

Face recognition:

face verification(1:1) VS face recognition(1:n)

one-shot learning problem: learn from one example to recognize the person again.

we can learn a similarity function:

$d(\text{img1}, \text{img2}) = \text{degree of difference between images}$  | face verification

Siamese network

$x1 \rightarrow f(x1)$ : the encoding generated from network

$x2 \rightarrow f(x2)$ : the encoding generated from the same network (same parameters)

$d(x1, x2) = \|f(x1) - f(x2)\|^2$

what you want to do is to learn parameters so that if any pair  $(x_i, x_j)$  are the same person, then  $d(x_i, x_j)$  is small. And what you can do is to use backpropagation to vary these parameters in order to make sure these conditions are satisfied.

Triplet loss

A: anchor, P: positive, N: negative

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

but if both are zeros or both are equal, then the network cannot be trained

so we add the hyperparameter  $\alpha$  as margin, pushing two sides further away from each other, like SVM

$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$

So the final loss function can be written as:

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

Minimize the loss function to update parameters so that the distance of two encodings of two images is small if two images are similar.

when you train the network, you need the training set like 10K pictures of a 1K person, during training, if A, P, N are chosen randomly,  $d(A, P) + \alpha \leq d(A, N)$  is easily satisfied. So you should choose triplets that are hard to train on, that is  $d(A, P) \sim d(A, N)$

After train the network, you can use it in a one-shot learning problem.

Face verification and binary classification

$x1 \rightarrow f(x1)$

$\rightarrow$  logistic regression classification (binary classification, alternative to triplet loss)

$x2 \rightarrow f(x2)$

rather than just take in the two encodings, you can compute the difference of two encodings like the following:

The image shows a handwritten formula for binary classification using face encodings. The formula is:

$$\hat{y} = \sigma \left( \sum_{k=1}^n w_k \underbrace{\frac{(f(x^{(1)})_k - f(x^{(2)})_k)^2}{f(x^{(1)})_k + f(x^{(2)})_k}}_{\text{difference of two encodings}} + b \right)$$

The formula uses the sigmoid function  $\sigma$  to output a probability  $\hat{y}$ . The input to the sigmoid is a weighted sum of the squared differences between the encodings of two images,  $x^{(1)}$  and  $x^{(2)}$ , across  $n$  dimensions. The weights  $w_k$  are learned parameters. The bias  $b$  is also a learned parameter. The denominator in the fraction is the sum of the two encodings,  $f(x^{(1)})_k + f(x^{(2)})_k$ .