

Position-wise Feed-Forward Networks

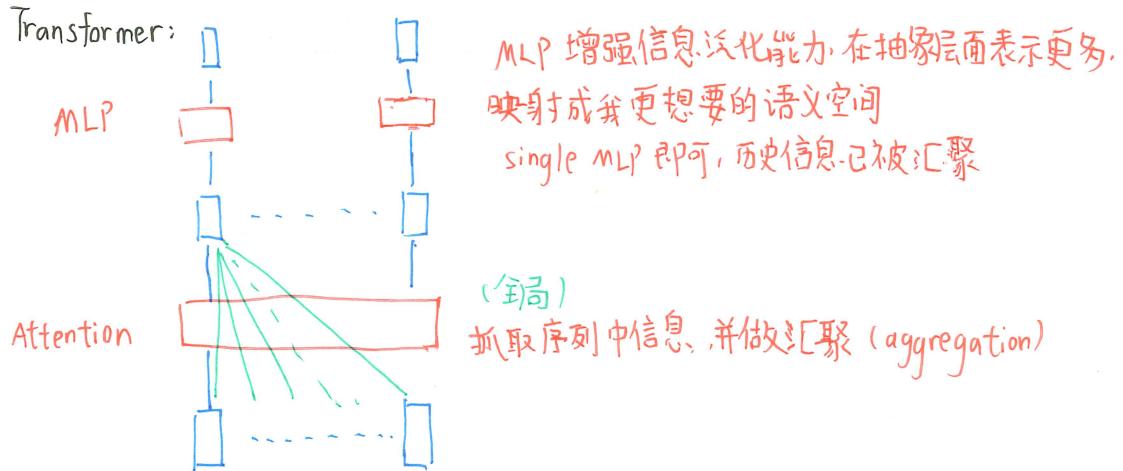
- Applied to each position separately and identically

$$FFN(x) = \max(0, xw_1 + b_1)w_2 + b_2$$

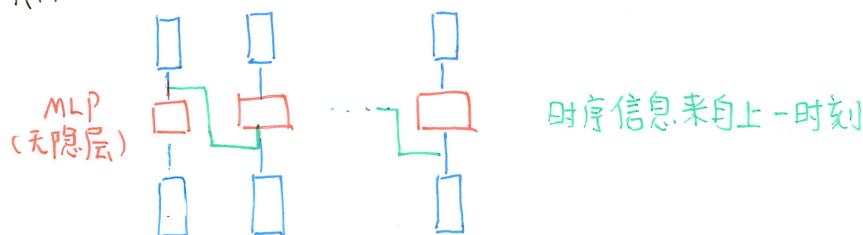
w_1 : project to $d = 2048$ 输入扩大四倍(中间隐藏层)

w_2 : project back to $d = 512$

Transformer:



RNN



Embedding

We use learned embeddings to convert the input tokens and output tokens to vectors of d model (512)

- same weight matrix between two embedding layers and the pre-softmax linear transformation.
- multiply weights by $\sqrt{d_{\text{model}}}$. 在学embedding时 L_2 long会比较小
乘以 d 以匹配 Positional Encoding 的 scale. (维度越大的向量归一化后单个值越小)

Positional Encoding

Attention 中无时序信息

$$PE(pos, 2i) = \sin(pos / 10000^{2i/d})$$

$$PE(pos, 2i+1) = \cos(pos / 10000^{2i/d})$$

用周期不同的 sin 和 cos 函数计算用来表示词的向量 ($d=512$)

与 Input embedding 直接相加, 两个都在 $[-1, 1]$ 附近

learn long-range dependency

Layer Type	Complexity / Layer	Sequential Operation	Max Path length
------------	--------------------	----------------------	-----------------

计算便利度(多少能并行计算) 信息糅合性

Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
----------------	------------------	--------	--------

Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
-----------	------------------	--------	--------

Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
---------------	--------------------------	--------	----------------

Self-Attention. (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$
------------------------------	------------------------	--------	----------

Campus Attention 缺点: 需要更多数据和更大模型, 费, 跑小训练集容易过拟合

Training

25000 source and target tokens

NVIDIA P100: 8

Base model: 0.4 sec / training step
100,000 steps

big models: 1.0 sec / training step
300,000 steps

Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$

$$\text{learning rate} = d_{\text{model}}^{-0.5} \cdot \min(\text{step_num}^{-0.5}, \text{step_num} \cdot \frac{4000}{\text{warmup_steps}^{1.5}})$$

模型越宽向量越长，学习率越低

从小到大

爬到最高后按 0.5 次方衰减

Regularization:

的输出

① Residual Dropout. 在子层(Multi-head, MLP) ^V 进入残差连接和进入 layer-alone 前做 dropout. $P_{\text{drop}} = 0.1 \rightarrow 10\% \text{ 的 } \times 0, \text{ 剩下 } 90\% \text{ 乘以 } \frac{10}{9}$

在 embedding & positional encoding 之后也 dropout 0.1
置信度

② Label smoothing: 对于正确的词 softmax 只需到 0.9 (实际很难逼近 1)
但会增加模型不确定性

Improving Language Understanding by Generative Pre-Training

No.

Date

generative

GPT：先在未标号 (unlabeled) 文件上生成一个预训练模型，再在有标号的子任务上训练分步的微调模型

在微调时构造任务相关 (task-aware) 的输入很少改变架构

pre-trained word embedding can improve performance in many NLP tasks
词嵌入模型

困难：①. 不知道用什么优化目标函数 (optimization objectives)
各有所长

②. 难以把学到的表示传递到下游子任务上面
子任务差异大，没有统一有效的方式

Semi-supervised — Combinations of unsupervised pretraining and supervised finetuning

架构：Transformer：

相比于RN，transformer 在迁移学习时学到的特征更稳健 (robust)

原因：Transformer 具有更结构化的记忆 (structured memory)

能处理更长文本，抽取句子和段落层面的语义信息

- 使用任务相关的输入表示

Unsupervised Pre-training

To maximize the following likelihood

$$L_i(u) = \sum_{j=1}^k \log P(u_i | u_{i-k}, \dots, u_{i-1}; \theta)$$

此处求和是因为 log 放回去就是 \prod
(hp) k: context window.

使用 Transformer 的解码器 (只有decoder 才有掩码)

The model applies a multi-headed self-attention operation over the input context tokens followed by position-wise feed forward layers

$U = (U_{-k} \dots U_{-1})$ is the context vector of tokens,

n is the number of layers.

W_e : Token Embedding Matrix W_p : Position Embedding Matrix

$$h_0 = U W_e + W_p$$

$$h_i = \text{transformer_block}(h_{i-1}) \quad \forall i \in [1, n]$$

$$P(u) = \text{softmax}(h_n W_e^T)$$

BERT: 做完形，用 Transformer 的编码器

GPT：用更难的目标函数，预测未来，解码器，更高天花板

Supervised fine-tuning

dataset C: each instance consists of a sequence of input tokens

x^1, x^2, \dots, x^m along with a label y

W_y : Linear output layer to predict y (parameter)

$$P(y|x^1, \dots, x^m) = \text{softmax}(h_m^T W_y)$$

★ To maximize the following objective.

$$L_2(c) = \sum_{(x,y)} \log P(y|x^1, \dots, x^m)$$

预测序列的标号

$$L_3(c) = L_2(c) + \lambda * L_1(c)$$

预测下一个词 超参

Finetuning on different tasks:

Classification.

start	Text	Extract
-------	------	---------

T

Linear

Entailment
蕴含

(互相推论)

Start	Premise	Delim	Hypothesis	Extract
-------	---------	-------	------------	---------

<S>

分隔定界

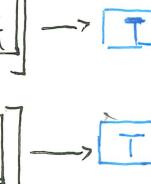
L

Similarity

Start	Text 1	Delim	Text 2	Extract
-------	--------	-------	--------	---------

\$

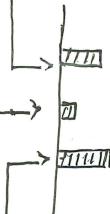
Start	Text 2	Delim	Text 1	Extract
-------	--------	-------	--------	---------



Multiple Choice.

Start	Context	Delim	Answer 1	Extract
-------	---------	-------	----------	---------

T



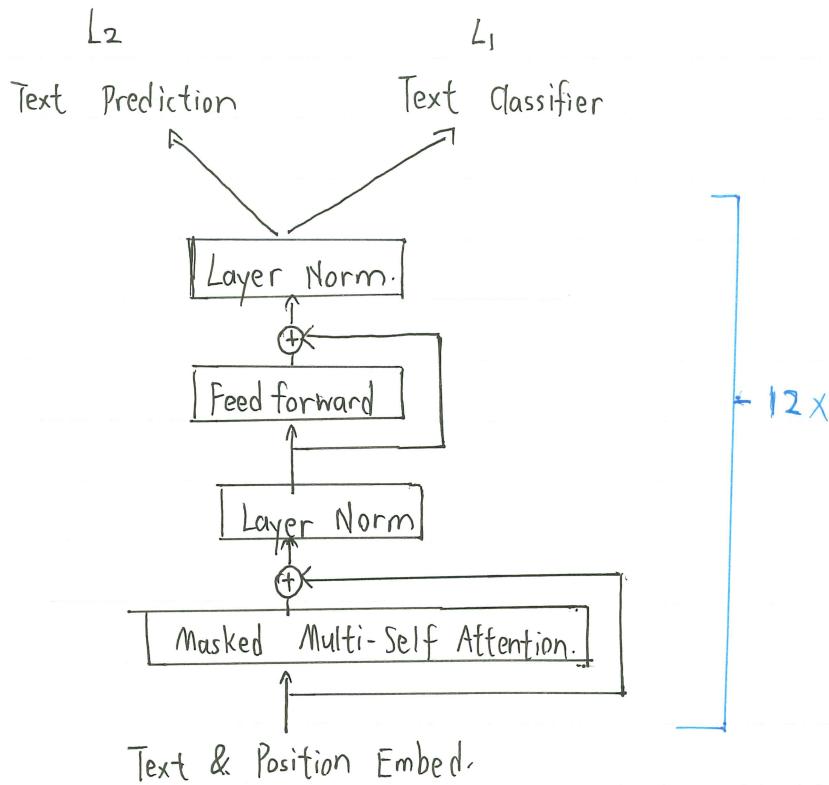
Start	Context	Delim	Answer 2	Extract
-------	---------	-------	----------	---------

T

Start	Context	Delim	Answer N	Extract
-------	---------	-------	----------	---------

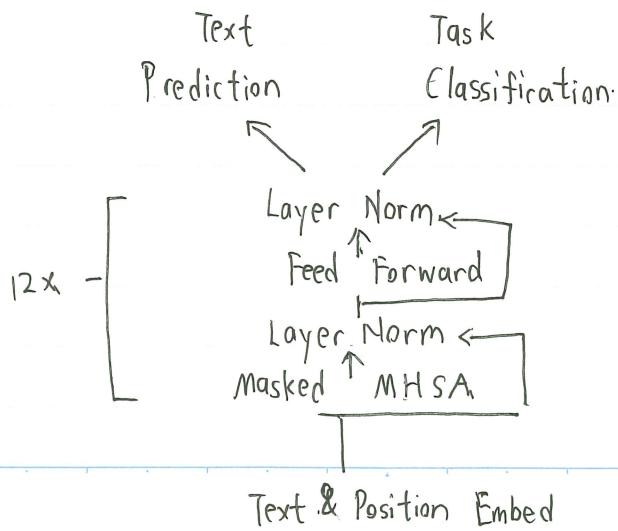
T

输入形式和输出构造无论怎么变，Transformer 模型不变



BERT 的数据集是 GPT 的 4 倍，精度更高

GPT-2 :



Language Models are Unsupervised Multitask Learners

Date

GPT2. 数据集 (wed text)

从新的角度看问题: Zero-shot, 更好泛化性

Estimate conditional distribution $p(\text{output} | \text{input}, \text{task})$

Language provides a flexible way to specify tasks, inputs and outputs all as a sequence of symbols.

e.g. (answer the question, document, question, answer)

Whether objective is supervised doesn't influence its global minimum.

爬了 reddit 上 karma > 3 的数据

Difference between GPT-2 Architecture. different from the transformer model such as.

	GPT-2	Transformer
Decoder Layers	12	6
Attention Heads	12	8
Dimensional	768	512
Context size	-	-
(for positional encoding)	1024	5000
Layer	GeLU	ReLU
Capture subwords	BytePair Encoding	Tokenization

Language Models are Few-Shot Learners

GPT3：回到 Few-shot，通过大量实验，

175 billion parameters, not sparse
contain a lot of 0

Applied without any gradient updates or fine-tuning.

微调效果好不代表大模型泛化好

* in-context learning - 在 multi task 中各任务内的上下文
无梯度更新；通过注意力机制处理长序列

Zero shot

Translate English to French: task description
cheese => prompt

One shot: task description

sea otter => butre de mer → example

prompt

缺点：难以载入长的新训练样本

每次都要往回全部抓取。

GPT2.3 相比于 GPT1 的改变：① modified initialization

②: pre-normalization

③. reversible tokenization 反转词元

* 计算复杂度与宽度是平方关系，和层数是线性关系

使用相对来说大的批量；计算性能↑，并行度↑，通讯量↓。

对于小模型降低批量大小可以增加噪音

训练数据基于过滤后的Common Crawl ① 基于 GPT2 数据分类过滤

② LSH 算法去重 (data retrieval)

③ 增加高质量数据集

训练过程：分布式训练、模型分割、数据分割

局限性：1. 文本生成新意度弱，没有推进，适合补全

2. 语言生成没有重点（模型会花很多时间学虚词）。能否分类以优化？

3. 缺乏多模态

4. 样本有效性不够

5. 贵

6. 决策方式未知（权重代表什么）

疑问：是基于样本学习还是在数据库中找相关？

Principle of prompt engineering

Principle 1: Write clear and specific instructions

- To avoid prompt injections

Tactic 1: Use delimiters.

triple quotes: """

triple backticks: ```

triple dashes: ---

Angle brackets: < >

XML tags: <tag> </tag>

2. Tactic 2: Ask for structured output

Tactic 3: Check whether conditions are satisfied
whether there're prompts?

Tactic 4: Few-shot prompting

Principle 2: Given the model time to think.

Tactic 1: specify the steps to complete the task

Tactic 2: Instruct the model to work out its own
solution before rushing to a conclusion
otherwise the model might neglect mistakes
in the process

Reducing hallucinations:-

First find relevant information, then

Answer the question based on evidence.

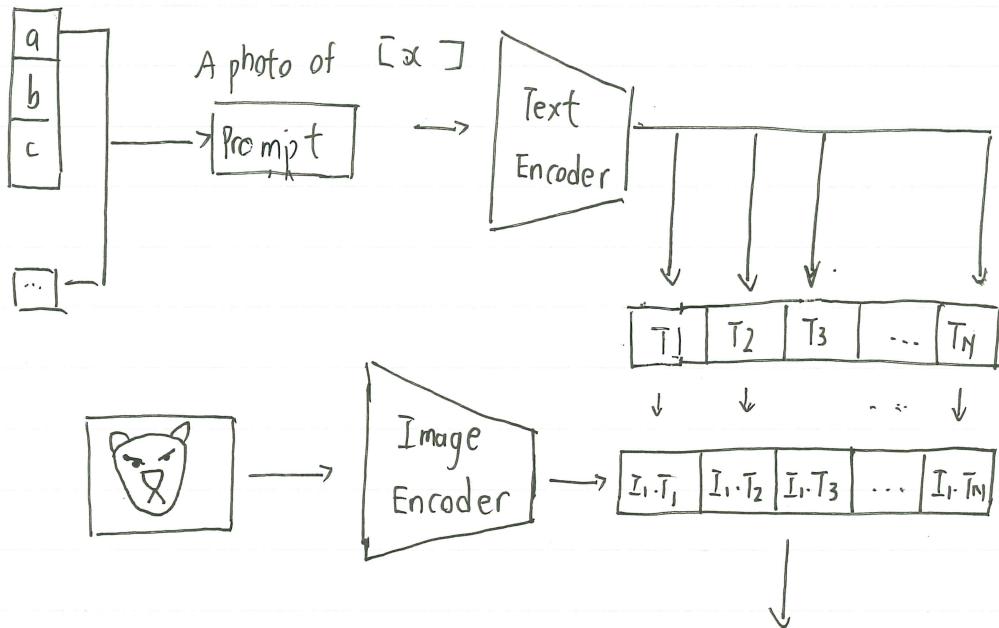
Prompt engineering is an iterative process

No.

Date

Campus

CLIP



- 通过对比式学习，省略生成式学习中繁琐实例细节，在抽象语义级别特征空间区分数据，拥有更简单的优化，更好的泛化
- 自回归预测 / 掩码完形填空 —> 自监督训练
目标函数与下游任务无关，而想要得到泛化能力强的特征

方法不新，但数据规模和模型规模实无前例

Transformer 的 Deep contextual representation learning
具备上下文语义环境学习方式，带来自监督范式下取之不尽用之不竭的文本监督信号

0.4 billion 数据(训练样本)

CV 与 NLP 的结合使得能够学习多模态的特征，
很容易做 zero-shot 的迁移学习

把训练任务变成对比任务 → 效率 × 4
(学习)

- Prompt engineering 的意义：
 - ① 多义性
 - ② 预训练时基本都是句子，只输入单词，抽取出来的特征可能不好
- CLIP 唯一的数据增强就是随机裁剪
- Linear Probe：冻住主体参数，抽特征，只微调分类头
- 过于难的任务还得 few-shot
- 推文 / 视频下的 hashtag 可以作为监督信号的来源
- Clip 把原超参数 temperature 优化为了可以学习的标量

伪代码：

# image-encoder	ResNet / Vision Transformer
# text-encoder	CBOW / Text Transformer
I [n, h, w, c] ^{224 224 3}	minibatch of aligned images
T [n, l]	minibatch of aligned texts
wi [d-i, d-e]	learned proj of image to embed
w-t [d-t, d-e]	learned proj of text to embed
t	learned temperature

extract feature representations of each modality

I-f = image-encoder(I)

T-f = text-encoder(T)

joint multimodal embedding 投影(单模态→多模态)

I-e = l2-normalize(np.dot(I-f, W_i), axis=1)

T-e = l2-normalize(np.dot(T-f, W_t), axis=1)

scaled pairwise cosine similarities [n,n]

logits = np.dot(I-e, T-e.T) * np.exp(t)

symmetric loss function

labels = np.arange(n) 正样本位于对角线, 与 MOCO 只在 0 号(第一位)正样本不同
ground truth

loss_i = cross_entropy_loss(clogits, labels, axis=0)

loss_t = cross_entropy_loss(clogits, label, axis=0)

loss = (loss_i + loss_t) / 2

局限性

- 想要达到专业领域或 SOTA 水平, 还需扩大 1000 倍 (Res 50 只是基线)
- 细分类 (如数个数, 异常分类) 情况下性能不好
- 泛化性在不属于 0.4 billion 的样本中很差, 如 MNIST
- 还得先给分类, CLIP 再判断 (对比式和生成式的目标函数能否融合?)

self-supervision | self-training 或许能提高数据利用率

自监督

伪标签

- 训练时重复使用已有监督训练数据集，可能有偏见
- 有时 one/ few shot 反而不如 zero-shot 好 —
复杂概念无法用语言描述

意义：

- 打破了固定标签种类的范式，只需搜集图片—文本的配对
用无监督的方式预测其相似性 / 生成
方便做推理、zero-shot

新颖度：100

有效性：100

问题大小：100 (灵活、多领域、高效)

Convolutional Neural Networks

Edge detection.

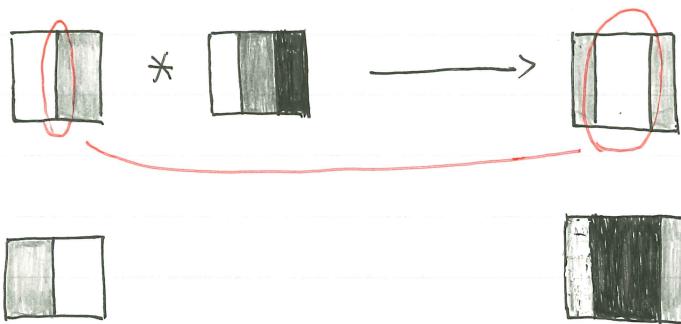
Vertical :

1	0	-1
1	0	-1
1	0	-1

Horizontal:

1	1	1
0	0	0
-1	-1	-1

filter kernel



1	0	-1
2	0	-2
1	0	-1

Sobel filter

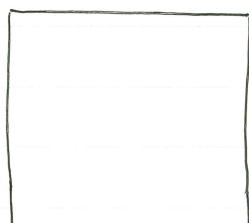
3	0	-3
10	0	-10
3	0	-3

Scharr filter

w ₁	w ₂	...
..
..

learnable ?

with
backpropagation



$f \times f$

$n-f+1 \cdot n-f+1$

Padding

$$n-f+1 \longrightarrow n+2p-f+1$$

0	0	0
0		
0		

Valid convolution : no padding $n-f+1$

Same convolution : (padding) $S_{\text{output}} = S_{\text{input}}$ $n+2p-f+1 = n \Rightarrow p = \frac{f-1}{2}$

By convention, f is always odd \rightarrow To align with "central pixel"

Strided convolution

$$\text{output: } \frac{n+2p-f}{s} + 1 \quad . \quad \frac{n+2p-f}{s} + 1$$

convolution in math textbook



Convolution in DL is Cross-correlation in Math

For convolution layer 1

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

Input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$n_H^{[l]} = \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1$$

$n_C^{[l]}$ = number of filters

Each filter: $f^{[l]}, f^{[l]} \cdot n_C^{[l-1]}$

Activation: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Weights: $f^{[l]}, f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$

bias: $n_C^{[l]} = (1, 1, 1, n_C^{[l]})$

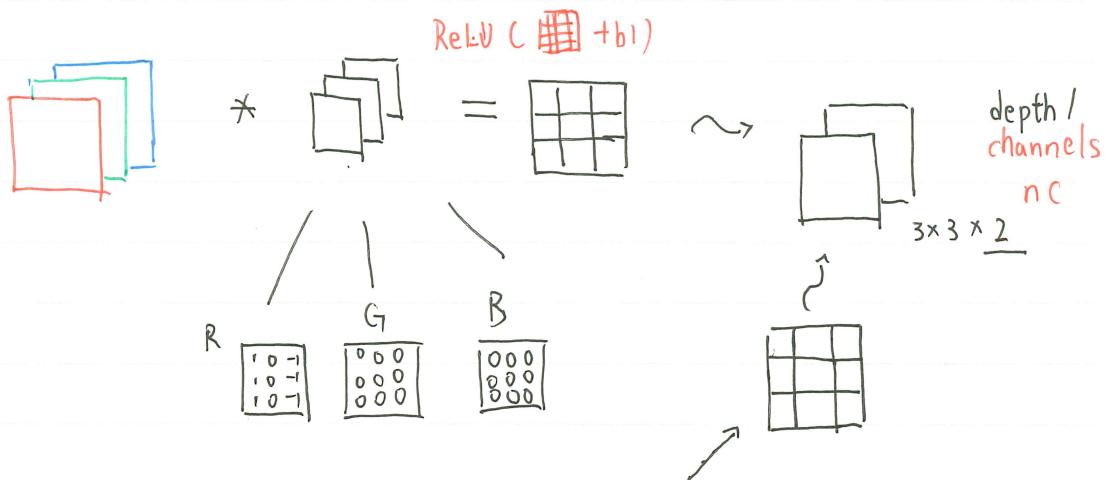
Type of layers

1. Conv

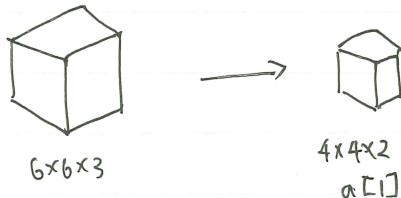
2. Pool

3. FC

Convolution on volume.



- We can have multiple filters focusing on different features.



N parameters remain same no matter how big input image is
— less likely to overfit

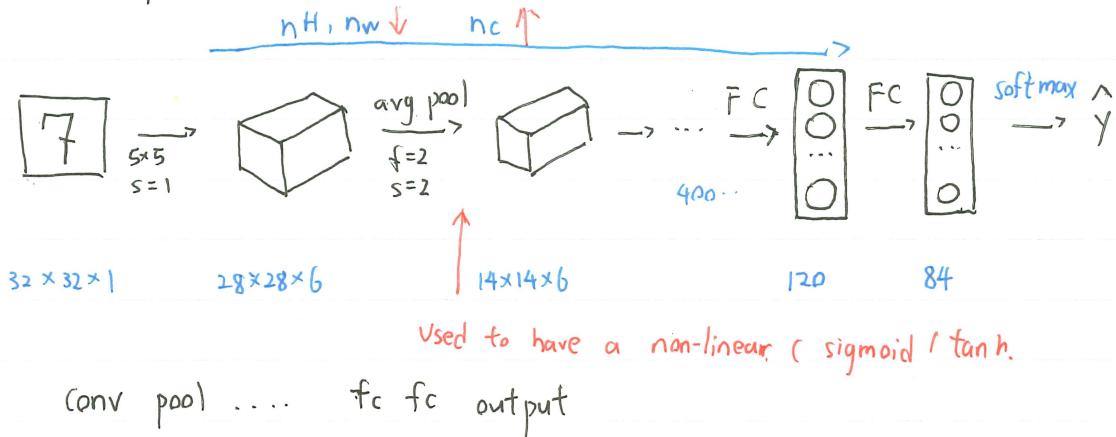
Max Pooling

Hyper Para: f.s.
 $2 \quad 2$ nothing to learn

Average Pooling (less likely to use)
Used when very deep

LeNet - 5

60k parameters.



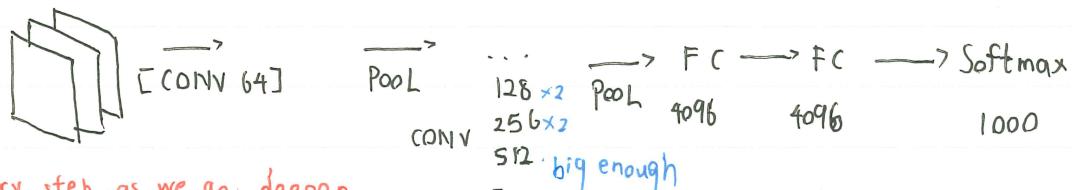
AlexNet. maxpool

- Similar to LeNet, but much bigger
- ReLU (solve gradient diffusion of sigmoid in deeper nets)
- Multiple GPUs
- Local Response Normalization
- Utilize Dropout



VGG-16 (simpler, less hp)

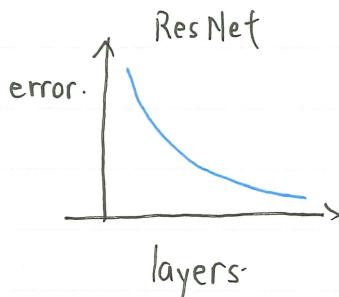
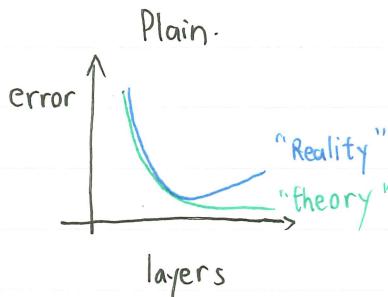
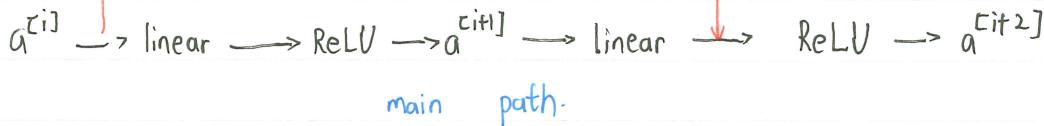
138 M parameters

 $\text{CONV} = 3 \times 3 \text{ filter}, s=1$. Same, $\text{MAX POOL} = 2 \times 2, s=2$  $nw \cdot nh / 2$ $nc \approx 2$ very systematic

uniform structure.

Res Net

Residual block

short cut / skip connection

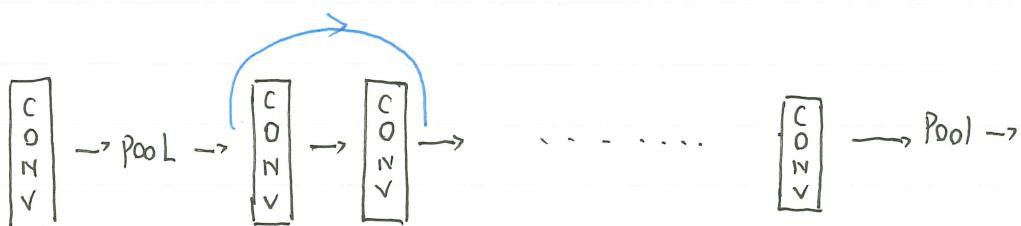
Allow for training deeper neural networks

Why Residual Network works

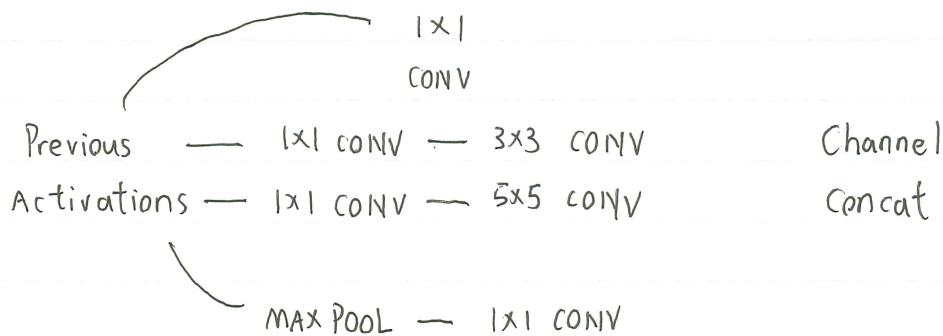
Normally $a^{[i]} = a^{[i+2]}$, it won't affect overall performance

More chances to choose correct parameter to improve.

* if dimensions of $a^{[i]}$ and $a^{[i+2]}$ are different, add a weight matrix



Inception block



- To prevent overfitting, we also add some softmax branches (like a regularizing effect)

Depthwise Convolution-

$$6 \times 6 \times 3 * 3 \times 3^{nc} \longrightarrow 4 \times 4 \times 3$$

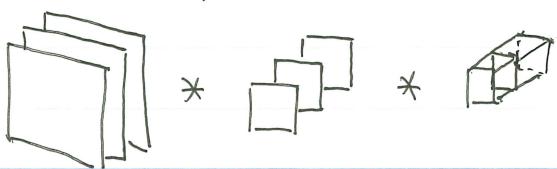
Cost = filter params \times filter positions \times # of filters

Pointwise Convolution

$$4 \times 4 \times 3 * 1 \times 1 \times 3^{nc} \longrightarrow 4 \times 4 \times 5^{nc}$$

Cost is similar to above.

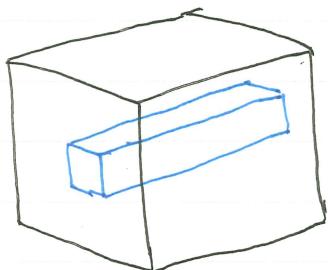
Depthwise Separable Convolution



most computational efficient
compared with normal convolution

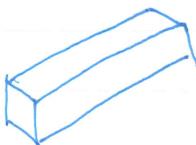
$\frac{1}{nc} + f^2$ more efficient

1×1 convolution: (network in network)

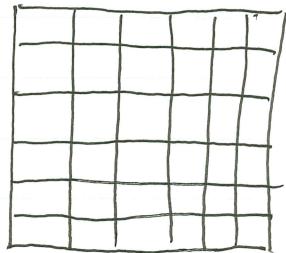


$6 \times 6 \times 32$

*



=



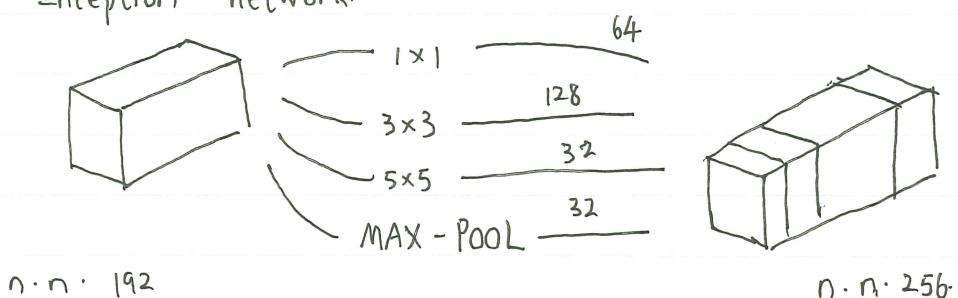
$1 \times 1 \times 32$

$6 \times 6 \times \# \text{filters}$

similar to a FC layer that applies to each positions

Adjust the channel. (升/降维特征矩阵)

Inception network.



No need to pick kernel size, let it learn

Problem: computational cost

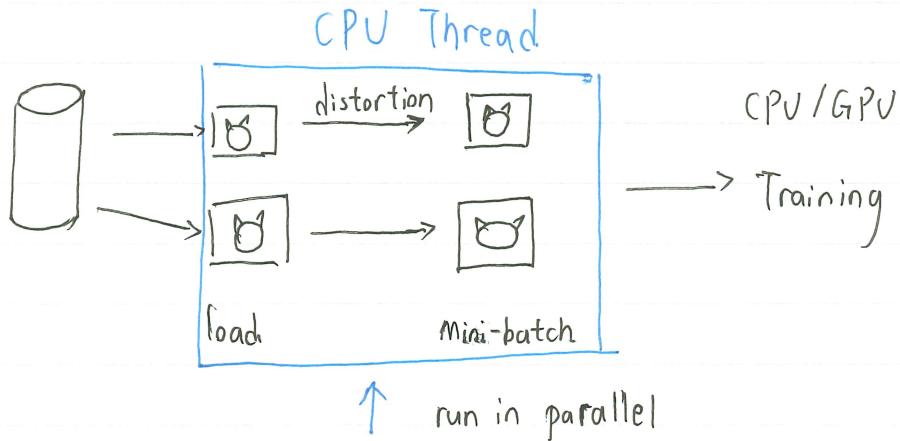
"bottleneck layer"



Computation 120 m \rightarrow 12.4 m

If bottleneck layer is implemented properly, you can shrink

Campus down representation size significantly while doesn't seem to hurt performance



Two sources of knowledge

- Labeled data
- Hand engineered features / network architecture

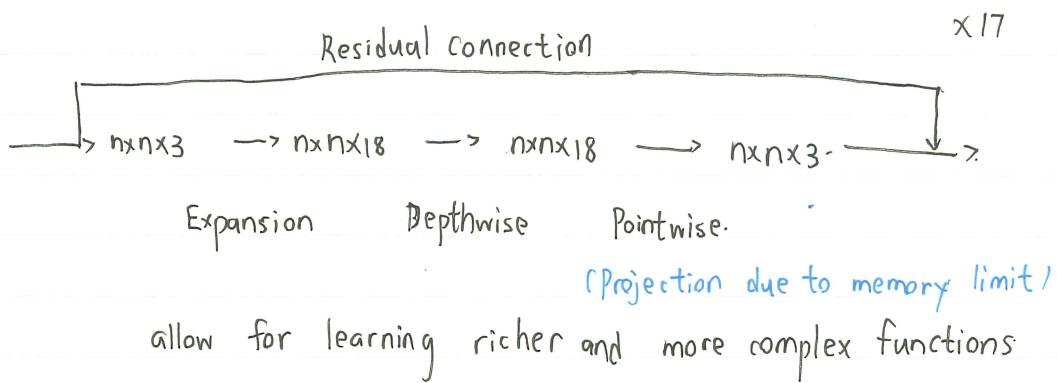
Tip. — Ensembling

- Train several networks independently and average their outputs

Multi-crop at test time.

- Run classifier on multiple versions of test images and average
Reduce speed: better not use it when building products results

Mobile Net V2 Bottleneck



Efficient Net: How to choose. r(resolution), d(depth) and w(width)?

Transfer Learning

- * Freeze previous parameters, modify the last layer
- * Use others' parameters as initialization.
- , More data \rightarrow less frozen layers

Data Augmentation.

- Mirroring
- Random Cropping
- Rotation
- Shearing
- Local Warping 局部约束
- Color shifting

PCA color augmentation: more changes on R.B. less on G

No.

Date

No.

Date

Campus