the cascade module and the context module. The `model.py` file therefore gives us the possibility to analyze not only the full context model, but also compare it with three slightly simpler ones: the model without cascade, without context and without both.

The presence of these alternative versions required some changes in the code:

- *Context ablation*: produces only one simple input, instead of three concatenated ones, before computing the losses in the training phase. This means we will train using only one cut of the photos, the outer one.
- *Cascade ablation*: in this case we lose the concept of two-point age representation. We replace the last dense layer W1 with a hidden layer with no particular meaning and then only considered one output, the age prediction, and ignore the second output, the age distribution. It follows that in the training phase we will only evaluate the MAE loss in this case, since the KL divergence loss will no longer be present.
- *Full ablation*: both of the above modifications apply.

The results of these ablation studies are discussed later in section 4.3: Ablation Study.

The training

The code relevant to the training phase can be found in `master/training.py`. Again, the choice of the parameters such as learning rate, batch sizes and optimizers follows the values indicated in the source paper when possible, with a couple of exceptions.

The dataset split ratio into training and validation was not specified for the datasets we used, so we chose a value of 88% - 12%.

In the original paper for a comparison with the state-of-the-art methods, each model has been trained for 600 epochs. In our implementation this number had to be cut down to 100, the reason being time constraints and hardware availability. All training was conducted locally on a machine with a Nvidia GTX 1660 SUPER GPU, and even with 100 epochs each training phase lasted for 10 hours on average. A discussion on how this cutback may not have affected the final results is present in the next chapter.

In this section of the code we also implemented the possibility to take an already trained model and continue its training on another dataset, instead of starting over from zero. An experiment taking advantage of this feature is described in ??: ??.

Chapter 4

Experiments

In this chapter we detail the experiments performed on the model described in chapter 3: The C3AE Model 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.

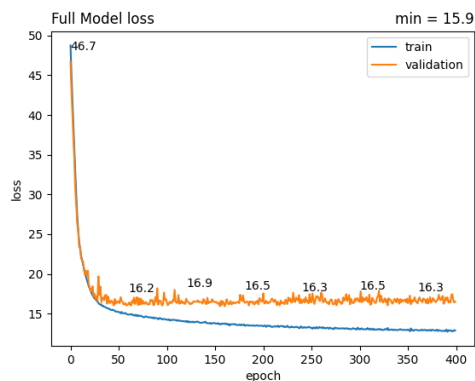


Figure 4.1: Initial experiment loss (400 epochs on Wiki)

4.2 Performance evaluation over multiple datasets

All the following experiments have been performed on datasets processed with the augmentation techniques described in subsection 2.2.1: Training data augmentation unless stated otherwise.

4.2.1 Wiki

The first and main dataset used in our experiments has been Wiki. As shown in the loss graph in Figure 4.2 the learning curve is correct and the best MAE value obtained in this experiment is 6.79 years. When compared to the result in the original paper, which claims to have reached a MAE of 6.44 on the same dataset, we could say that we have pretty much achieved a state-of-the-art performance.

However, this result only takes in consideration Wiki images for both training and validation. So we tested the output of this experiment with the FGNET test set, and found a much higher MAE of 18.1 years. The main reason for this is that the Wiki dataset is mostly composed by images of adults and elders and features almost no children, while FGNET has a much lower age on average and comprises even newborns with a declared age of 0 years. The result is that the model in this experiment always overshoots the age of the subjects and thus we obtain a high error.

4.2.2 UTK

In this experiment the starting conditions and parameters are the same as the previous one, but the dataset this time is UTK. The number of images is one third of Wiki, but it covers better the whole range of ages from 0 to over 100. As a matter of fact the validation MAE (computed on the UTK set itself) is slightly higher than before at 8.67, but the test MAE on FGNET is almost halved, at 9.79 years.

The conclusion is that UTK is a better dataset for our experiments, but we can still use Wiki, as seen on the next experiment.

4.2.3 Wiki + UTK

This third experiment combines the previous two. We started by pre-training a model from scratch with the Wiki dataset, and then we took the output of this process and further trained it for another 100 epochs on the UTK dataset. In this way we hoped to combine the scale of the first dataset with the completeness of the second to obtain a model that outperforms the previous two. And indeed, with a final validation MAE of 8.23 and a test MAE of 8.64 years on FGNET this has proven to be our best result yet.

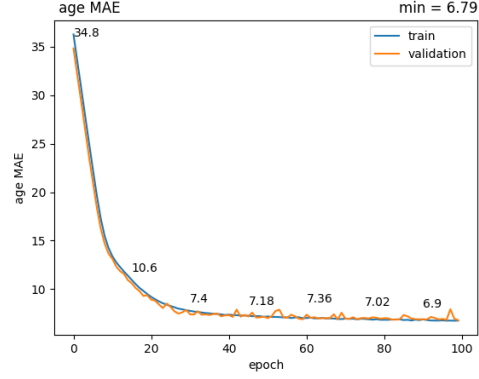


Figure 4.2: MAE on Wiki dataset (100 epochs)

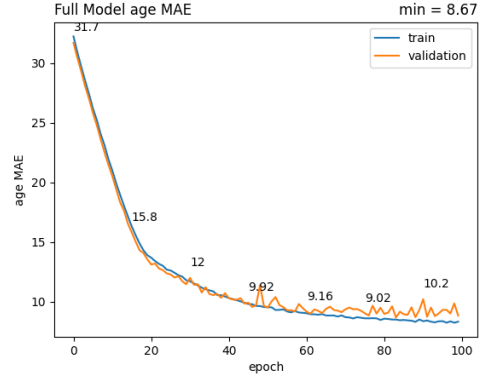


Figure 4.3: MAE on UTK dataset (100 epochs)

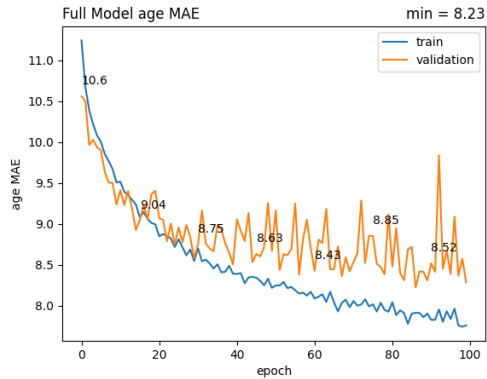


Figure 4.4: MAE on Wiki+UTK (100 epochs)

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 is given as input to the model.
- *No cascade*: the cascade module of C3AE is excluded. Consequently, this variant does not compute any KL Divergence and outputs only the final age estimation.
- *No context and no cascade*: both modules are disabled.

Each variant was trained for 100 epochs on the Wiki dataset, so every ablation experiment results must be compared to the baseline results of the one described in subsection 4.2.1: Wiki.

4.3.1 No augmentation

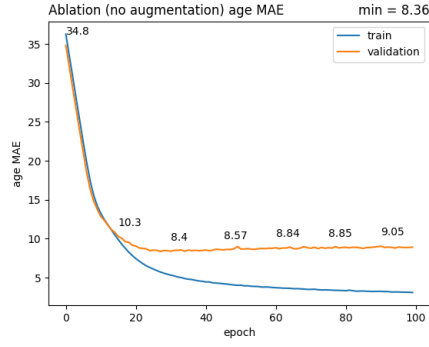
In this first ablation study we removed the whole augmentation process described in subsection 2.2.1: Training data augmentation. The validation MAE increases by 23% compared to the the Wiki full model experiments, and the test MAE by not less that 30%. We conclude that the augmentation phase implemented by us has been highly beneficial to the full model performance.

4.3.2 No context

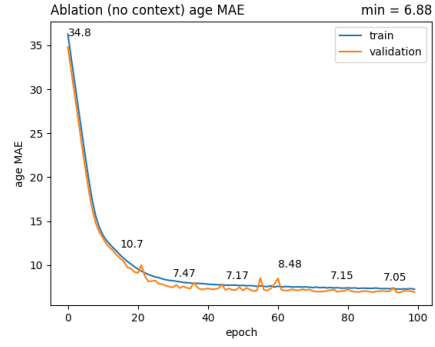
After the removal of the context module, the model retained its performance if we consider only the validation results ($MAE = 6.88$). However when examining the test results on FGNET they are rather worse then before, with the test MAE going from 18.11 to 20.73. This performance difference proves us that considering image cuts with three different levels of details instead of only one and combining their results is useful to the full model performance.

4.3.3 No cascade

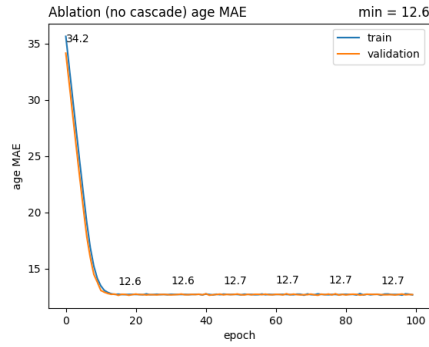
When we take out the cascade module we the only metric that contributes to the total loss is the MAE, since we lose the output with the age distribution and also the KLD loss. The direct exstimation of an age from the image in input is quite a harder task, judging by the results of this experiment. The validation MAE is by far the highest found so far, at 12.6 years, and the test MAE on FGNET isn't great either, reaching 18.67. Once again, we can interpret this performance drop as a proof that also the cascade module is essential to obtain a good result in the full model.



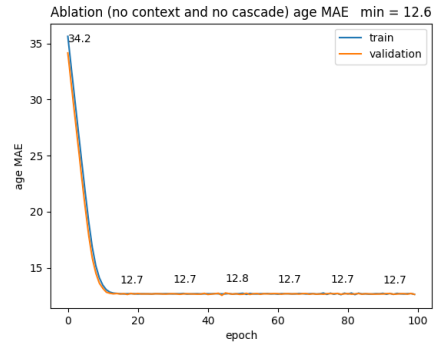
(a) MAE without augmentation



(b) MAE without context



(c) MAE without cascade



(d) MAE without either

Figure 4.5: Ablation Study MAE

4.3.4 Full ablation

The absence of both the previous modules results unsurprisingly in a performance comparable to the worst of the two. The validation and test MAE values are in fact aligned to the MAEs returned by the "no cascade" model.

Chapter 5

Conclusions

Section	Model	Dataset	Augm.	Epochs	Valid. MAE	Test MAE
First	Full	Wiki	Some	400	6.82	22.72
Wiki	Full	Wiki	✓	100	6.79	18.11
UTK	Full	UTK	✓	100	8.67	9.79
Wiki+UTK	Full	Wiki+UTK	✓	100+100	8.23	8.64
No augm.	Full	Wiki	×	100	8.36	24.41
No context	No context	Wiki	✓	100	6.88	20.73
No cascade	No cascade	Wiki	✓	100	12.60	18.67
Full ablation	Full ablation	Wiki	✓	100	12.60	18.69

Bibliography

- [1] Chao Zhang, Shuaicheng Liu, Xun Xu, Ce Zhu. *C3AE: Exploring the Limits of Compact Model for Age Estimation*. CoRR, 2019.
URL: <https://arxiv.org/pdf/1904.05059v2.pdf>
- [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.