

0: 引言

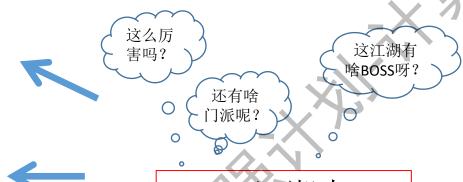






CNN能干啥?

CNN到底是啥?

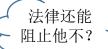


cv江湖中 天下武功出**卷积!**

天下武功出少林?

CNN凭啥这么厉害?

还有啥CNN解 决不了的吗?



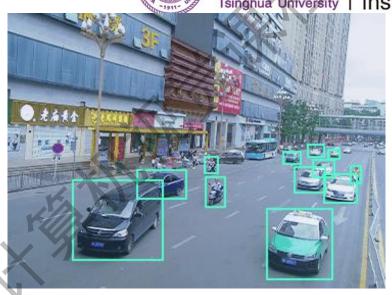
1.1: CNN应用概览



分类 (classification)



Institute for Data Science



探测 (detection)



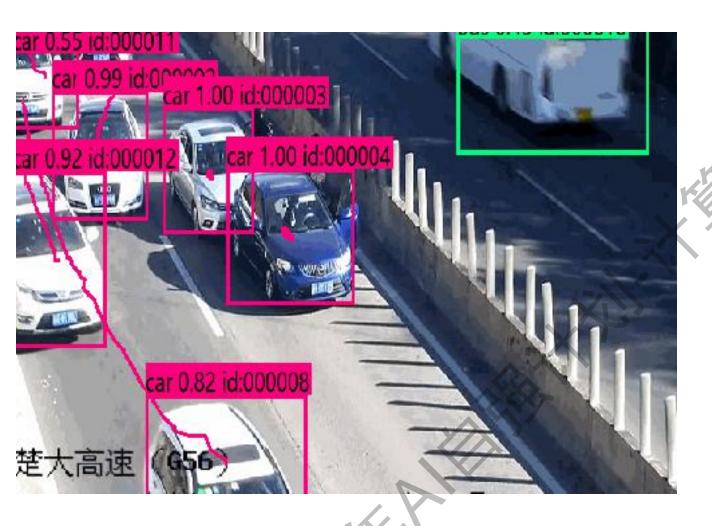
语义分割 (semantic segmentation)



实例分割 (Instance segmentation)

1.1: CNN应用概览







跟踪 (Tracking)

0CR (字符识别)

1.1: CNN应用概览-什么落地最好?

Tsinghua University

Institute for Data Science

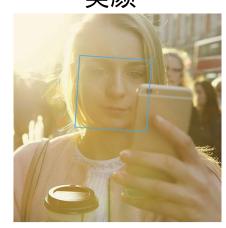




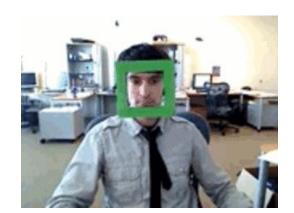


美颜





人脸解锁







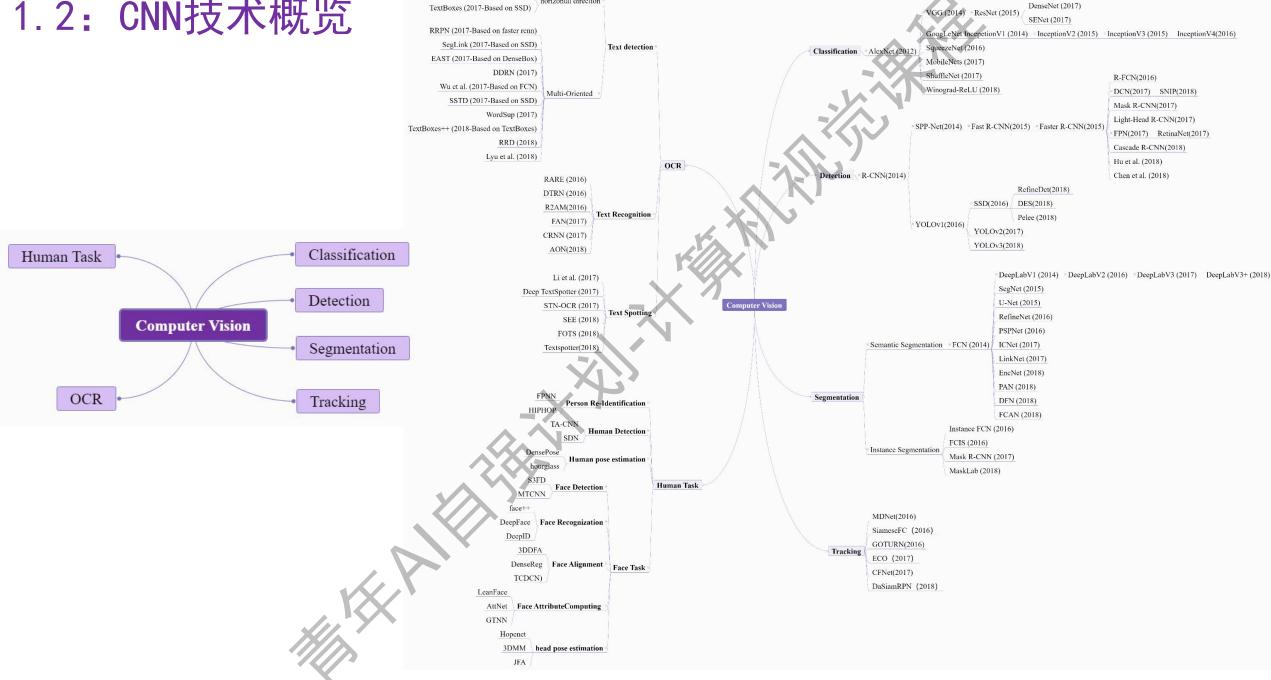






抖音-尬舞机

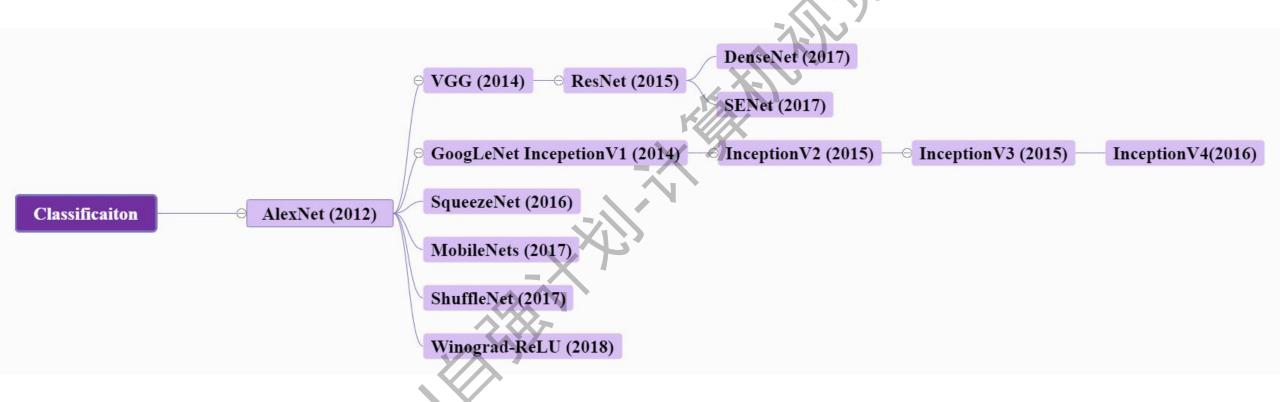
人体姿态识别



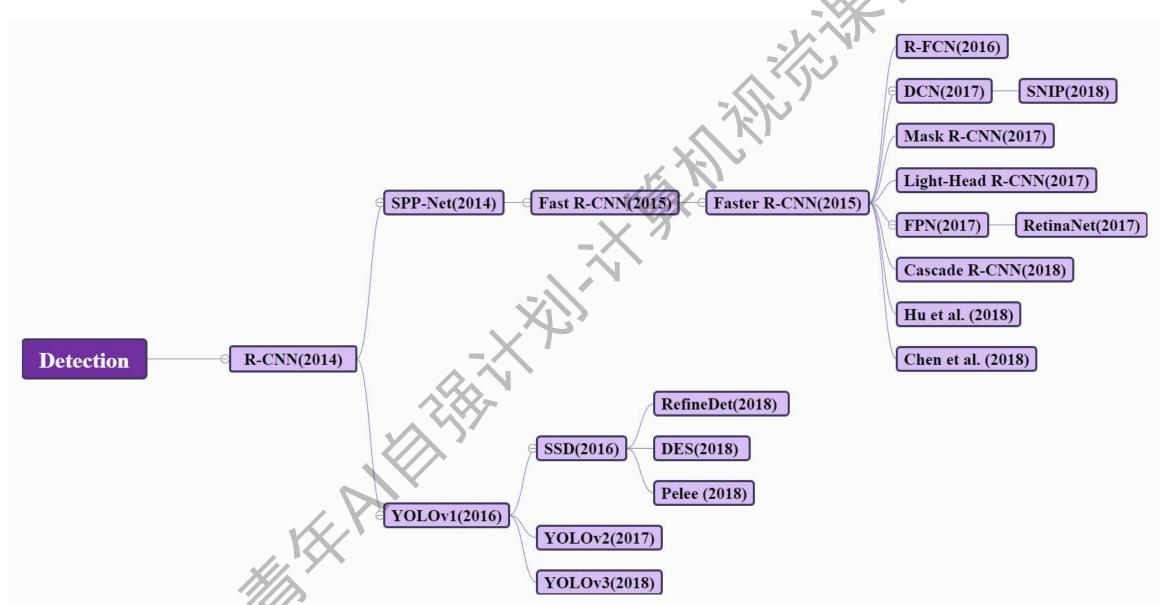
CTPN (2016-CNN+RNN)

horizontal direction





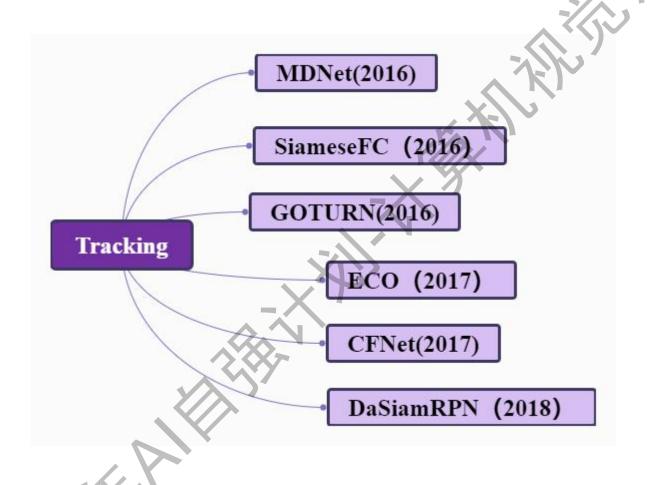




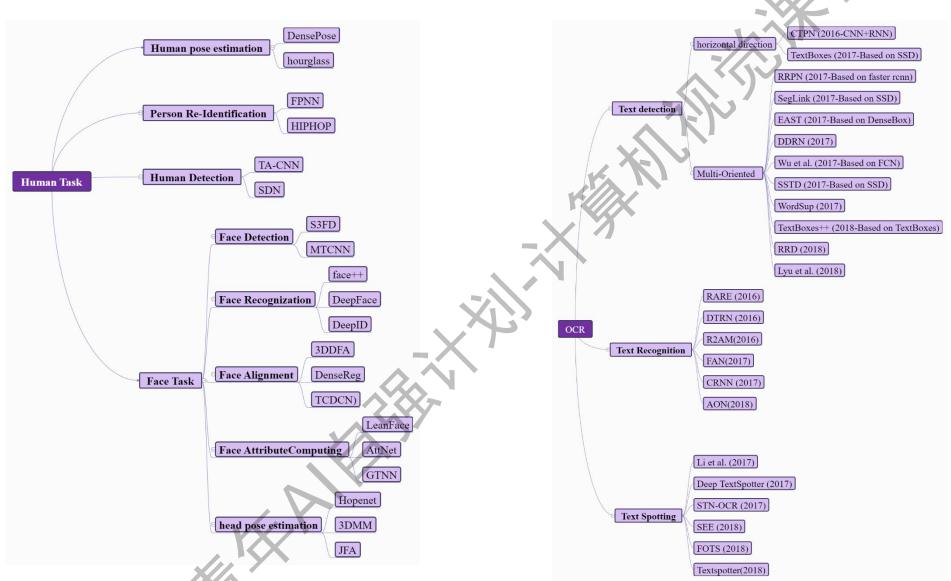














讲一个不太高兴的事让大家高兴一下





情绪的跌宕起伏

来到保安室后



|数据科学研究院

Institute for Data Science

- 1、嫌犯可能出现在这70个小时中的任意时刻
- 2、如果按照1:1的速度浏览,显然耗时太久并不合算
- 3、最大加速比是16倍,加速之后的视频时间是4.3小时
- 4、回顾一下真实可能的肇事场景,肇事时间可能只有 不到半分钟
- 5、在16倍加速的情况下,整个肇事过程只留给我不到2 秒的时间观察,2秒是0.0005小时,只占视频总长的0.01%
- 6、经过试验,人会高度集中注意力到重复性极强的画面 后,会开始出现幻觉[捂脸]



真实监控图像



这根本不是给人干的活, 决定花300块钱解救自己

刚刚的问题 现在的算法无能为力

一张图片|| 一段视频





一段逻辑: 谁 剐蹭 我车了?



把算法当作黑盒



数据科学研究院 Institute for Data Science

标签labels

现在的输出

一个逻辑结果 及证明材料

生活中的通用CV需求

教导主任无话可说



有谁借着出早操的名义 牵手谈恋爱的?



地中间的粑粑是谁拉的?



何时算法能够读懂 通过肢体语言表述的逻辑?

2.1: 标准DNN的CV局限性-参数爆炸



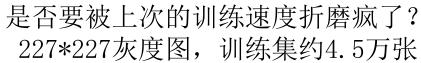
数据科学研究院

Institute for Data Science









上次简化训练情况

输入层约有 227×227 = 51,529个神经元 链接首个隐藏层的权值w个数为: 51,529×1,024 = 52,765,696个 用1个样本更新首层w,需要至少计算 105,531,392次,跑1个epoch需计算

105,531,392×45,000≈4兆7千5百亿次 CPU并行数约为核数,所以训练约需要3-5h



真实的图片要更加复杂 4160*2336<mark>彩色图</mark>,千万像素

A1 = num_units=1024 A2 = num_units=2048 A3 = num_units=2048 A4 = num_units=1024

结合上次作业的 模型结构算一笔账

真实情况

输入层约有

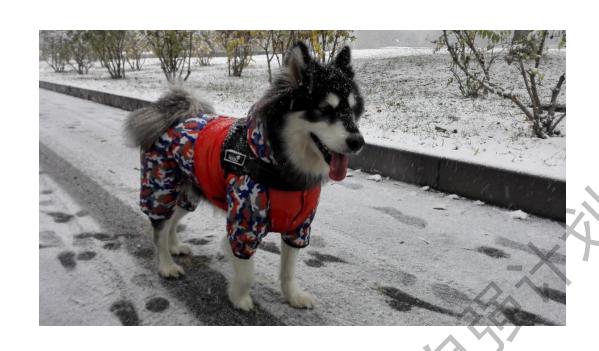
4160×2336×3 = 29,153,280个神经元 链接首个隐藏层的权值个数为: 29,153,280 ×1,024 = 29,852,958,720个 用1个样本更新首层w,需要至少计算 59,705,917,440次,1个epoch需要计算

59,705,917,440 × 45,000 = 2,686,766,284,800,000 约为2700兆次,如果仍用CPU训练大概需要3-5个月

2.1: 标准DNN的CV局限性-鲁棒性差



大家一起来找茬





大家眼中看起来没区别的两张图在计算机看来却又翻天覆地的变化

标准DNN提取特征的方式与人(平移不变性)有本质的差距



积: Conv

化: Pooling

全连接: FC

CNN的三种主要结构

好消息 FC结构 大家已经掌握 先说卷积是啥

1、发生关系的两个变量是啥

2、运算规则是啥

3、运算结果是啥

1、变量是啥?



What We See



What Computers See

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1.,	1	0
0	1	1	0	0

×1	×0	×1
×0	×1	×0
×1	×0	×1

变量1: 输入图像

变量2: 卷积核



3、运算结果是啥?

2、运算规则是啥?

1,	1,0	1,1	0	0
0,0	1,	1,0	1	0
0,,1	0,0	1,1	1	1
0	0	1	1	0
0	1	1	0	0

4	0.00	



Image

Convolved Feature

- 1、用卷积核从左到右,从上到下将需要处理的图像依次扫一遍;
- 2、每次被卷积核扫到的部分,与卷积核对应相乘再相加

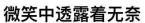
一个新的矩阵 可以被称为Feature map

恭喜大家3分钟学会卷积运算



Institute for Data Science





真的管用?

这玩意 分三个方面为大家解释 卷积如何提取特征

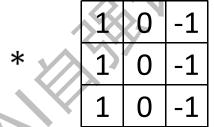
2、与标准DNN相比的优点

3、与传统机器视觉算法相比的优点是啥

怎么感觉高兴不起来?

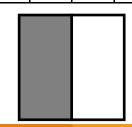
1、卷积如何提取图像中的特征

0,	010	0_1	10	10	104
0,	0	ŌŪ	100	10	10
0,	010	0_0	100	100	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

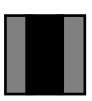


0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0

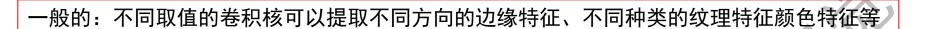
原画面中间的 边界特征被提取出来了













输入图像



卷积核2

这与人类识别物体的原理是相仿的 我们人类也同样是<mark>通过识别物体的边缘来确定其形态进而对其分类的</mark> 比如同样出现在墙上,圆的一般是表,而方的一般是门窗

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	4
	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	9
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	4

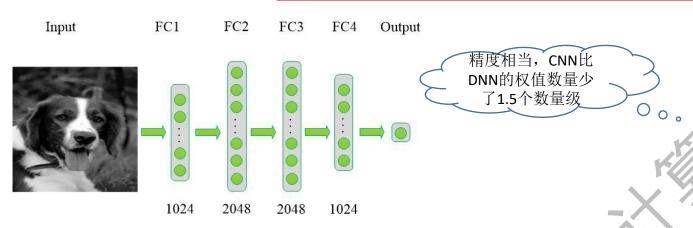




数据科字研究院

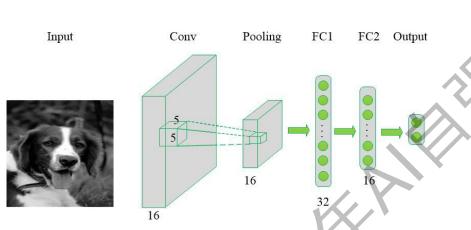
Institute for Data Science

2、卷积与标准DNN相比的优势:一个对比实验



甘土女件边里	对比结果		
基本条件设置	DNN	CNN	
数据集: dot & not dog			
Batch-size: 128			
Lr: 0.1 (learning rate decay)	权值数量: 61,161,473	权值数量: 1,549,810	
Loss function:交叉熵 优化方法:Adam	检测精度: 93%	检测精度: 94%	
Batchnorm: 使用 激活函数: relu			

DNN的网络结构



同等精度下简单CNN的网络结构

结果对比&结论

第一节课小实验, 体会DNN和CNN 的感受野

参数效率高:

随着输入图像尺寸的增大

- a. 标准DNN参数指数型增加
 - b. CNN参数可以线性增加

鲁棒性强:

- a. 标准DNN首层1个参数感受1个像素点 (1列像素,天差地别)
- b. CNN一组参数在整个画幅感受一类特征 (无垂直边缘改变,无差别)



"提取图边缘信息"来做特征表示,并不是CNN首创的,在他之前的传统机器视觉算法, 大多符合这种指导思想。标准DNN曾经就一度败在他们手下

d+ /T /+r/L	-LARIHA N	效果	展示
特征名称	功能描述	特征提取图	原图
LBP	一种常用的纹理描述算子,用于度量和提取图像局部的纹理信息,对光照具有不变性。 LBP 结合 BP 神经网络 已经用于人脸识别等领域。		
HOG	通过计算和统计图像局部区域的梯度方向直方图来构成特征。Hog 特征结合 SVM 分类器已经被广泛应用于图像识别中,尤其在行人检测中获得了极大的成功。具有光照不变性,不具有尺寸和旋转不变性。HOG+SVM 进行行人检测的方法是法国研究人员 Dalal 在 2005 的 CVPR 上提出的。		
Haar	该算法是一种用于目标检测或识别的图像特征描述子,Haar 特征通常和AdaBoost 分类器组合使用,且由于 Haar 特征提取的实时性以及AdaBoost 分类的准确率,使其成为人脸检测以及识别领域较为经典的算法。		
Canny	该算法是目前最常用的边缘提取算法,通过计 算梯度、非极大值抑制、双阈值和边界跟踪,可以 完美的勾勒图像轮廓。		



数据科学研究院

Institute for Data Science

之

前

例子

的卷积核

可还记得NLP 领域的学派大战



CNN的本质是神经网络,同样有FP和BP。卷积核中的权值都是"学"来的

CNN是通过对大量经验数据的学习, 来对物体进行分类的"经验派"

0	w_1	w_2	w_3
	W_4	w_5	w_6
	$\overline{w_7}$	w_8	W ₉

7 0

➤ 传统机器视觉算法是希望通过规则, 描述物体的特征来对其进行分类的"规则派"



2012年的AlexNet开始,CV进入CNN时代

3、CNN与传统机器视觉算法的优势对比

优势一:源自NN的先天优势 CNN的特征提取时自动学习的, 每一个卷积核负责提取一类特征 无论在特征提取的全面性或者鲁棒性上都更优



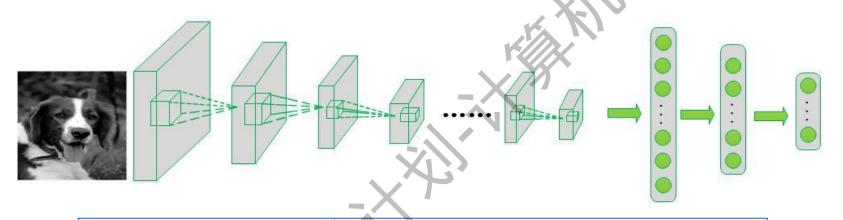
Institute for Data Science

就是意思一下, 大家切勿较真

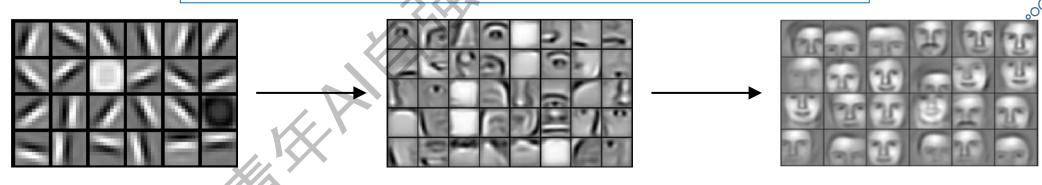
3、CNN与传统机器视觉算法的优势对比

优势二:源自CNN深度网络的先天优势 CNN的网络层数可以任意加深

Conv Pooling FC1 Input Conv Pooling Conv Pooling FC2 Output



以人脸识别为例, 每层卷积负责提取不同抽象等级的特征



临近输入层

中间层

临近输出层

2.3: 池化 (pooling) 略解

了 Tsinghua University

数据科学研究院

Institute for Data Science

深思熟虑

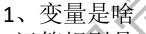


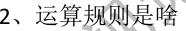




Max

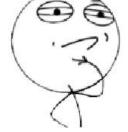
pooling



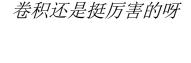


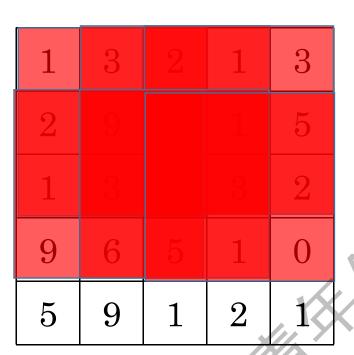
3、运算结果是啥





就这么扫一下能管用?





1. 变量:输入的 feature-map

2.运算规则: 扫一遍挑出最大的

没有参数需要更新,很满意吧?

995955

1. 变量: 池化核

3.结果:输出新的feature-map

2.3: 池化 (pooling) 略解



池化的作用1:减少参数总量 之前叫降采样(down sampling)



这么随便就把参数丢了吗?

Pooling 类型	描述	示意图
Average pooling	对邻域内的特征点求平均。	31 15 28 184 0 100 70 38 12 12 7 2 12 12 45 6
Sum pooling	对邻域内的特征点求和 (有三种不同版本)	
Stochastic pooling	对 feature map 邻域内的元素按照其概率值大小随机选择,即元素值大的被选中的概率也大。而不像 max-pooling 那样,永远只取那个最大值元素。	0 1.1 2.5 0.9 2.0 1.0 0 1.5 1.0
S3Pool	首先采用最大池化得到激活值,然后随机行和列按照 grid 区域进行下采样。	12 1 5 1 max pool 2 2 falter viside p 7 6 8 8 0 Downsampling 6 8 (c) S3Pool, pooling window $k = 2$, stride $s = 2$, grid size $g = 2$
Lp-Pooling	使用高斯核进行池化,具体操作如公式所示。O 表示池化输出,I 表示 feature map,G 表示高斯核。 $O = (\sum\sum I(i,j)^P \times G(i,j))^{1/P}$ $P=1$ 表示简单的高斯平均。	kout and constituent of the state of the sta



结果

除了一些特殊情况外 普遍采用的都是Max pooling

许多算法科学家与大家有同样的想法

2.3: 池化 (pooling) 略解



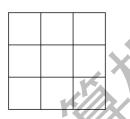


Institute for Data Science

池化的作用2: 平移不变性---小实验

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
9	6	5	1	0
5	9	1	2	1

Max pooling



- 7				
1	9	9	5	
J	9	9	5	
	9	9	5	

5	9		2	1
1				3
2				5
1	3		3	2
9	6	5	1	0

1.变量: 向下平移了1行

Max pooling

-	

9 9 5 9 5 9 9 5 9

3.结果: 并没有改变

2.3: 池化(pooling)略解



Institute for Data Science

池化的作用2: 平移不变性---结论



这跟人类处理图像的原理再一次不谋而合了,

当我们关注的目标物体在我们的视野中发生位移的时候,我们不会因此懵逼。

例如你看着一只羊在撒欢,他在你的视野中只是从右边跑到左边,你肯定不会因此而怀疑他是不是原来的那只羊 背后的原因是不管他在你的视野中平移与否,你对他的特征提取是完全一致的

maxpooling在一定程度上实现了这一点

3.1: 数据集简介



本次作业数据集介绍



为了能让大家不浪费太多的时间 在等待训练上,我们这次的作业 还是选择一个入门级的数据集。

The CIFAR-10 dataset

共60000幅图像 图像尺寸: 32 x 32 RGB

10类,每类6000幅图像。

训练集: 50000幅图像

测试集: 10000幅图像

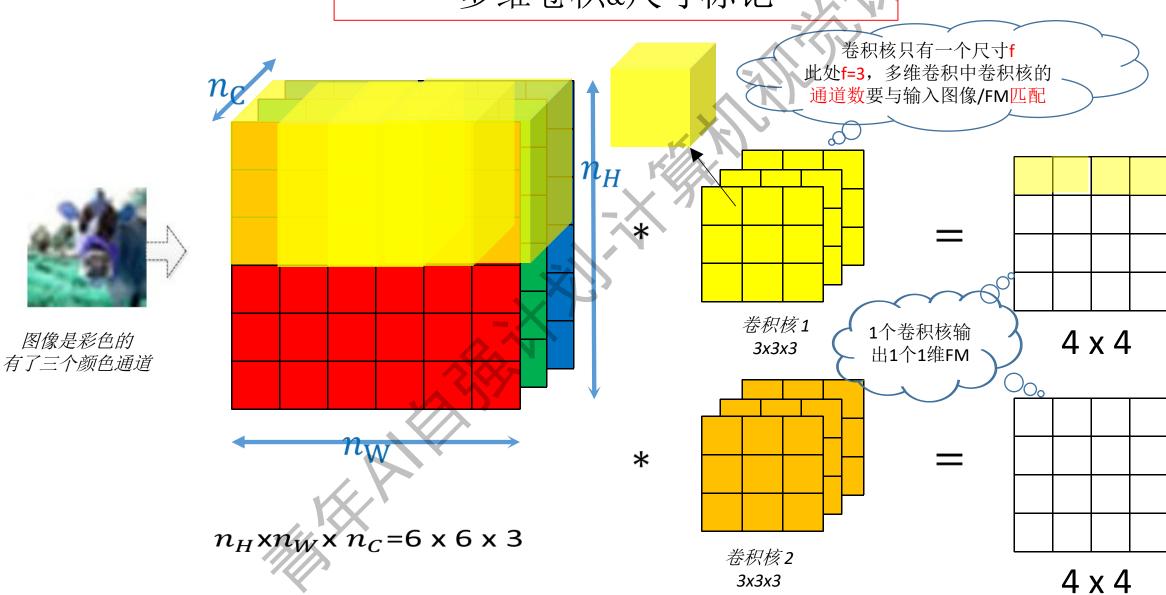


还有一些小问题

- 1. 输入图像是三维的(RGB), 我们只学过平面卷积呀
- 2. 输入图像和卷积核的尺寸怎么标记呢?
- 3. 每扫一次必须只走一格吗? 走两格行不行?
- 4. 你会发现这个"扫一遍"的过程对于角落里的元素不太公平, 角落里的元素只被扫到了一次, 但在画幅中心处的元素却被扫到了好几次。
- 5. 刚刚说到CNN会有好多隐藏层,那么不同隐藏层之间的 卷积核怎么区分标记?



多维卷积&尺寸标记





数据科字研究院 Institute for Data Science

stride

1,	1,0	1,	0	0
0,0	1,	1,0	1	0
0,1	0,0	1,	1	1
0	0	1	1	0
0	1	1	0	0

Image

Convolved Feature 每扫一次走过的像素数量用s表示 S越大,运算后得到的FM越小,假设卷 积核尺寸为f*f,输入图像尺寸为n*n 则得到FM的尺寸为:

$$\left[\frac{n-f}{s}+1\right]*\left[\frac{n-f}{s}+1\right]$$

在做卷积/池化的时候,一次只能走一步吗?

[]为向下取整符号



数据科学研究院

Institute for Data Science

padding

角落里的元素只被扫到一次不太公平? 越靠近边界,被扫到(特征表示)的几率越小

1,	1,0	1,	0	0
0,0	1,	1,0	1	0
0,1	0,0	1,,1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3		
		- (3) (3)	
3 3	3)	23	

Convolved Feature

0	0	0	0	0	0	0
0		4//	X			0
0	1					0
0						0
0						0
0						0
0	0	0	0	0	0	0

在SAME padding且s=1时, $p=\frac{f-1}{2}$,所以f必为奇数

→

填补的主要策略有2种:

SAME: 保持FM不缩小

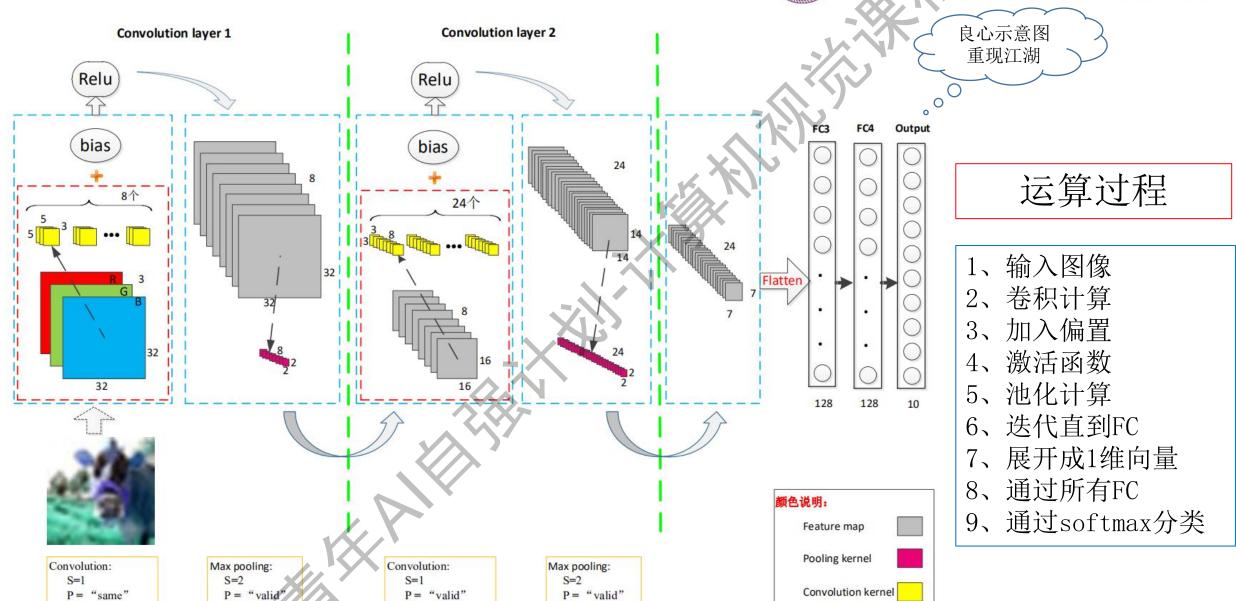
VALID: p=0.

加入padding后,输出FM尺寸变为 $\left[\frac{n + 2p - f}{s} + 1 \right] * \left[\frac{n + 2p - f}{s} + 1 \right]$

在边界处填补(padding)一些像素块 边界向外拓展的像素个数用p表示,此处p=1

3.3: CNN前馈传播结构





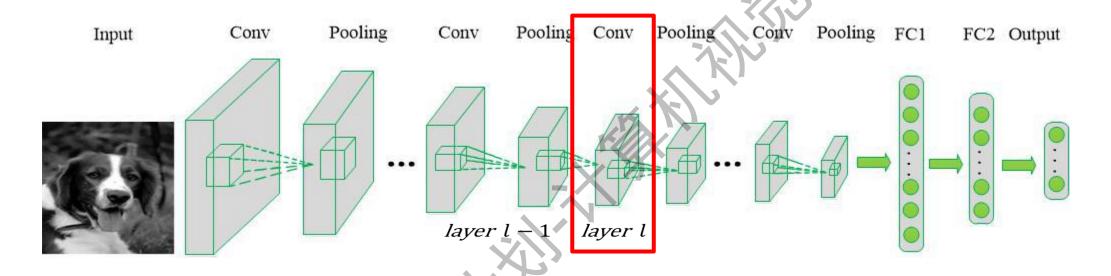
3.3: CNN前馈传播结构



数据科字研究院 Institute for Data Science

Institute for Data Science

Notation汇总-以第1层卷积为例



Input FM size (for layer l): $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$ $f^{[l]} = 第 l$ 层卷积核尺寸值 $p^{[l]} = 第 l$ 层的padding取值 $n_C^{[l-1]} = 第 l$ 层输入FM的通道数 同时也是第l 层卷积核的通道数 第l 层卷积核*权值总数*: $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$

Output FM size (for layer
$$l$$
): $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

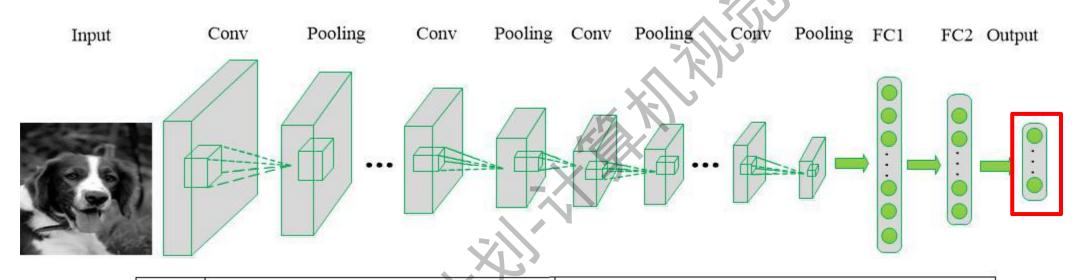
Bias size $b^{[l]} = 1 \times 1 \times n_C^{[l]}$
 $a^{[l]}$ size $= n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$
 $z^{[l]}$ 经过激活函数后变为 $a^{[l]}$
 $n_H^{[l]} = \left[\frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right]$

3.4: CNN反向传播及softmax



Institute for Data Science

多分类任务,如何做预测?



	sigmoid —	softmax
函数	$f(z) = \frac{1}{1 + \exp(-z)}.$	$S_i = rac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$ 其中, z_i 表示第 i 个输入值, S_i 表示第 i 个输出值, N 为类别数
描述	sigmoid 将一个值 <mark>映射</mark> 到(0,1)的区间 这样可以用来做 <mark>二分类</mark> 。	softmax 把一个 k 维向量(a1,a2,a3,a4…) 映射 成另一个 K 维向量(b1,b2,b3,b4…),其中 bi 是一个 0~1 的常数,并且 b 向量的和为 1。然后可以根据 bi 的大小来 进行多分类 的任务。

3.4: CNN反向传播及softmax



Institute for Data Science

取自真实例子 输入可以是任何数

所有输出的和为1 多分类与sigmiod相比的优势

Softmax输入	[-1.467, -0.309, 0.423, -0.884, 0.131, 0.942, 0.806, 1.097, -2.917, 0.589]				
Softmax funtion	$S_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$ 其中, z_i 表示第i个输入值, S_i 表示第i个输出值,N为类别数,这里是 10				
Softmax输出	[0.017, 0.054, 0.111, 0.030, 0.083, 0.187, 0.163, 0.219, 0.004, 0.132]				
Label	[0 1 0 0 0 0 0 0 0]				
Loss function	$loss = -\sum_{i=1}^{N} label_i * log(S_i)$ 其中, $label_i$ 表示标签中的第i个数, S_i 表示softmax输出的第i个数 N表示类别数量,这里为10				
loss	[2.918]				
	要求大家定性掌BP比较复杂				

带入Label会发现 只剩一项了

握即可

不要求掌握

3.5: 训练结果



Institute for Data Science

基本条件设置。	对比结果。		
至 本家什以且。	DNN @	CNN	
数据集: CIFAR-10。			
数据增强: 左右翻转。			
Batch-size: 256			
Lr: 0.1 (learning rate decay)			
Loss function: 交叉熵。	权值数量: 11,550,730。	权值数量: 170,818。	
优化方法: Adam	检测精度: 56.25%。	检测精度: 70.18%。	
训练次数: 5000-			
模型层数: 4。			
Batchnorm: 使用。	XY		
激活函数: relu。			

精度不是太 高?

标准DNN的 局限性

什么是 卷积

卷积的优势

什么是 池化

多维卷积与 池化

实例练兵

3.6: 尾声&作业说明



Institute for Data Science

第一节课时的两个小实验





















人类是通过对样本的学习来获取知识的

- 2、人类可以自动提取特征并进行分类
- 3、人类可以通过调整注意力来调整视角

CNN已经具备的





- 1、人类可以通过小样本学习知识
- 2、人类学习之后的泛化能力极强
- 3、人类具备联想、抽象、举一反三的能力

CNN尚且缺少的

个是虎皮猫,哪个是小老虎呢?

8.2: 尾声&作业说明



作业说明: 本次只有1个作业 用简单的CNN拟合CIFAR100数据集 帮助大家彻底搞懂CNN的前馈传播结构







扫码加好友进群

关注直播间公告