

# Face Recognition

Dienstag, 21. Juli 2020 17:11

face verification vs. face recognition

1:1

1:k

One shot Learning - learning from one example to recognize again

-- not actually using CNN to give prediction

-- learn the similarity -> degree of difference between img1 and img2

use only one image -> CNN -> softmax according to each person --> won't perform well because dataset is too small; furthermore, if you have new image in order to predict the new image the whole CNN have to be trained again.

$d(\text{img}_1, \text{img}_2) = \text{degree of difference}$  if  $d(\text{img}_1, \text{img}_2) \leq \tau$  same ✓ } verification  
 $> \tau$  different }

Siamese network

$x^{(1)} \rightarrow \text{CNN} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{428} \rightarrow \text{encoding of } x^{(1)}$   
 $x^{(2)} \rightarrow \text{CNN} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{428} \rightarrow \text{encoding of } x^{(2)}$   
 $d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$

$\Rightarrow$  parameters of NN define an encoding of  $f(x^{(1)})$

Learn parameters to:

if  $x^{(1)}, x^{(2)}$  are same person  $\Rightarrow \|f(x^{(1)}) - f(x^{(2)})\|_2^2$  is small  
 else  $\Rightarrow$  large

Loss for siamese network -- triple loss

A: anchor (example) P: positive N: negative

want:  $\|f(A) - f(P)\|_2^2 \leq \|f(A) - f(N)\|_2^2$

$\Rightarrow d(A, P) \leq d(A, N)$

$\Rightarrow \|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 \leq 0 - 2 \uparrow$

$\Rightarrow \dots \dots \dots + 2 \leq 0$

given 3 images A P N

given 3 images A P N

2 1 0

$$f(A, P, N) = \max [ \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + 2, 0 ]$$

$$J = \sum_{i=1}^m f(A^{(i)}, P^{(i)}, N^{(i)})$$

training 10 k images of 1 k person

After training can be applied to one short problem

How to choose A P N

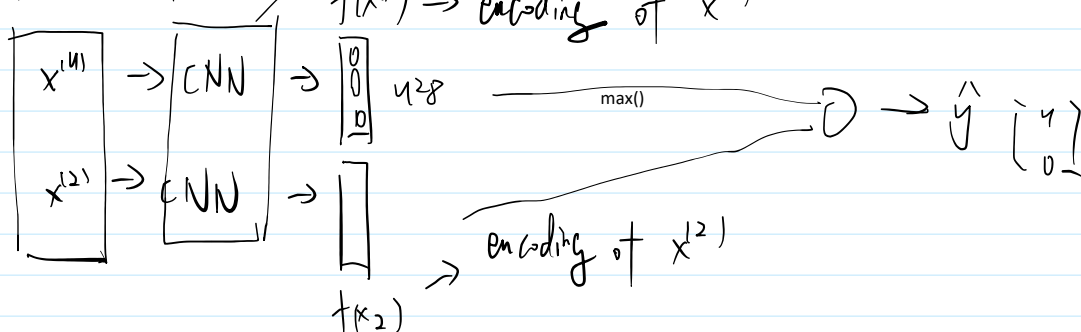
if choosed randomly, it is easy to satisfy  $d(A, P) + 2 \leq d(A, N)$  because A N distance large

-> learning algorithms will not learn much, will not perform well on one short problem

Choose triplets hard to train on  $\Rightarrow \underline{d(A, P)} \approx \underline{d(A, N)}$

Face verification as binary classification  $\rightarrow$  same CNN

input  $\rightarrow$  image pair



$$\hat{y} = \Delta \left( \sum_{k=1}^{128} w_k |f(x^{(1)})_k - f(x^{(2)})_k| + b \right) \Rightarrow \text{learn } w_k, b$$

k-th element

$$\text{or } \frac{(f(x^{(1)})_k - f(x^{(2)})_k)^2}{f(x^{(1)})_k + f(x^{(2)})_k} \rightarrow \chi^2 \text{ similarity}$$

Training also uses image pairs and labels (1 or 0).

Neural Style Transfer

how close are these two images' content

how close are these two images' style

Content C + Style S  $\rightarrow$  new image G

$$J(G) = 2J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

Content C + Style S → new image G

$$J(G) = 2J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

1. initiate G randomly
2. Use gradient descent to minimize J(G)

$$G = G - \frac{\partial}{\partial G} J(G)$$

## Content cost

use hidden layer l to compute content cost (l is usually mid-layer)  
use pre trained ConvNet (e.g. VGG)

$a^{[l]}(C)$  → activation of layer on content image

$a^{[l]}(G)$  → — — — — — generated image

if  $a^{[l]}(C)$   $a^{[l]}(G)$  similar ⇒ similar content on images

$$J_{\text{content}}(C, G) = \frac{1}{2} \|a^{[l]}(C) - a^{[l]}(G)\|^2 \Rightarrow \text{elementwise sum of square difference}$$

## Face verification and face recognition

### Siamese network

Neural Style Transfer  
Cost function of NST  
Content Cost Function  
Style Cost Function

## Take away from programming assignment:

### Neural Style Transfer:

1. NST uses a previously trained convolutional network, and builds on top of that. The pre-trained NN can be different, such as VGG in original paper.
2. The first step is usually load pre-trained model. The model is stored in python dictionary, the key-value represents each layer.
3. feed a image as input to the model by using: `model["input"].assign(image)`
4. access a specific layer when the network is running on image: `sess.run(model["conv4_2"])`
5. The shallower layers of a ConvNet tend to detect lower-level features such as edges and simple textures. The deeper layers tend to detect higher-level features such as more complex textures as well as object classes.
6. Make generated image G match the content of image C:
  - a. Choose a "middle" activation layer  $a[l]$
  - b. Forward propagate image "C"
  - c. Forward propagate image "G"
  - d. Compute the "content cost"

$$e. J_{\text{content}}(C, G) = \frac{1}{4 \times n_H \times n_W \times n_C} \sum_{\text{all entries}} (a^{(C)} - a^{(G)})^2$$

### 7. G(gram)i,j: correlation

1) The result is a matrix of dimension (nC,nC) where nC is the number of filters (channels). The value G(gram)i,j measures how similar the activations of filter i are to the activations of filter j.

### 8. Compute the style cost on one layer

### 9. Combine style cost with different layers with weights:

$$J_{style}(S, G) = \sum_l \lambda^{[l]} J_{style}^{[l]}(S, G)$$

10. compute style cost on each layer with weights
11. Total cost function

$$J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

After all these, in total:

1. Create an Interactive Session
2. Load the content image (load and reshape)
3. Load the style image (load and reshape)
4. Randomly initialize the image to be generated (slightly correlated to content image, which helps to match content image quickly)
5. Load the VGG19 model
6. Build the TensorFlow graph:
  - 1) Run the content image through the VGG19 model and compute the content cost:
    - Assign the content image to be the input to the VGG model.
    - Set a\_C to be the tensor giving the hidden layer activation for specific layer.
    - Set a\_G to be the tensor giving the hidden layer activation for the same layer. -- (assigned but not evaluated)
    - Compute the content cost using a\_C and a\_G,
  - 2) Run the style image through the VGG19 model and compute the style cost
  - 3) Compute the total cost, using initialized alpha and beta
  - 4) Define the optimizer, the learning rate, training step (optimizer.minimize(J))
  - 5) Initialize the TensorFlow graph, assign generated image to the model and run the train step tensor for a large number of iterations, updating the generated image at every step.

Face Detection:

1. face verification: learn a neural network that encodes a face image into a vector of 128 numbers. By comparing two such vectors, you can then determine if two pictures are of the same person. ---- "is this the claimed person?", this is a 1:1 matching problem
2. face recognition: "who is this person?" ---- this is a 1:K matching problem
3. Walk through face verification:
  - 1) Encoding face image into 128-D-vector
    - i. Use pre-trained ConvNet (weights): input batch of m faces images (m, n\_c, n\_h, n\_w); output encoder (m, 128)
    - ii. Evaluation criteria, if an encoding is a good one then:
      - 1) The encodings of images of same person are similar  $\leq$  threshold
      - 2) The encodings of images of different persons are different
    - iii. Using triplet loss function to define the criteria:

$$J = \sum_{i=1}^m \left[ \underbrace{\|f(A^{(i)}) - f(P^{(i)})\|_2^2}_{(1)} - \underbrace{\|f(A^{(i)}) - f(N^{(i)})\|_2^2}_{(2)} + \alpha \right]_+$$

$$\|x - y\|_2^2 = \sum_{i=1}^N (x_i - y_i)^2$$

- 2) train model to minimize the triplet loss
- 3) run forward propagation to encode each image (face) and store the f(x\_i) encoding in a database: key(name/id), value(128 encoding)
- 4) input new image (image to be verified) in to model to run the forward propagation to encode.
- 5) compute the distance between the encoding of new input image and the corresponding encoding in database
1. Walk through face recognition (takes an image as input, check if it is one of the authorized persons; and if so, who ):
  - 1) based on previous face\_verification, compute the input encoding
  - 2) find encoding from database that has smallest distance
    - i. Initialize a temp variable with a large enough number to update, used to keep track of what is the closest encoding to the input's encoding.
    - ii. Loop over the database dictionary's names and encodings.
      - a) Compute the L2 distance between the target "encoding" and the current "encoding" from the database.
      - b) If this distance is less than the temp, then update