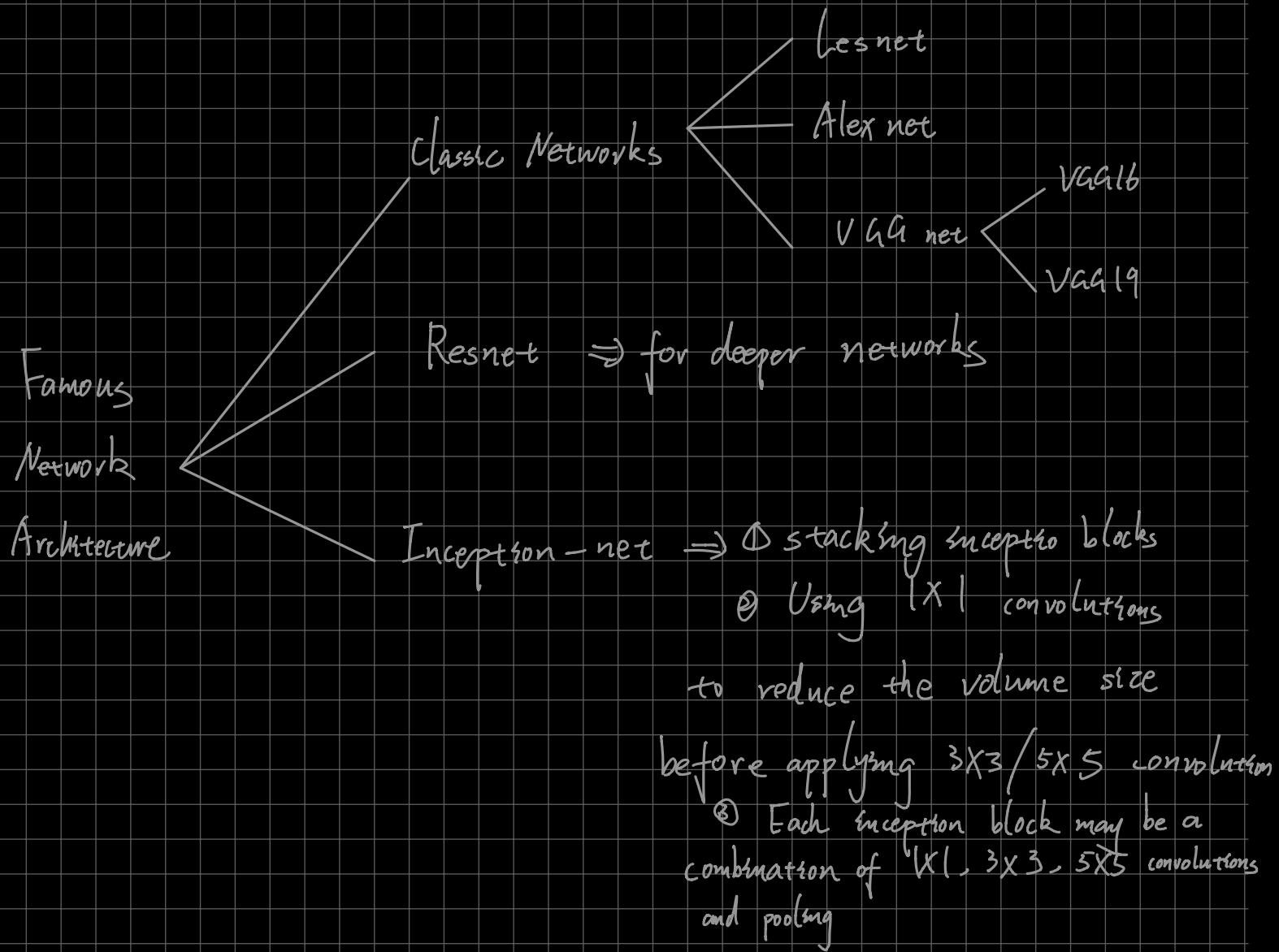
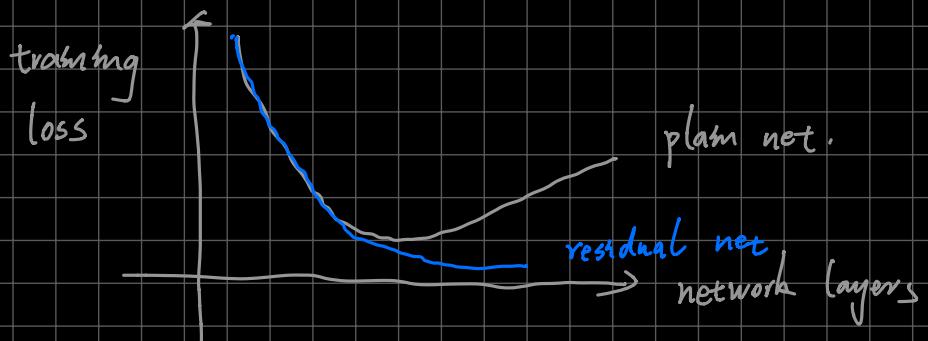


Plain networks: non-residual networks. \Rightarrow Theoretically - deeper ^{larger} the network, the training loss will be lower. But in fact, this is not true for plain network.



扩展详细知识见 CU: YOLO. Object Detection 最简单

Resnet \Rightarrow The way to build a deep network

How: By stacking several Resnet blocks, each block has skip connections inside it.

Explanation: Resnet blocks with the shortcut make it very easy for one of the blocks to learn identity function. That's why you can stack on additional Resnet blocks with little risk of harming performance.

\triangle There is some evidence that the ease of learning an identity function accounts for Resnet's remarkable performance more than they help vanishing gradients do.

\triangle The skip connections compute a complex non-linear function of the input to pass to a deeper layer in the network.

Some facts about YOLO (You Only Look Once)

Why it is popular ?

High accuracy + Quick (run in real time)

Key thoughts \Rightarrow NMS (Non-Max Suppression)

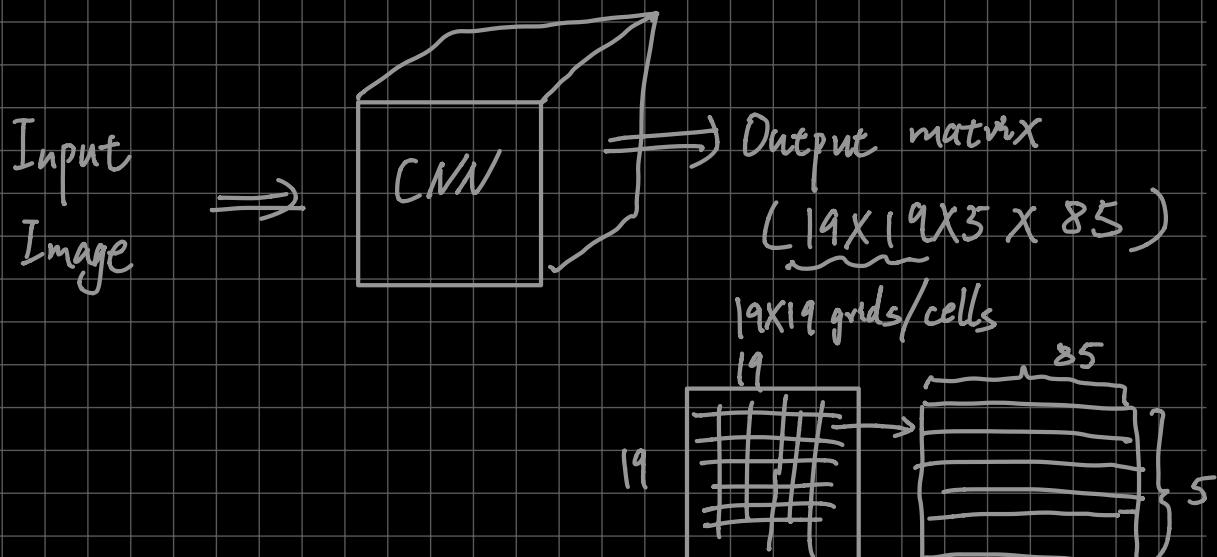
Step 1: Get rid of boxes with a low score

Step 2: Select the box that has the highest score.

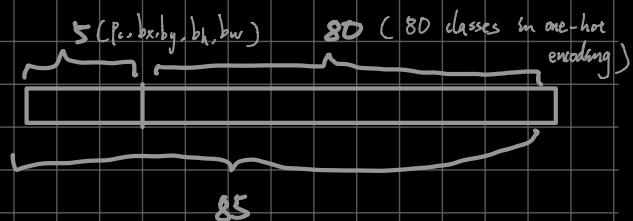
Get rid of the boxes that: $\text{IoU} \geq \text{threshold}$

Intersection over Union

Network structure:



5 anchor boxes , each have :



Face Recognition

Different Problems:

Face verification

Input: image / name / ID

Output: Whether the input is
that of the claimed person

(Sometimes we call it an
one-to-one problem)

E.g.: 笔记本电脑登录人脸识别
(只识别是否是你)

Face recognition

image

ID, if the input is any
of the K persons in the
database (this is the difference)
Not recognized, otherwise.

公司签到人脸识别

(识别所有公司的员工)

One Short Learning Problem:

For most face recognition APPs, you can only get one training
example (one short image of a person)

Difficulty: ① You have too little data to train the model
② Once you need to add a new person, you need
to change model's parameters.

How to solve? Calculate the similarity of an image.

Intuition:

- ① Similarity 能较大地提取出一张图片的特征
相比 CNN，它提取的更彻底
- \Rightarrow solve ①
- ② Similarity 针对每张图片，因此无需改变参数
 \Rightarrow solve ②

Siamese Network: 连体网络

What?

Parameters of NN is a good encoding of input

images $X^{(i)}$!

How to train?

Learn parameters so that:

- If $X^{(i)}, X^{(j)}$ are the same person. $\|f(X^{(i)}) - f(X^{(j)})\|^2$

is small.

- If $X^{(i)}, X^{(j)}$ are different person. $\|f(X^{(i)}) - f(X^{(j)})\|^2$
is large. even slightly different

Loss: triplet loss

anchor image \swarrow positive image \downarrow negative image \searrow

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

↓
margin, to prevent the network get trivial
solutions $\Rightarrow f(A) \equiv f(N) \equiv f(P) \equiv 0$

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

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

△ In training state, you need to ensure that there are lots

of training examples (e.g. 10k of 1k person).

While in applying state, you just need at least one shot of a person.

△ Need to choose triplets that're hard to train on. $d(A, P) \approx d(A, N)$

random choosing triplets are easily satisfied.

\Rightarrow push the siamese network to gain a high accuracy.

△ One way to do face recognition is to perform binary

classification (Input: a pair of images

Network: two similar sub-network

Output: 0 \Rightarrow not 1: verified

Neural Style Transfer

Cost function:

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

↑
content image
↓
Style image

use GD to minimize

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[L](C)} - a^{[L](G)} \|^2$$

Using pre-trained
⇒ VGG net

output of network L th layer
using content image as
input

activation of L th layer using generated
image as input.

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

$$J_{\text{style}}(S, G) = \sum_L \lambda^{[L]} J_{\text{style}}^{[L]}(S, G)$$

decay factor ↓

normalization term

$$\frac{1}{(2 n_H^{[L]} n_W^{[L]} n_C^{[L]})^2} \cdot$$

$$\sum_k \sum_{k'} \left(\sum_{i,j} G_{kk'}^{[L](S)} - G_{kk'}^{[L](G)} \right)^2$$

different channel's activation correlation

$$G_{kk'}^{[L](S)} = \sum_{i=1}^{n_H^{[L]}} \sum_{j=1}^{n_W^{[L]}} a_{ijk}^{[L](S)} \cdot a_{ijk'}^{[L](S)}$$

Gram matrix:

(compares how similar two matrix are.)

$$G_{kk'}^{[L](G)} = \sum_{i=1}^{n_H^{[L]}} \sum_{j=1}^{n_W^{[L]}} a_{ijk}^{[L](G)} \cdot a_{ijk'}^{[L](G)}$$

If they are highly similar, $G_{kk'}$ will be large.

Gram Matrix:

$$G_{\text{gram}} = \underbrace{A}_{\text{unrolled}} * \underbrace{A^T}_{\text{unrolled}}, \text{ shape: } (nC, nC)$$

$$A = (m, nh, nw, nC) \xrightarrow{\text{tf. transpose}} (m, nC, nh, nw) \xrightarrow{\text{tf. reshape}} \underbrace{(m, nC, nh \times nw)}$$

$G_{\text{gram}}(ij)$:

how similar the activation of channel i with
the activation of channel j

$G_{\text{gram}}(i, i)$: (对角线上元素)

how active a channel i is. (large \Leftrightarrow the
channel i is highly activated by input image)

△ In neural style transfer, we optimize the cost function $J(h)$ to update
generated image's $\underbrace{\text{pixel values. (not the network parameters)}}$

△ We apply transfer learning in NST (use VGG19 as
pre-trained network)

△ The layers chosen in computing J_{content} and J_{style} are middle
layers (not shallow layers, not too deep layers)

