

MARCH 2021

EFFICIENT FACE DETECTION + TRACKING IN VIDEO SEQUENCES BASED ON DEEP LEARNING

Notes

Face tracking is Hard. ^{deviation}
I will always have some error rate. But its
effectiveness can be influenced by a few
things

- 1) target ^{whole individual} moves too fast screws up the
update to the tracking
- 2) the face it's self rotates or some
details about the face change
again messing up the tracking
- 3) The environment impact (change of
lighting, etc.) cause difficulty

Deep ML has exploded in popularity
It has shown promising results for
face detection.

In this paper they will propose a efficient face detection & tracking Deep ML Framework.

1. they first propose model for initial face extraction
2. they propose efficient model that can deal with Movement, lighting changes, etc
3. they take care of this Scale deviation Problem

1. Face detection

People have tried many things.
Most recently, YOLO.

Face Detection Falls into 2 categories

- a) Region Based
- b) Sliding Window Based

Proceeds to explain relevant work in both

Part 4 continued But Now on
face tracking

3 types of face tracking

a) feature-Based

b) model - Based

c) learning - Based

a) has poor effects Because of the
difficulty to extract features
from target-in-motion.

b) also not ~~that~~ good due to afformation
problems caused by the real world. & they
have high latency

2. there face Detection model

so this is a Rather Deep
convolutional Net with residual
connections. Inbetween the conv
layers they note that they add
something called a SEN

Layer, this layer is supposed to fuse the features of each conv layer. It "recalibrates" Learns the importance of each feature. it "enhances" useful info + "suppresses" useless info. It works by after the conv first do Squeeze then Excitation.

Squeeze - why? in early layers conv has small receptive field so this leads to poor understanding of global feature.

So to fix this we compress the whole feature map for each channel. $H \times W \times C \Rightarrow 1 \times 1 \times C$

this is now the summary statistic for the feature globally + we will use it in the next step.

Excitation - This is where we do the "recalibration", we will figure out which feature relate, which are important, which are redundant etc.

it does this via a gating mechanism where it assigns importance via learned weights. If you look at the equation in the paper it feels a little like an attention mechanism. & after these linear layers (2 of them) is passed through a sigmoid to get importance score.

LOSS functions

so sigmoid is very popular in image classification. But it can't meet the needs of inter-class distance increase + intra-class variance decrease.

They go through iterations on the softmax function improving it to fit the needs better till they reach Additive-margin-softmax the best one

Neural Nets

Now for the face tracking model

So they will be using a regression based model where we will extract facial features from 2 adjacent frames then predict the next spot.

the Network has 4 parts

P1 -

P2 - P4

P3 -

P1 + P2 are the earlier RNFT models. P4 is a FC layer + P3 is a correction Network.

the correction Network helps with more stable learning since there can be scale deviation this layer is supposed to adjust the size of the face detected.

Model training:

if we know position of face

in last 2 frames then it can't move much to the next frame. This is good prior knowledge. Because of that

new center is $c'_x = c_x + w \cdot \Delta x$ \leftarrow Random
 (c_x, c_y) = center now $c'_y = c_y + h \cdot \Delta y$ \leftarrow Var
 w, h = width, height \leftarrow sampled from Laplace Distribution

Input to the model is

1. video sequence
2. target face in the initialization window

we detect (c_x, c_y) from #2

For Evaluation

2 Data sets

1. images of faces
2. videos from VTB

Results show their models very good.