

BLG 561 E FALL 2021

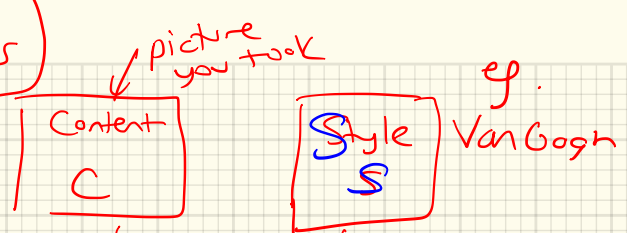
Deep Learning

16.11.2021

Görde ÜNAL

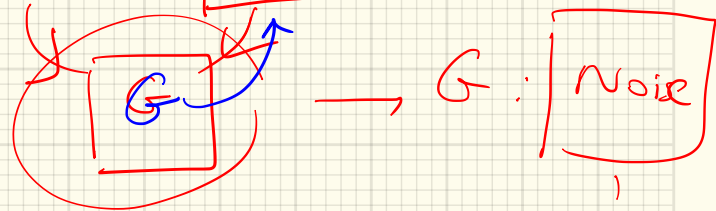
Designing Cost/Losses

for Style Transfer:



$$J(G)_{\text{overall}} = ?$$

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



$$\hat{G} = \arg \min_G J(G)_{\text{overall}}$$

your picture in the style of Van Gogh

1 Content Loss:

pretrained ConvNet (VGG)

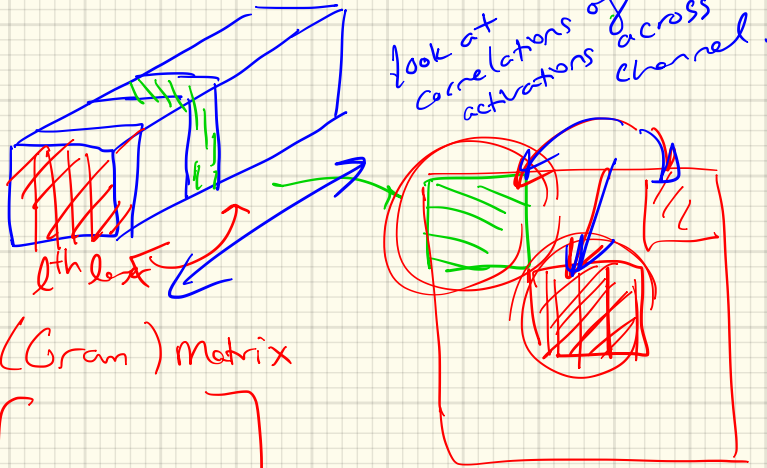
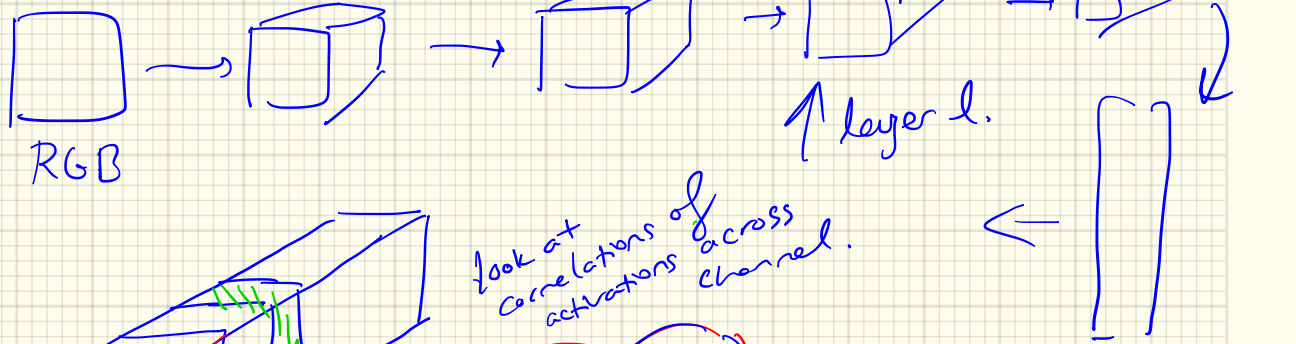
Let $a^{(l), C}$, $a^{(l), G}$: l : typically middle layers

similar content

(feature loss) (perceptual loss)

$$J_{\text{content}}(C, G) = \sum_{l=k_1, k_2, \dots} \left\| \frac{a^{(l), C}}{\downarrow} - \frac{a^{(l), G}}{\uparrow} \right\|_2^2$$

② Style Loss :



look at correlations of activations across channel.

Assumption:
Certain set of activations coexist to produce a certain style

Style (Gram) matrix

$$G^{[l],s} = \begin{bmatrix} G_{kk'}^{[l]} \end{bmatrix}$$

$$G_{kk'}^{[l]} = \sum_{ij} \sum a_{ij,k}^{[l],s} a_{ij,k'}^{[l],s}$$

$$J_{\text{style}}^{(l)}(G, S) = \frac{1}{\# \text{pixels}} \left\| \underline{G}^{(l), S} - \underline{G}^{(l), G} \right\|_F^2$$

$$J_{\text{style}}^{\text{overall}}(S, G) = \sum_l J_{\text{style}}^{(l)}(S, G)$$

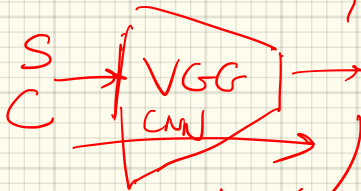
hyperparameter set $l = 1, \dots, \# \text{layers}$

Overall Loss: in Style Transfer:

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

$\{a^{l, S}\}_{l=1}^{\# \text{layers}}, \{a^{l, C}\}$

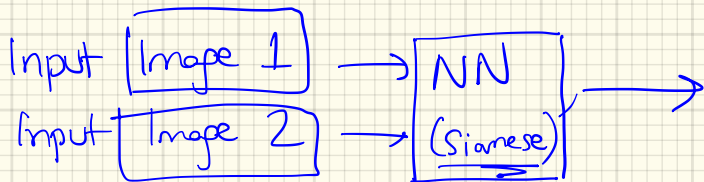
$\min J(G)$



no training

One-Shot Learning / Metric Learning

Use CNNs / NNs,

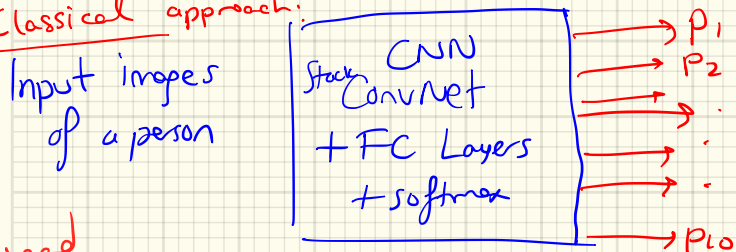


Learn a "similarity"

How similar are the two input images.

Say you have a Face Recognition System, say you have 10 employees:

→ Classical approach:



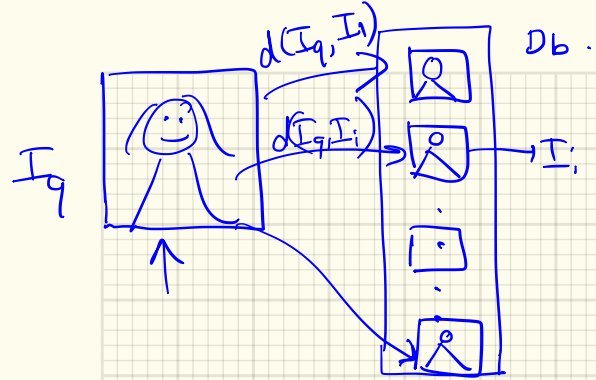
Cons

- 1) To train such a system, you require many different images of each of the 10 persons
- 2) When a new employee joins, you have to retrain!

Instead
→ One-Shot Classification

NN learns: $\begin{cases} d(I_1, I_2) \leq \tau & \text{for the "same" person} \\ d(I_1, I_2) > \tau & \text{for different persons} \end{cases}$

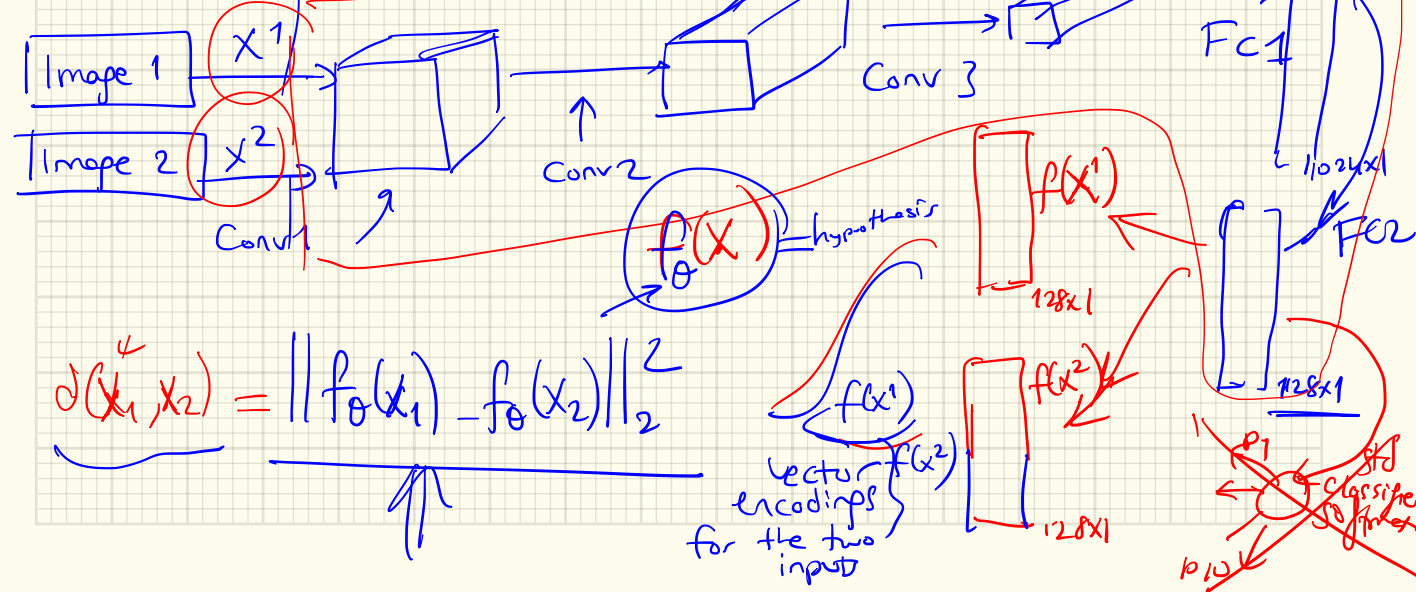
Similarity fn.



For the matching picture, we want $d(I_q, I_j) \approx \text{small number}$.

goal.

Siamese Networks:

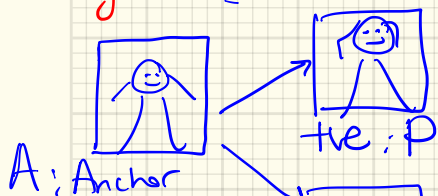


Goal: Learn parameters θ s.t.

If $x^{(i)}$ & $x^{(j)}$ are of the "same" person $\|f_{\theta}(x^i) - f_{\theta}(x^j)\|_2^2$ "small"
 " " " " "different" " $\|f_{\theta}(x^i) - f_{\theta}(x^j)\|_2^2$ "large"

What is the loss?

eg. Triplet Loss (FaceNet, 2015, Shreff)



Triplet loss minimizes:

$$\min L(A, P, N)$$

A, P: different pictures of the same person

$$L = \max(0, \underbrace{\|f_{\theta}(A) - f_{\theta}(P)\|_2^2}_{\text{smaller}} - \underbrace{\|f_{\theta}(A) - f_{\theta}(N)\|_2^2}_{\text{margin}})$$

as long as this part is -ve / loss $\rightarrow 0$

α : margin parameter to avoid learning trivial $f(x) = 0$

α : margin
 hyperparameter

Overall Triplet loss: $J = \sum_{i=1}^m \mathcal{L}(A^i, P^i, N^i)$

— Say you have 30K pictures, choose/sample your triplets;
You need multiple pictures of the same person in the training set to form (A, P) pairs.

→ Choosing triplets is tricky: $d(A, N) \gg d(A, P)$ when:

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

→ $d(A, P) + \alpha \leq d(A, N)$ easily satisfied.

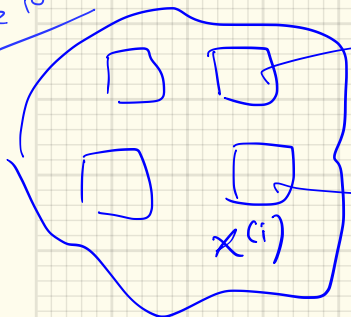
→ Choose triplets that are hard to train on:

$$d(A, P) \approx d(A, N) \wedge d(A, P) + \alpha \leq d(A, N)$$

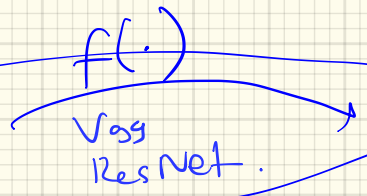


Hard Negative Mining Problem!

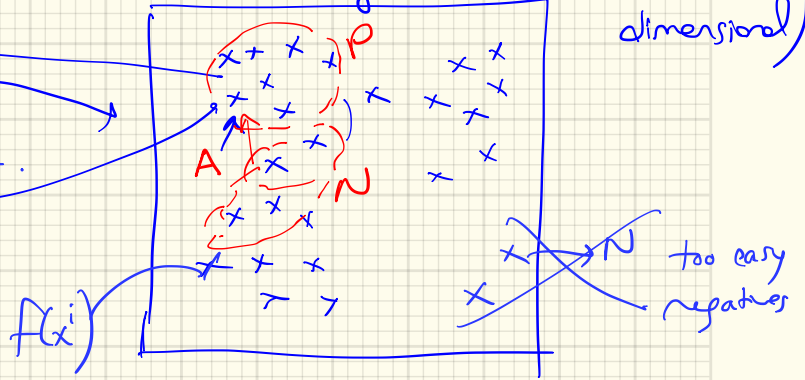
one idea



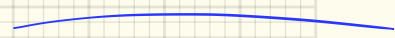
High-Dim input data manifold



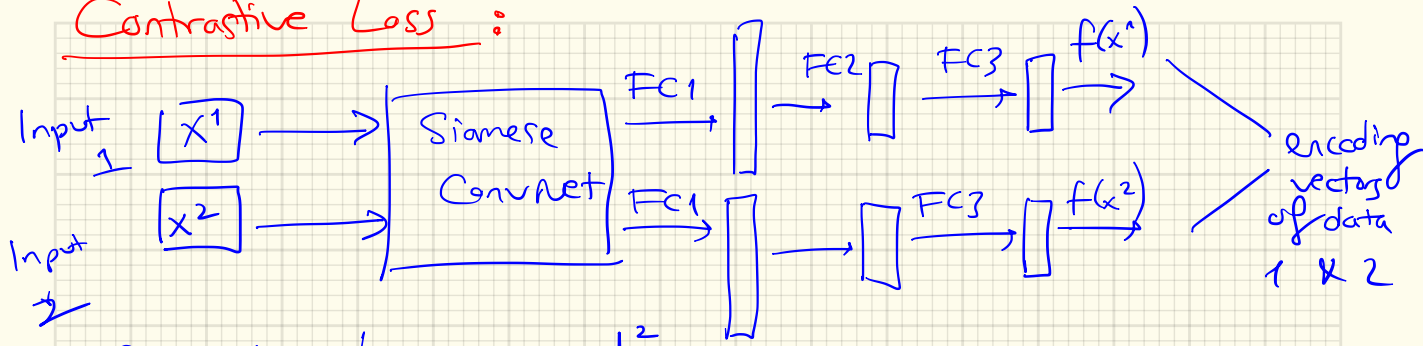
Embedding Space (Lower-dimensional)



Cluster in a low-dim embedding space, they use L^2 norm to pick (A, P, N) triplets.
eg.



Contrastive Loss :



Contrastive Loss

$$d(x^1, x^2) = y \cdot \overbrace{\|f(x^1) - f(x^2)\|_2}^d^2$$

$$+ (1-y) \cdot \max(0, \alpha - d)^2$$

Label:

$y=1$; if the two data are similar

$y=0$; o/w

margin

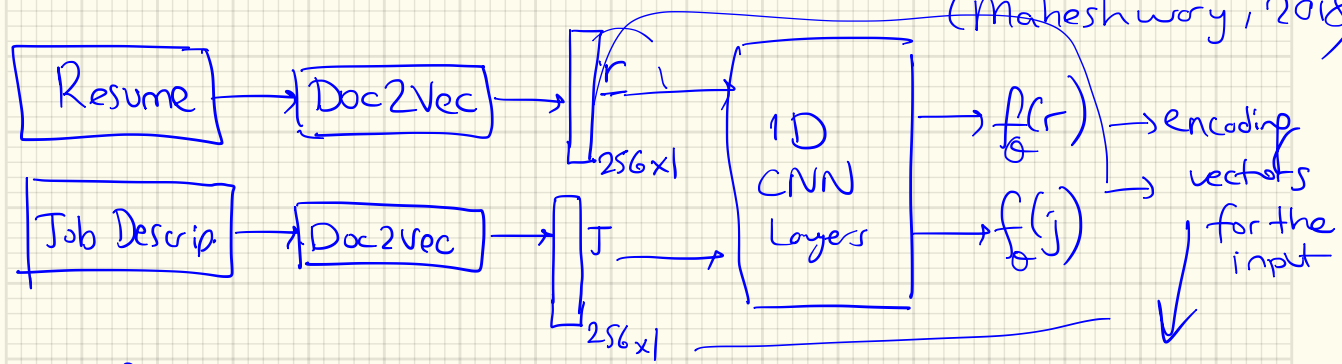
$$\text{if } \|f(x^1) - f(x^2)\| > \alpha$$

0 loss.

\therefore If x^1 & x^2 are dissimilar,
then their distance should be at least α , o/w a +ve loss is incurred.

Ex: Matching Resumes to Jobs via Deep Siamese Networks

(Maheshwary, 2018)



(r, j) pairs : 1 match $y=1$
0 no match $y=0$

$$d = \left\| \frac{f(r)}{\theta} - \frac{f(j)}{\theta} \right\|_2$$

use contrastive loss

→ many hyperparameters to tune !!