

Face Vision : Real - Time Face Detection From Video Streams

Aim → To built a Deep Learning model which will take input as video video as input and draw a bounding box around face wherever the location of the face in video is.

Data Collection → We will require a high amount of data to train the neural network deep learning model. We have captured a total ~~384~~^{total 384} photos from laptop camera which is 1280-width and 720-height size. We will augment this data to create 30 times photos which we will store in the format of augmented versions of each image. Initially we have only 384 face images captured from the laptop front camera.

Data Setting and Labelling of Images →

Dataset was uploaded to Kaggle. This dataset was copied to the working directory of the Kaggle note book. Then we created "labels" directory to store the labels of each image. For labelling the face in images we used "labeling" rebuild library. References are as follows → linkhub link = (), [Video-link] = (). When I followed the video and on the labeling tool, it allows us to label the image with square bounding box around face which stores different information about face [object] in xml format. The output label of the image [image link] is given as [label link] in the xml format. We converted all the xml formatted labels into json formatted labels because we will be using Albumentations for a data augmentation. Then we created three folders called "train", "test" and "val" and inside each of them we created images and labels folders to store images and labels. (**FOLDERS STRUCTURE ON NEXT PAGE**)

stagge working /

- face customer image dataset
 - o Data (originally captured images)
 - o Train
 - o Images
 - o labels
 - o Test
 - o Images
 - o labels
 - o Val
 - o Images
 - o labels
 - o Labels (xml formatted labels)
 - o labels (json formatted labels)

Now we will split the total data/images available in 70% for training, and 15% each for testing and validation. So we have ~~388~~ 384 images, after splitting we will be having 268 images in train, 57 images in test and 59 images in val. We will store these respective number of images in train/images, test/images and val/images respectively. Now we will have to store the json formatted labels of each image in train, test and val in train/labels, test/labels and val/labels directory respectively.

Applying Image Augmentation on Images and Labels using Albulmations -

Image Augmentation - It is the process of creating new training images by applying random transformations to existing images.

Albulmation - It is a python library for fast and flexible image augmentation, this is especially useful in computer vision projects.

Then we defined a augmentation pipeline called "augmenter" using augmentation library which randomly alters images and bounding boxes to improve model generalization. Then we sanity checked whether image file is accessible and correctly read before doing further processing like augmentation or training. We parsed down the bounding box coordinates and class of the regions. Now we will apply augmentation on train and validation data. To store this augmented data, we created one folder at ~~src~~ facecustomizedatadataset/aug-data and reside that created 3 subfolders train, test and val and inside each of these 3, we created images and labels folders to store augmented images and their corresponding labels. Now we will apply augmentation on train, test and val data and store the images in aug-data folders along with their corresponding labels.

Then we resized the image to 120 * 120 pixels for standardizes input size to the neural network and normalized the pixel values to be in the range of [0, 1] by dividing by 255. This will give us all resized images in (120, 120, 3) with all pixel values scaled between 0 to 1.

Now we have to prepare the labels for augmented images, and we will do it with the help of loading class labels and bounding box for each image using its corresponding label file (.json).

Now we have train-image - ~~src~~ facecustomizedatadataset/train-images, train-labels, test-image, test-labels and test-labels and val-image, val-labels. We have to combine train, test and val image with train, test and val labels. We will shuffle the images and define batch size at 8.

Now before training the model and defining the model, we will previsualize the training image whether the class and bounding box is drawn properly or not.

Building Deep Neural Network Model using Keras Functional API

Keras functional API \Rightarrow it is a way to build more flexible and complex neural network architectures compared to the simple sequential API. Here API means Application Programming Interface provided by keras. We will use pretrained model of vgg16 which is pretrained on large dataset called "Imagenet". We will freeze the classification layer of vgg16 and use only convolution layers of vgg16 and we will add our custom classification and regression model to the output of vgg16 convolution layer. We have 2 problems to solve here, one is regression \Rightarrow Identifying the coordinates of the bounding box and another one is classification \Rightarrow Detecting whether face is available in the video / image or not.

Here valid model architecture \rightarrow [link].

Building Custom CNN model

- Define input layer with shape of $(120, 120, 3)$ because we have input pixels with 120×120 with RGB.
- Passed input layer through vgg16 imagenet pre-trained model with frozen classification layer.

- Now we will define two different models, one for classification and one for regression. F_1 is for classification and F_2 is for regression
- While output is passed through both F_1 and F_2
- F_1 output is passed through the class 1 dense layer with 2048 nodes and ReLU activation function.
- F_2 output is passed through both the regress 1 dense layer with 2048 nodes and ReLU activation function.
- class 1 output is passed through class 2 layer which is also the output for classification model.
- regress 1 output is passed through regress 2 layer which is also the output for regression model.
- At the end, we will combine the output of both classification and regression model.

Loss function & Training of the model -

For fine-gained learning rate scheduling which will help the model converge more smoothly, we calculate learning rate decay per batch. We set Adam optimizer for training our neural network in TensorFlow with learning rate decay strategy.

For classification model, we will use loss function as binary crossentropy. We define localization loss explicitly for predicting where an object is located within an image. Predicting bounding boxes is not an a classification task, so it requires different loss than cross-entropy. Localization loss is a custom loss function that penalize that penalizes the model based on how far off its predicted bounding box is from the ground truth.

Now we will define custom keras model class `FaceTracker`, designed for a multi-task learning problem of face detection where we perform both classification and regression.

`FaceTracker` class → combines classification and regression in a custom model.

`train_step()` → custom training logic, forward pass, compute losses, backpropagate

`test_step()` → custom testing logic, forward pass and loss calculation.

`compile()` → adds support for custom losses (classification + localization).

`call()` → Defines model behavior when called.

Now after defining custom keras model class `FaceTracker` we will create an instance of the custom `FaceTracker` class by passing the `facebox` model to it. We will compile the model with adam optimizer with learning rate and decay, binary crossentropy loss for classification and localization loss for regression. We will train the custom `FaceTracker` model using `fit()` method for 10 epochs with `val` as validation loss and `call-back`.

Output plots and testing model on test data

We have drawn three plots after training the deep learning model and those are → total loss (1) epochs, classification loss (2) epochs and regression loss (3) epochs. Then plotted the visualized the bounding box on test data.

Testing Model on Video Stream

We captured a video of nearly 1 min and saved it in local machine. Then uploaded this video in kaggle data set directory. We passed this video frame through the model and detected faces in a video every 0.3 seconds, drawn a bounding box around them if detected, and display sampled frames with annotations.