# One-Shot Learning: Face Recognition with Siamese Neural Network

Aveen

March 2020

Abstract Although recent research in the ground of face recognition, achieving face verification and recognition accurately at extent presents deliberate objections to current techniques. Here we present a structure, called FaceNet, that precisely grasp an outline from face images to a compress Euclidean space where distances precisely coincide to a measure of face closeness. When this area has been formed, functions as face recognition, verification and clustering can be conveniently carry out adopting typical techniques with FaceNet embedding's as feature vectors. This model uses a deep convolutional network trained to precisely utilize the embedding itself. To train, we have used triplets of roughly aligned matching / non-matching face patches generated using triplet mining method. The benefit of this approach is much greater representational efficiency: we get face recognition performance using only 128-bytes per face. On the VGG2 Face dataset, this model achieves an accuracy of 99.63

## 1 Introduction

One-shot learning is a classification task where one, or a few, examples are used to classify many new examples in the future. This define tasks examine in the field of face recognition, such as face identification and face verification, where people can be differentiated accurately with different facial expressions, lighting conditions, accessories, and hairstyles given one or a few template images. Here we are using it for the comparison of two people's identities in images. Current face recognition models access the dilemma of one-shot learning along face recognition by training a low-dimensional feature representation, called a face embedding, which can be determined for faces efficiently and correlated for verification and identification. Embedding were used to learn for one-shot learning complications using a Siamese network. The training of Siamese networks with comparative loss functions appear in better performance, next leading to the triplet loss function used in the FaceNet system by Google that achieved then state-of-the-art results on standard face recognition tasks. Here, you will discover the challenge of one-shot learning in face recognition and how comparative and triplet loss functions can be used to learn high-quality face embedding for One-shot learning.

## 2    Face Recognition

The growth in the implementation of facial recognition organizations has been remarkable in current years; yet, it came below a lot of examination recently. There are numerous Artificial Intelligence enthusiasts who ponder that the practice of any type of facial recognition system should be appropriately structured in demand to avoid despicable undertakings. Face recognition is an all-purpose assignment of classifying and verifying individuals from pictures of their appearance. In face recognition structures, we need to be capable to recognize an individual's identity by just serving one picture of that person's face to the structure. And in case of failure to recognize the person, it shows the results accordingly. To resolve this issue, we cannot use only a convolutional neural network for two reasons: 1) CNN doesn't work on a minor training set. 2) It is not suitable to retrain the model every time we add a picture of a new person to the system. However, we can use Siamese neural network for face recognition.

## 3    One-Shot Learning

Characteristically, classification contains fitting a model given various samples of each class, then using the fit model to make calculations on many examples of each class. We might have thousands of quantities of plants from three altered classes. A model can be fit on these instances, simplifying from the commonalities among the capacities for a given class and contrasting alterations in the capacities across class. The consequence with any luck is a healthy model that, given a new set of capacities in the future, can precisely predict the plant class. One-shot learning is a classification assignment where one instance (or a very small number of instances) is given for each class that is used to make a model that in turn must make predictions about numerous unidentified samples in the future.

"In the case of one-shot learning, a single exemplar of an object class is presented to the algorithm." — Knowledge transfer in learning to recognize visual objects classes, 2006.

This is a comparatively relaxed unruly for human being. An individual may see a Ferrari sports car once and in the future be capable to identify Ferraris in new circumstances, on the road, in movies, in books, and with different lighting and colors.

"Humans learn new concepts with very little supervision – e.g. a child can generalize the concept of "giraffe" from a single picture in a book – yet our best deep learning systems need hundreds or thousands of examples." — Matching Networks for One Shot Learning, 2017.

Explicitly in the circumstance of face identification, a model or system may only have one or a few examples of a given individual's face and must appropriately recognize the individual from new pictures with deviations to appearance, hairstyle, illumination, accessories, and more. In the case of face verification, a model or system may only have one example of an individual's face on best and

must appropriately authenticate new photos of that person, possibly each day. As such, face recognition is a common example of one-shot learning.
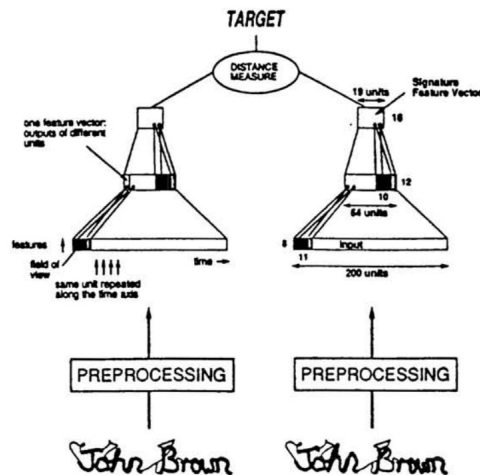
# 4    Siamese Neural Network

A system that has been popularized given its use for one-shot learning is the Siamese network. Siamese network is a structural design with two equivalent neural networks, each captivating an altered input, and whose outputs are combined to provide some prediction.

It is a system intended for authentication tasks, first projected for signature verification by Jane Bromley et al. in the 2005 paper titled "Signature Verification using a Siamese Time Delay Neural Network."

"The algorithm is based on a novel, artificial neural network, called a "Siamese" neural network. This network consists of two identical sub-networks joined at their outputs." — Signature Verification using a "Siamese" Time Delay Neural Network, 2005.

Two equal systems are used, one taking the known signature for an individual, and another taking a candidate signature. The results of both systems are collective and counted to specify whether the candidate signature is real or a counterfeit.

"Verification consists of comparing an extracted feature vector with a stored feature vector for the signer. Signatures closer to this stored representation than a chosen threshold are accepted, all other signatures are rejected as forgeries." — Signature Verification using a "Siamese" Time Delay Neural Network, 2005.



Siamese systems were used more lately, where deep convolutional neural networks were used in parallel image inputs in a 2015 paper by Gregory Koch, et al. titled "Siamese Neural Networks for One-Shot Image Recognition."

The deep CNNs are first trained to distinguish among instances of each class. The impression is to have the models study feature vectors that are active at extracting intellectual features from the input images.
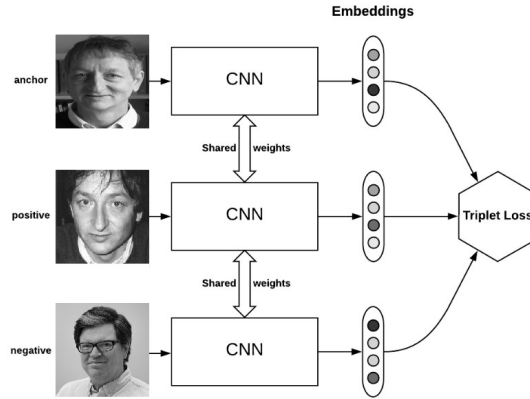


**Verification tasks (training)**

The models are then re-purposed for verification to forecast whether new samples match a pattern for each class. Precisely, each system produces a feature vector for an input picture, which are then matched using the L1 distance and a sigmoid activation. The model was applied to benchmark handwritten character datasets used in computer vision.

# 5 FaceNet Model

The intimidations like data seepage, confidentiality destruction, initiating from careless usage of facial recognition systems are actual and appropriate actions should be taken to evade them, but even after all the recent disapproval, we have to acknowledge that it is still a pretty valuable application which can be broadly used to make publics' lives better. FaceNet delivers an exclusive strategy for performing tasks like face recognition, verification and clustering. It uses deep convolutional networks lengthways with triplet loss to accomplish state of the art precision. FaceNet delivers a combined embedding for face recognition, verification and clustering tasks. It plots each aspect picture into a Euclidean space such that the detachments in that space resemble to face similarity . An image of person A will be positioned nearer to all the other images of person A as linked to images of any other person present in the dataset.

Alteration between FaceNet and other methods is that it acquires the mapping from the images and creates embeddings rather than using any block layer for recognition or authentication tasks. Once the embeddings are shaped all the other tasks like verification, recognition can be done using normal practices of that specific domain, using these anew produced embeddings as the feature vector. For instance we can use k-NN for face recognition by using embeddings as the feature vector and likewise we can use any clustering technique for clustering the expressions together and for verification we just want to describe a threshold value.

FaceNet uses deep convolutional neural network (CNN). The network is trained such that the squared L2 distance between the embeddings resemble to face similarity. The imageries used for training are ascended, altered and are closely cropped around the face area. Additional significant feature of FaceNet is its loss function . It uses triplet loss function. In order to compute the triplet loss, we need 3 imageries namely anchor, positive and negative.

# 6   Model Training

FaceNet by David Sandberg that delivers FaceNet representations constructed and proficient by means of TensorFlow, the project looks established although at the time of inscription does not provide a library-based installation. Usefully, David's project delivers a number of high-performing pre-trained FaceNet models and there are numerous of schemes that port or change these models for use in Keras. Keras FaceNet by Hiroki Taniai this project delivers a writing for adapting the Inception ResNet v1 model from TensorFlow to Keras this also provides a pre-trained Keras model ready for use. For pre-processing face alignment is done by using MTCNN. Problem with the method appears to be that the Dlib face detector false step some of the hard instances (partial occlusion, silhouettes etc). Which kinds the training set laidback that bases the model to perform poorer on other benchmarks. To solve this, other face landmark detectors has been tested. One face landmark detector that has proven to work precise in this situation is the Multi-task CNN. A Matlab/Caffe implementation has been used for face alignment with very good results. The CASIA-WebFace dataset has been used for training. Training set consists of total of 4,53,453 imageries over 10,575 individualities after face detection. Some presentation upgrading has been seen if the dataset has been sieved before training. The best carrying out model has been trained on the VGGFace2 dataset consisting of 3.3M faces and 9000 classes. The accuracy on LFW for the model 20180402-114759 is 0.99650+-0.00252. Note that the input images to the model need to be standardized using fixed image standardization

# 7  Contrastive Loss

Contrastive loss receipts the consequence of the system for an optimistic instance and computes its distance to an instance of the same class and differences that with the distance to negative instances. In other words, the loss is truncated if positive trials are encoded to alike or closer depictions and adverse instances are programmed to diverse or farther depictions.

This is proficient by taking the cosine objectivities of the vectors and giving the subsequent distances as forecast prospects from a distinctive categorization network. The vast idea is that you can extravagance the distance of the constructive instance and the detachments of the adverse examples as output probabilities and use cross entropy loss. In supervised classification, the system productions are characteristically track through a softmax function then the negative log-likelihood loss. Learning a vector depiction of a compound input like a picture is an example of dimensionality reduction.

"Dimensionality reduction aims to translate high dimensional data to a low dimensional representation such that similar input objects are mapped to nearby points on a manifold." — Dimensionality Reduction by Learning an Invariant Mapping, 2006.

The objective of operative dimensionality reduction is to acquire a novel lower dimensional demonstration that conserves the construction of the contribution such that distances between output vectors expressively seizure the alterations in the input. However, the vectors must seizure the invariant features in the input.

"The problem is to find a function that maps high dimensional input patterns to lower dimensional outputs, given neighborhood relationships between samples in input space." — Dimensionality Reduction by Learning an Invariant Mapping, 2006.

Dimensionality reduction is the tactic that Siamese networks practice to report one-shot learning.

In 2006 research paper titled "Dimensionality Reduction by Learning an Invariant Mapping," Raia Hadsell, et al. discover by means of a Siamese network for dimensionality reduction with convolutional neural networks with picture data and suggest training the models using contrastive loss. Nothing like other loss functions that may possibly appraise the recital of a model across all input instances in the training dataset, contrastive loss is considered between pairs of contributions such as between the two contributions provided to a Siamese system.

Sets of instances are provided to the system, and the loss function disciplines the model otherwise based on whether the modules of the examples are similar or diverse. Specifically, if the classes are the same, the loss function encourages the models to output feature vectors that are more similar, whereas if the classes differ, the loss function encourages the models to output feature vectors that are less similar.

"The contrastive loss requires face image pairs and then pulls together positive pairs and pushes apart negative pairs.However, the main problem with the

contrastive loss is that the margin parameters are often difficult to choose." —
Deep Face Recognition: A Survey, 2018.

Loss function requires that a margin is designated that is used to regulate the limit to which examples from different pairs are castigated. Selecting this boundary requires cautious consideration and is one problem of using the loss function.

$$\ell_{i,j} = -\log \frac{\exp(\mathrm{sim}(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau)}{\sum_{k=1}^{2N} {}_{[k \neq i]} \exp(\mathrm{sim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau)} \ ,$$

Contrastive loss appearances questionably like the softmax function. Because it is with the count of the vector resemblance and a temperature normalization influence. The similarity function is just the cosine distance as mentioned before. The other alteration is that standards in the denominator are the cosign distance from the constructive instance to the negative examples, unlike from Cross Entropy Loss. The instinct here is that we want our alike vectors to be as close to 1 as possible, since -log(1) = 0, that's the optimal loss and the negative examples to be close to 0. Since any non-zero values will reduce the value of similar vectors.

# 8   Triplet Loss

The impression of comparative loss can be extra extended from two examples to three called triplet loss. Triplet loss was introduced by Florian Schroff, et al. from Google in their 2015 paper titled "FaceNet: A Unified Embedding for Face Recognition and Clustering." Relatively calculating loss based on two instances, triplet loss involves an anchor example and one positive or corresponding example and one adverse or non-matching example.

The loss function disciplines the model such that the detachment between the corresponding instances is abridged and the detachment between the non-matching examples is increased.

"It requires the face triplets, and then it minimizes the distance between an anchor and a positive sample of the same identity and maximizes the distance between the anchor and a negative sample of a different identity." — Deep Face Recognition: A Survey, 2018.

# 9  Results of Model

# 10  Conclusion

# 11  References

https://machinelearningmastery.com/one-shot-learning-with-siamese-networks -contrastive-and-triplet-loss-for-face-recognition/

https://towardsdatascience.com/one-shot-learning-face-recognition-using- siamese-neural-network-a13dcf739e

https://sorenbouma.github.io/blog/oneshot/ https://medium.com/intro-to-artificial-intelligence/one-shot-learning- explained-using-facenet-dff5ad52bd38 https://blog.floydhub.com/n-shot-learning https://www.cs.cmu.edu/ rsalakhu/papers/oneshot1.pdf http://cis.csuohio.edu/ sschung/CIS660 $//arxiv.org/pdf/1711.06025v2.pdf\,https://arxiv.org/pdf/1605.06065v1.pdf$