

Survey on DL@UAV vision

CONTENT



What is DL?



Application



Road Map



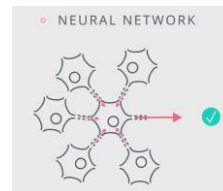
Start a Project

01

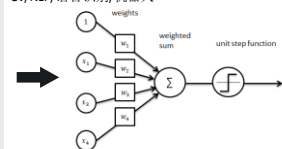
What is DL?

1_ What is DL

From Neurons to Perceptron

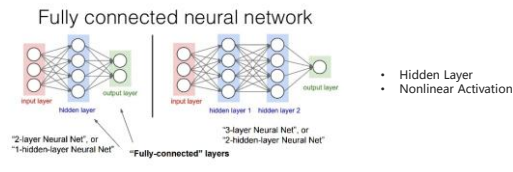


1957 年美国康乃尔大学的心理学家罗森布拉特 (Rosenblatt) 提出了感知器或感知机 (Perceptron) 的概念, 并试图用它来模拟动物和人脑的感知与学习能力, 而且用一个电子线路设计了著名的感知器神经网络模型, 称为仿脑机, 能够识别英文字母印刷体。深度学习应用的四个大的领域: CV/NLP/语音识别/机器人。



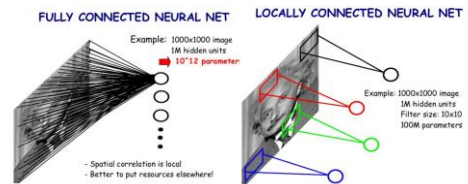
1_ What is DL

FC – Fully Connection



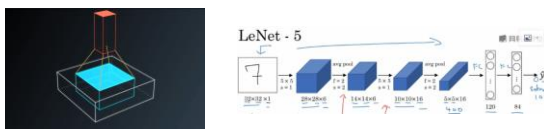
1_ What is DL

CNN – Convolutional Neural Network



1_ What is DL

CNN – Convolutional Neural Network



- Convolutional Kernel dimension
 $\text{width} * \text{height} * \text{input channel} * \text{filter number}$

1_ What is DL

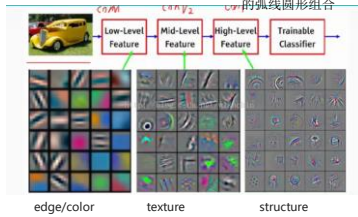
Why Convolution(compare with FC)



- Parameter sharing:**
A feature detector(such as a vertical edge detector) that is useful in one part of the image is probably useful in another part of the image
- Sparse connection**
In each layer, each output value depends only on a small number of inputs
- Less Parameters, Low Overfitting Risk**

1_ What is DL

CNN as Feature Extractor







CNN Visualization

CNN就是一个特征提取的工具,通过特征可视化的手段,我们可以还原出对某一层激活函数激活值最大的一些图像,可以看到从低到高,CNN在逐渐抽取越来越抽象的特征,开始时边缘/颜色块等,然后就是纹理,然后就是些结构,比如下面的蜂巢,上面的像眼睛一样的弧线圆形组合

02

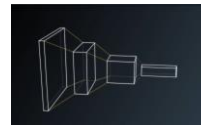
Application

2_ Application

-  Object Recognition
-  Object Detection
-  Visual Object Tracking(VOT)
-  Visual Slam

2_ Application

Object Recognition



Classification

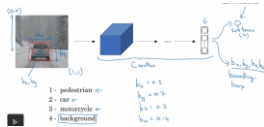
CNN用于OR的一般架构就是在空间上越来越小,深度上越来越深,就是把图像的信息从空间上抽象到深度上的过程,所以我们可以看到,CNN逐层空间大小在变小,深度在增加



Dogs vs. Cats in Kaggle

2_ Application

Object Detection



Classification + multi-object localization

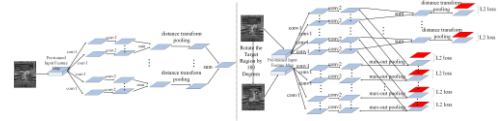


VEDAI UAV Object Detection Dataset

当用于OD时,其实和OR是相同的思路,这个时候抽象出来的信息不只是对物体类型的识别,还有位置的描述,通常会用 x,y,h,w 来描述一个bb,现在主流的做法是通过回归的方法让网络学习bb的伸缩和偏移,从而确定物体的位置

2_ Application

VOT

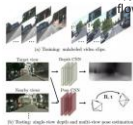


LSART - VOT 2017 Championship

- Exploit both the KRR and CNN as two complementary regressions for visual tracking
- CNN focuses on the small localized region
- KRR focuses on the holistic target

2_ Application

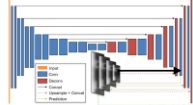
Visual Slam



SfM-Learner - CVPR 2017 from Google

- Use two isolated network (Depth CNN & Pose CNN) to predict depth & camera pose
- Base on photometric consistency principle
- See more on SfM-Net (also compute optic flow / scene flow / 3D point cloud)

- 它用两个网络分别/独自无监督地估计单帧的深度, 和视频序列中的camera的pose变化
- 光度一致性, 就是对于同一个物体的点, 在不同两帧图像上投影点, 图像灰度应该是一样的
- SfM-Net论文的核心思想也是利用 photometric constancy 来计算 pose, depth. 除此之外, 作者还计算了光流, scene flow, 3D point cloud 等. 可以说是 SfM-Learner 的升级版

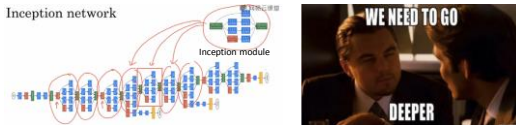


03

Road Map

3_ Road Map

Inception(GoogLeNet) – 2014 ImageNet Championship (Top-5 Error 9.2%)



- Much less Parameters ~5,000,000(1/12 of AlexNet, 1/27 of VGGNet)
- 1x1 convolution(Network in Network)
- Group Convolution(Inception Module)
- Replace FC with Average Pooling

3_ Road Map

数据标准

Single-crop, Single-model experimental results(一个样本截取一张图的, 单一模型不用集成学习的)

On ILSVRC 2012 dataset validation set

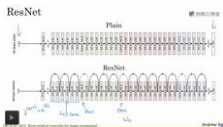
目前在ImageNet数据集上人眼能达到的错误率大概在5.1%, 这还是经过了大量训练的专家能达到的成绩, 一般人要区分1000种类型的图片是比较困难的

Inception Family

- Inception-V2 (Top-5 Error 7.8%)
 - Batch Normalization/two 3x3 replace 5x5
- Inception-V3 (Top-5 Error 5.6%)
 - Conv Factorization 3x3 -> 1x3+3x1
- Xception (Top-5 Error 5.5%)
 - Depth-wise separable convolution(separate depth conv & spatial conv)

3_ Road Map

ResNet– 2015 ImageNet Championship (Top-5 Error 4.49%)
2016 CVPR best paper



那么这里可以看到, 本来从上一层传过来的梯度为1, 经过这个block之后, 得到的梯度已经变成了0.0001和0.01, 也就是说, 梯度流过一个block之后, 就已经下降了几个量级, 传到前一层梯度将会变得很小!

这就是梯度弥散。假如模型的层数越深, 这种梯度弥散的情况就更加严重, 导致浅层部分的网络权重参数得不到很好的训练, 这就是为什么在Resnet出现之前, CNN网络都不超过二十几层的原因

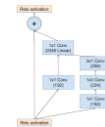
- Shortcut connection/solve Gradient vanishing problem
- Break the depth limitation

数据来自ResNet论文(Single-Model)

0.0001和0.01能够限制层数的只有内存大小了 – KM He KM He 2017年还有两篇ICCV best paper, 与object detection有关, 下个专题再讨论

3_ Road Map

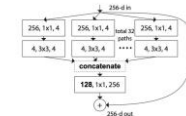
Inception-ResNet & ResNeXt(Top-5 Error 4.8% & 4.4%)



Inception Module + shortcut connection

天下大势, 分久必合, 合久必分
Cardinality和Group Conv的区别在于Cardinality是把一个Conv操作分成了拓扑结构完全相同的几个path, 这样方便加速

2017年最后一届ImageNet分类的冠军, Momenta的SENet也要研究一下
以及2017 CVPR best Paper, DenseNet, 可以看做是对ResNet的又一次大升级



Residual block + Group conv(cardinality)

3_ Road Map

Cost Comparison

CNN模型在不断逼近计算机视觉任务的精度极限的同时，其深度和尺寸也在成倍增长。巨大的模型无法在嵌入式平台上落地，就算想通过网络传输，较高的带宽也让很多用户望而生畏。另一方面，大尺寸的模型也对设备功耗和运行速度带来了巨大的挑战。所以，模型的小型化和加速成了亟待解决的问题。

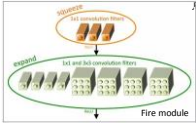
Model	Model Size(MB)	Mult-Adds	Million Para
AlexNet	200	720	60
VGG16	500	15300	138
GoLeNet	~50	1550	6.8
Inception-V3	90-100	5000	23.2

3_ Road Map

- Introduction of CNN Development
- Small CNN for mobile application
- Network Compression

3_ Road Map

SqueezeNet(ICLR-2017, Stanford)



1. 用1x1替换部分的3x3,减少参数
2. 减少进入3x3 filter的通道数,也就是先用1x1去压缩通道,就是squeeze
3. 少用pooling或者是conv stride,保留比较大的特征图,信息保留的多,最后的准确性更高,但这一点使模型解释的成本更高,运算代价更大了

- Replace part of 3x3 filters with 1x1 filters(small network)
- Decrease the number of input channels to 3x3 filters(small network)
- Downsample late in the network, conv layers has large activation map(high accuracy)

3_ Road Map

SqueezeNet(ICLR-2017, Stanford)

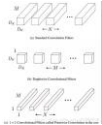
Table 2: Comparing SqueezeNet to model compression approaches. By model size, we mean the number of bytes required to store all of the parameters in the trained model.

CNN Architecture	Compressive Approach	Top-1 Accuracy (%)	Model Size (MB)	Top-1 Accuracy (%)	Model Size (MB)
AlexNet	None (Baseline)	56.7	248MB	56.7	248MB
AlexNet	HYD (Huang et al. 2016)	56.7	248MB	56.7	248MB
AlexNet	Network Pruning (Han et al. 2015)	56.7	248MB	56.7	248MB
AlexNet	Deep Compression (Han et al. 2015)	56.7	248MB	56.7	248MB
SqueezeNet (ours)	None	56.7	5MB	56.7	5MB
SqueezeNet (ours)	Deep Compression	56.7	5MB	56.7	5MB

- Small Network, but high compute cost(large activation map)
- Nothing Fresh in architecture

3_ Road Map

MobileNet(CVPR-2017, Google)



- 1 conv filter for 1 channel(totally separate spatial conv & depth conv)
- Low compute cost(compare with SqueezeNet)
- Almost equivalent accuracy level(SqueezeNet, AlexNet)

1. 将空间conv和深度conv完全分离开,提速
2. 相比SqueezeNet,计算量方面有明显提升(SqueezeNet似乎并没有考虑这个)
3. Accuracy上, MobileNet其实和AlexNet差不了太多,稍有提升

Table 3. Smaller MobileNet Compared to Popular Models				
Model	ImageNet Accuracy	Params	Multi-Add	FLOPs
MobileNet V1	71.0%	5.3M	1.1B	1.1B
SqueezeNet	71.0%	1.0M	1.2B	1.2B
AlexNet	57.0%	237M	2.3B	2.3B

3_ Road Map

ShuffleNet(CVPR-2017, Face++)

ShuffleNet的几个创新点

1. 提出了一个类似于ResNet的BottleNeck单元。借鉴ResNet的旁路分支思想，ShuffleNet也引入了类似的网络单元。不同的是，在stride=2的单元中，用concat操作代替了add操作，用average pooling代替了1x1stride=2的卷积操作，有效地减少了计算量和参数。
2. 提出将1x1卷积采用group操作会得到更好的分类性能。在MobileNet中提出，1x1卷积的操作占据了约95%的计算量，所以作者将1x1也更改为group卷积，使得相比MobileNet的计算量大大减少。
3. 提出了核心的shuffle操作将不同group中的通道进行打散，从而保证不同输入通道之间的信息传递。解决多个group conv叠加出现的边界效应问题。Channels Shuffle操作会导致内存不连续这个影响有待评估

- Channel shuffle for less cross-talk between channels(group conv issue)
- 1x1 group conv to reduce compute cost
- Use Residual block to raise accuracy

3_ Road Map

ShuffleNet(CVPR-2017, Face++)

Model	Complexity(MFLOPs)	Ch. err. (%)	Δ err. (%)
1.0 MobileNet-224	569	29.4	-
ShuffleNet 2x (g=3)	524	26.3	3.1
ShuffleNet 2x (with SF[3], g=3)	527	24.7	4.7
0.75 MobileNet-224	328	31.6	-
ShuffleNet 1.5x (g=3)	292	28.5	3.1
0.5 MobileNet-224	149	36.3	-
ShuffleNet 1x (g=8)	140	32.4	3.9
0.25 MobileNet-224	41	49.4	-
ShuffleNet 0.5x (g=4)	38	44.6	7.8
ShuffleNet 0.5x (shallow, g=2)	40	42.8	6.6

- Lower Complexity than MobileNet
- Higher Accuracy than MobileNet

MFLOPs(Million Floating-point Operations per Second) 衡量计算复杂度的指标

3_ Road Map

Small CNN Conclusion

- Separate convolution thoroughly (X-conv + Y-conv + Depth-conv)
- Raise utilization of channel information
- Residual block + Group Conv would be standard architecture

3_ Road Map



3_ Road Map

Deep Compression (ICLR-2016 best paper, Stanford)

1. 剪枝, 去掉一些不重要的连接(训练出来几乎为0权重的连接), Han Song在ICLR颁奖的演讲这儿说的很有意思, 剪枝相当于一个人喝大了, 但又没有丧失意识, 可以做一些简单的任务, 所以对简单任务, 我们不需要那么多的冗余连接
2. 权值共享是用少的bit数来存之前都是浮点数的权重, 而且只存权重的shift和code book, 就相当于训练这边编密码, 解释的时候用code book来解码
3. Huffman按照符号出现的概率来进行变长编码, 可以进一步减轻存储参数的Memory压力, 过去图像压缩经常会用到

- Pruning/Sparse Connection
- Quantization/Weight sharing
- Huffman Encoding

3_ Road Map

Deep Compression (ICLR-2016 best paper, Stanford)

1. 几乎没有准确率的损失
2. 如果和小模型搭配使用, 将更有利于在移动端部署
3. 虽然模型被压缩了, 但模型的解释成本增加了, 这一点还要double confirmed一下

因为模型的稀疏连接, 需要特殊的库来处理, 所以这种压缩方法在模型解释的时候一定是增加了成本

Network	Original Size	Compressed Size	Ratio	Accuracy	Compression
GoogLeNet	522MB	274KB	46%	89.36%	80.14%
VGGNet	274MB	149MB	45%	89.43%	45.54%
ResNet	165MB	8.36MB	95%	92.21%	95.18%
ResNeXt	199MB	1.73MB	91%	93.04%	91.36%
MobileNet	145MB	2.46MB	98%	93.87%	98.32%
ShuffleNet	1.66MB	2.47MB	98%	93.87%	98.32%

- Almost no accuracy lost
- Combining with small network would be more useful in embedded platform
- High model interpretation cost? Need to be verified

3_ Road Map

Channel Pruning (ICCV-2017, Face++)

1. 结构化剪枝, 不会造成稀疏连接, 所以解释的成本要更低
2. 把求解通道mask β 当做一个优化问题求解
3. 对于ResNet进行通道剪枝的时候, 增加一个sampler block, 使两个input的通道一样



- Structure Pruning. Lower model interpretation cost
- Compute mask β by solving a optimization problem
- Add sampler block for multi-branch networks(e.g. ResNet)

$$\min_{\beta} \frac{1}{N} \sum_{i=1}^N \left\| \sum_{k=1}^K \beta_k \mathbf{A}_k \mathbf{W}_k \right\|_F^2 \quad \text{s.t.} \quad \beta_k \in \{0, 1\}, \quad \forall k$$



3_ Road Map

Channel Pruning (ICCV-2017, Face++)

Decrease of top-1 error (%) when keeping 80% of			
Structure	1.25	1.5	1.75
Jacobson et al. [13] (1.5x mag.)	9.7	20.7	—
Alvarez [14]	8.38	13.34	—
Filter pruning [11]	—	0.8	0.6
Other method (our mag.)	—	—	—
Other (without fine-tune)	2.7	7.9	22.8
Other (fine-tune)	6	13.7	13.7

Table 1: Accelerating the VGG16 model [11] using a top-1 error of 2.7% (0.1% smaller in binary).

1. 和Deep Compression不一样, channel pruning会造成精度损失,但是结构化剪枝将来模型解释的成本低
2. 剪枝完fine-tune一下效果会好很多
3. 增加了很多hyper-para,比如正则化系数 λ ,还有压缩的比例

- Accuracy Lost
- Need fine-tune after pruning
- More hyper-para(regular coefficient λ , compress ratio) need be manually set, increase tuning difficult

3_ Road Map

Squeeze the Last Bit Out (AAAI-2018, Alibaba)

1. 基于低bit网络表示
 2. ADMM算法来解决
 3. 增加了很多hyper-para,比如正则化系数 λ ,还有压缩的比例
- 看了Squeeze the last bit之后,把这里丰富一下
这里还应该有一个对于压缩方法的总结

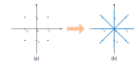


Figure 1: To extend spatial accuracy, we first convert the quantized weights into binary values (0 or 1) and then use the ADMM algorithm to solve the problem.

- Low Bit Networks
- ADMM algorithm
- Even higher accuracy in VGG-16, Compression has regularization effect

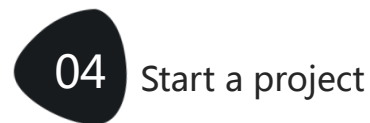
Model	Top-1 Error (%)	Top-5 Error (%)	Top-10 Error (%)	Top-20 Error (%)	Top-30 Error (%)	Top-40 Error (%)	Top-50 Error (%)	Top-60 Error (%)	Top-70 Error (%)	Top-80 Error (%)	Top-90 Error (%)	Top-100 Error (%)
VGG-16	7.37	21.02	28.12	35.12	39.12	42.12	44.12	45.12	46.12	47.12	48.12	49.12
VGG-16 (Squeeze the Last Bit Out)	7.37	21.02	28.12	35.12	39.12	42.12	44.12	45.12	46.12	47.12	48.12	49.12

(16)

3_ Road Map

Network Compression Conclusion

- Pruning
- Quantization
- Low-Bit Network(Binary Network, Ternary Network)
- Low rank decomposition
- Teacher-student Framework



4_ Start a project

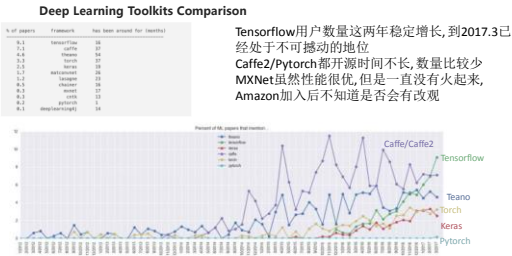
- Deep learning Toolkits
- Distribution Learning
- Deployment on Mobile
- Issues when Training

4_ Start a project

Deep Learning Toolkits Comparison

- Caffe从UC Berkeley的实验室出来, 最开始的作者是贾扬清, 后来被Facebook收购, 因为Caffe灵活性不足, 贾扬清在Facebook内部又搞了Caffe2, 估计Caffe用户一两年内都要迁移到其他框架上去
- PyTorch目前也由Facebook团队维护
- Theano开创了深度学习框架符号编程的先河(相比于caffe的命令编程), 17年停止维护
- Tensorflow由Google团队维护, 所以我们看到前面Google发的paper框架一定都是tensorflow的
- MXNet之前是一帮大学的研究人员维护(Xgboost作者华盛顿大学陈天奇在内, 美国, 加拿大, 新加坡的多所高校), 17年被Amazon接盘

4_ Start a project



4_ Start a project

- Deep Learning Toolkits Comparison
- Tensorflow的优势在于社区庞大, 对移动端支持较好, 可以快速地开发一个部署在Android系统上的应用, 缺点是性能比较差, 据说google内部对于内存和GPU优化还是留了后门的, 而且调试困难
 - Facebook同时支持Caffe2和Pytorch两个框架, 意图很明显, Pytorch用于研究, 架构优雅简洁, 方便快速的验证想法, Caffe2主攻工业应用, 性能已经可以和MXNet比肩了, 并且重视移动端的支持
 - MXNet最大的优势在于对内存的优化比较极致, 相同的内存可以训练更大的网络, 但这点估计很快要被Caffe2超越





Toolkit	Modeling Capability	Interface	Model Deployment	Architecture Complexity	Training Speed (VGG-style CNN on CIFAR-10)	Distribution Training
TensorFlow	CNN/RNN/LSTM/MLP	Python/C++/Go/JS	ARM/M	★★★★★	170s	
Caffe2	CNN/RNN/LSTM/MLP	Python/Python/Matlab	ARM/M	★★★★★	140s	
MXNet	CNN/RNN/LSTM/MLP	Python/C++/Java/Shell	ARM/M	★★★★★	140s	
Pytorch	CNN/RNN/LSTM/MLP	Python/Java	ARM/M	★★★★★	180s	

4_ Start a project





还有一些针对特定任务的框架,比如纯C写的darknet,还有Facebook新出的针对目标检测的detectron

- Deep Learning Toolkits Conclusion**
- **TensorFlow** is a safe bet for most projects, Not perfect but has huge community, wide usage.
 - **Pytorch** is best for research. However still new, there can be rough patches
 - Use **TensorFlow** for one graph over many machines
 - Consider **Caffe2** or **TensorFlow** for production deployment
 - Consider **TensorFlow** or **Caffe2** for mobile
1. TensorFlow虽然不完美,但是利于部署,有稳定社群,同时上层API如Keras支持Tensorflow,便于开发
2. PyTorch很适合研究,但是它很新,因此用的时候有很多坑要自己填
3. 除了Tensorflow, Caffe2也可以考虑用于产品部署
4. 手机端可以考虑TensorFlow或Caffe2
- 总结: 框架相当于刻刀, 雕塑的好坏还是要取决于匠人, 框架就是框架, 关键是要出活儿





4_ Start a project

-  Deep learning Toolkits
-  Distribution Learning
-  Deployment on Mobile
-  Issues when Training

4_ Start a project

-  Deep learning Toolkits
-  Distribution Learning
-  Deployment on Mobile
-  Issues when Training

4_ Start a project

-  Deep learning Toolkits
-  Distribution Learning
-  Deployment on Mobile
-  Issues when Training

4_ Start a project

Issues when Training

- **Lack of Training Data** -> Data Augmentation/Transfer Learning
- **Low Accuracy** -> Ensemble learning/Feature stacking
- **Model convergence Issue** -> Gradient Descent Optimization/Activation function
- **Local Optima** -> Weight Initialization/Learning rate auto-tuning
- **Overfitting** -> Regularization



Thanks
