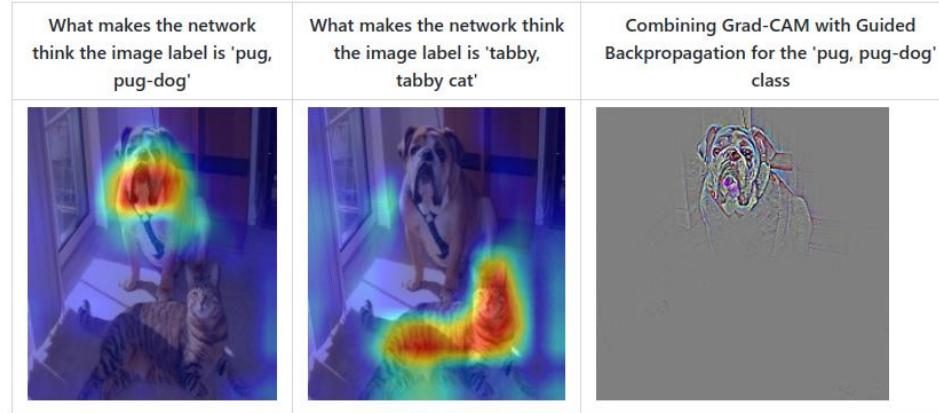


Grad-CAM

Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization

ICCV'17

Ramprasaath R. Selvaraju · Michael Cogswell · Abhishek Das · Ramakrishna Vedantam · Devi Parikh · Dhruv Batra



力场梯度映射
CAM

(Class Activation Mapping)

梯度网络,加权梯,透明度更高
Grad-CAM

(Gradient-weighted Class Activation Mapping)

论文下载地址: <https://arxiv.org/abs/1610.02391>

推荐博文: https://blog.csdn.net/qq_37541097/article/details/123089851

推荐代码 (Pytorch) : <https://github.com/jacobgil/pytorch-grad-cam>

Grad-CAM

6.3 Identifying bias in dataset

做一个医生二分类器，
Input Image



Ground-Truth Nurse

偏置很大的模型
Grad-CAM for
Biased model



Predicted Nurse

→训练的数据集中大部分为女性。

无偏置模型
Grad-CAM for
Unbiased model



Predicted Nurse



Ground-Truth Doctor



Predicted Nurse

预测偏置很大的模型是通过人脸进行的。根据服装、使用器械来区分

→训练的数据集中大部分为女性。

Grad-CAM

- ① 输入图像，经过 CNN 提取特征 A，再接下流子网。
- ② Grad-CAM：输入图像高感兴趣类别，图像经过前向传播得到该类分类概率。
把其它类别的幅度都置 0。只保留感兴趣类别置为 1。之后反向传播到特征 A(幅度 A')
③ 将 A' 乘以值，与 A 加权求和，通过 ReLU 即可得到 heatmap 意思是拿到网络输出的那个
类的值，放缩
生效。

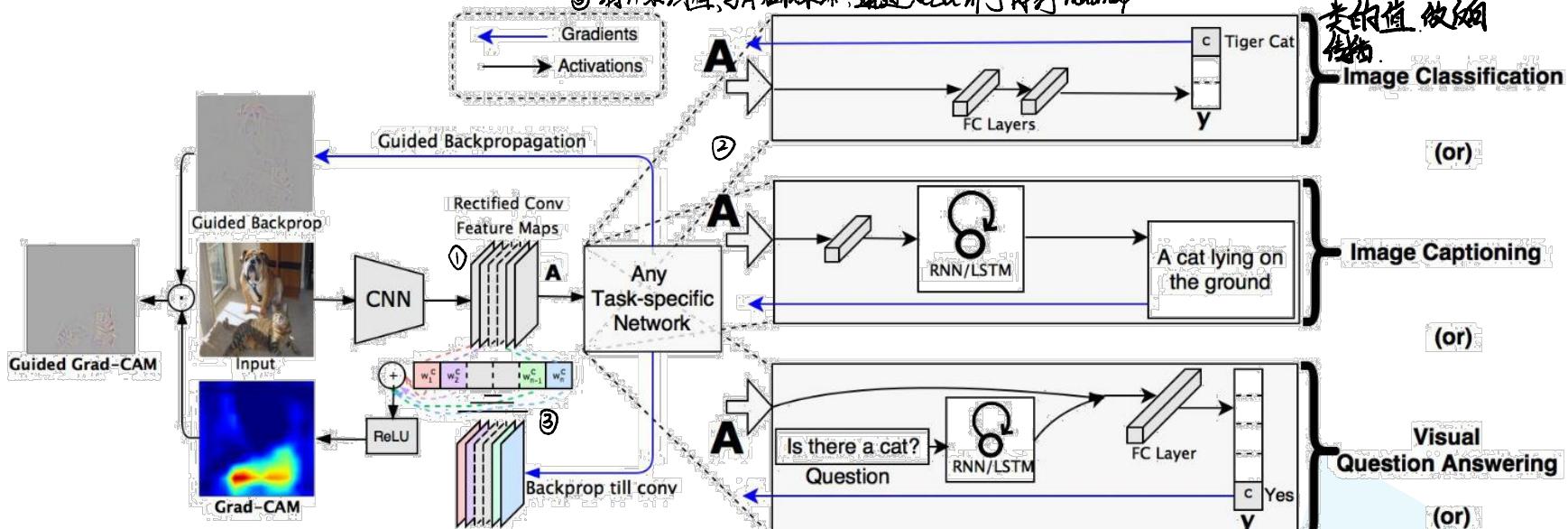


Fig. 2: Grad-CAM overview: Given an image and a class of interest (e.g., ‘tiger cat’ or any other type of differentiable output) as input, we forward propagate the image through the CNN part of the model and then through task-specific computations to obtain a raw score for the category. The gradients are set to zero for all classes except the desired class (tiger cat), which is set to 1. This signal is then backpropagated to the rectified convolutional feature maps of interest, which we combine to compute the coarse Grad-CAM localization (blue heatmap) which represents where the model has to look to make the particular decision. Finally, we pointwise multiply the heatmap with guided backpropagation to get Guided Grad-CAM visualizations which are both high-resolution and concept-specific.

Grad-CAM

最好的可视化结果就是最深那个卷积层。

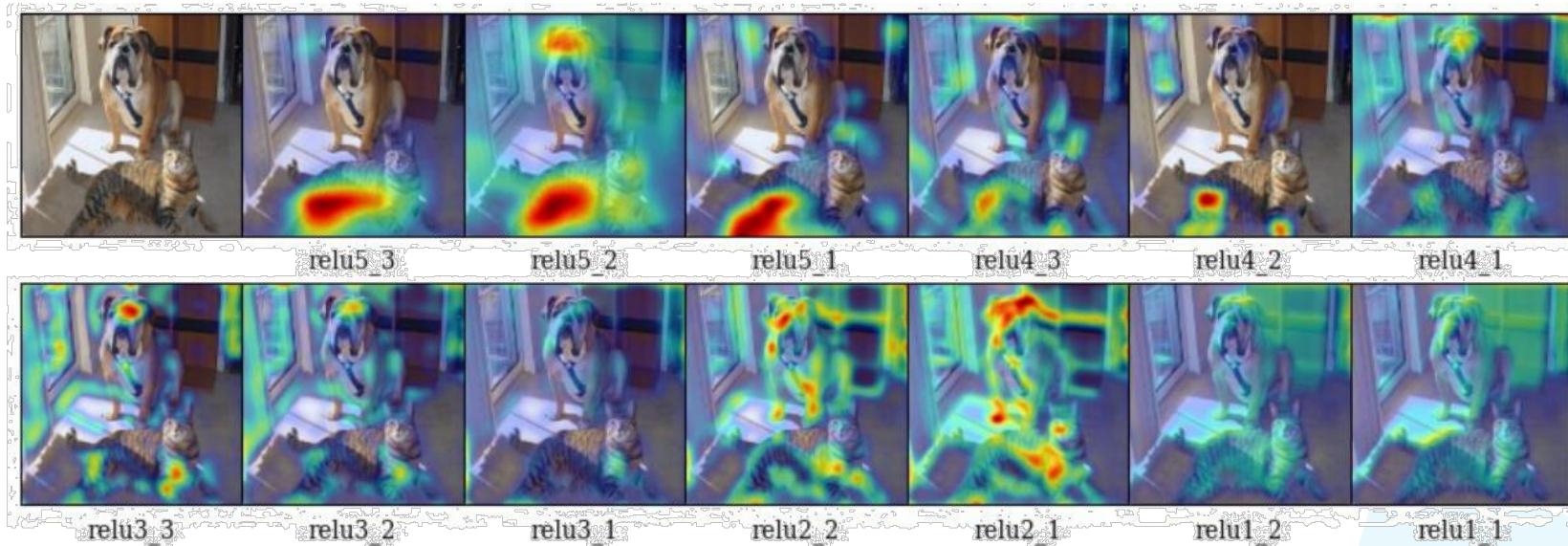


Fig. 13: Grad-CAM at different convolutional layers for the ‘tiger cat’ class. This figure analyzes how localizations change qualitatively as we perform Grad-CAM with respect to different feature maps in a CNN (VGG16 [52]). We find that the best looking visualizations are often obtained after the deepest convolutional layer in the network, and localizations get progressively worse at shallower layers. This is consistent with our intuition described in Section 3 of main paper, that deeper convolutional layer capture more semantic concepts.

Grad-CAM

$$L_{\text{Grad-CAM}}^c = \text{ReLU}\left(\sum_k \alpha_k^c A^k\right) \quad (1)$$

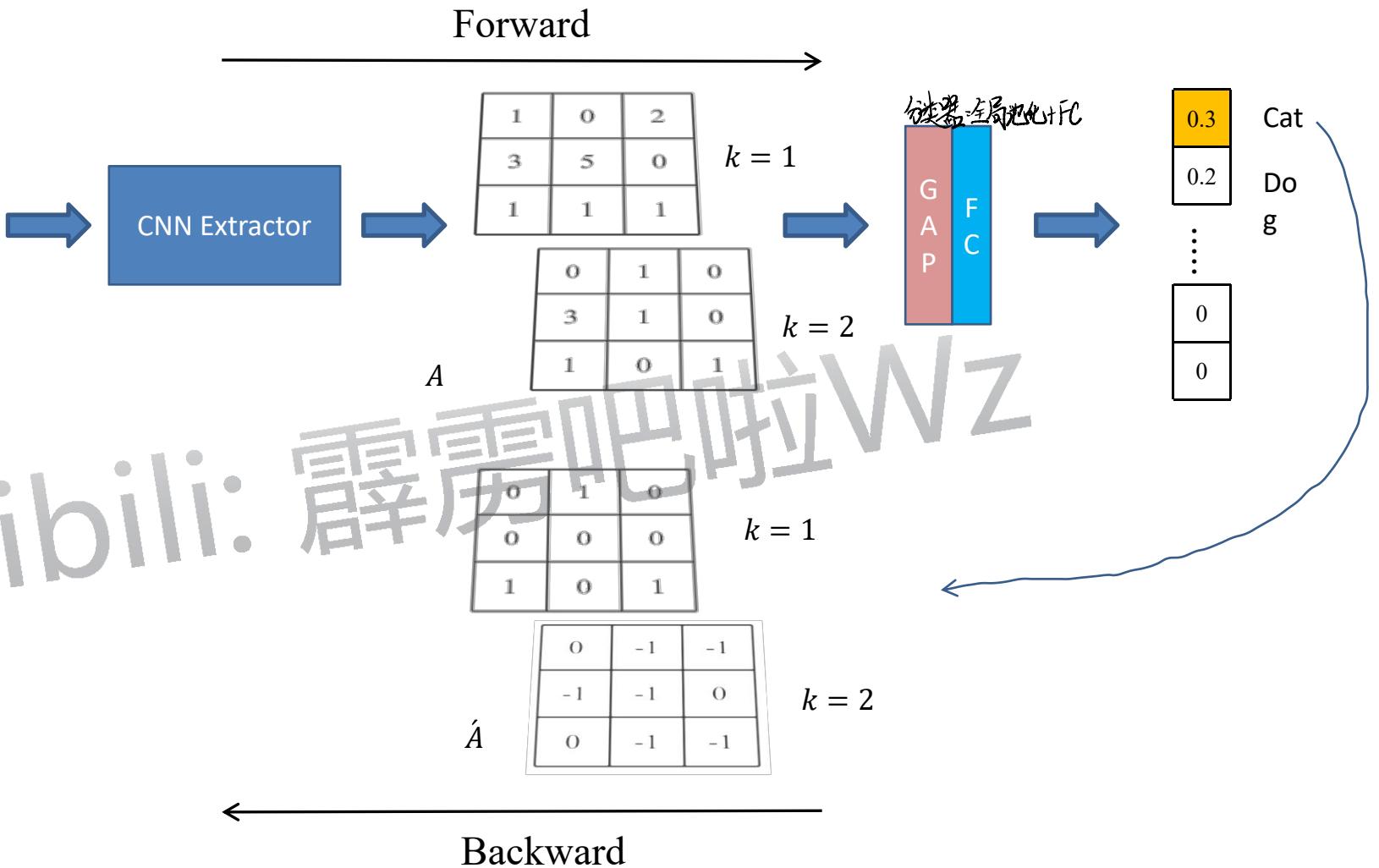
特征图与权重加权求和，并通过ReLU

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k} \quad (2)$$

求平均

对数据求导

- A 代表某个特征层（一般取最后一个卷积层的输出）
- k 代表特征层 A 中第 k 个通道（channel）
- c 代表关注的类别 c
- A^k 代表特征层 A 中第 k 个通道的数据
- α_k^c 代表对于类别 c ，特征层 A 第 k 个通道的权重
- y^c 代表网络针对类别 c 预测的分数(score)，注意这里没有通过 softmax 激活
- A_{ij}^k 代表特征层 A 在通道 k 中，坐标为 ij 位置处的数据
- Z 代表特征层的宽度乘以高度



Grad-CAM

1	0	2
3	5	0
1	1	1

$k = 1$

0	1	0
3	1	0
1	0	1

$k = 2$

A

0	1	0
0	0	0
1	0	1

$k = 1$

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

$k = 2$

\hat{A}

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k} \quad (2)$$

$$\alpha^{\text{Cat}} = \begin{pmatrix} \alpha_1^{\text{Cat}} \\ \alpha_2^{\text{Cat}} \end{pmatrix} = \begin{pmatrix} \frac{1}{3} \\ -\frac{2}{3} \end{pmatrix}$$

$\nearrow \frac{3}{9}$
 $\searrow \frac{-6}{9}$

Bilibili: 霹雳吧啦Wz

Grad-CAM

1	0	2
3	5	0
1	1	1

$k = 1$

0	1	0
3	1	0
1	0	1

$k = 2$

A

0	1	0
0	0	0
1	0	1

$k = 1$

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

$k = 2$

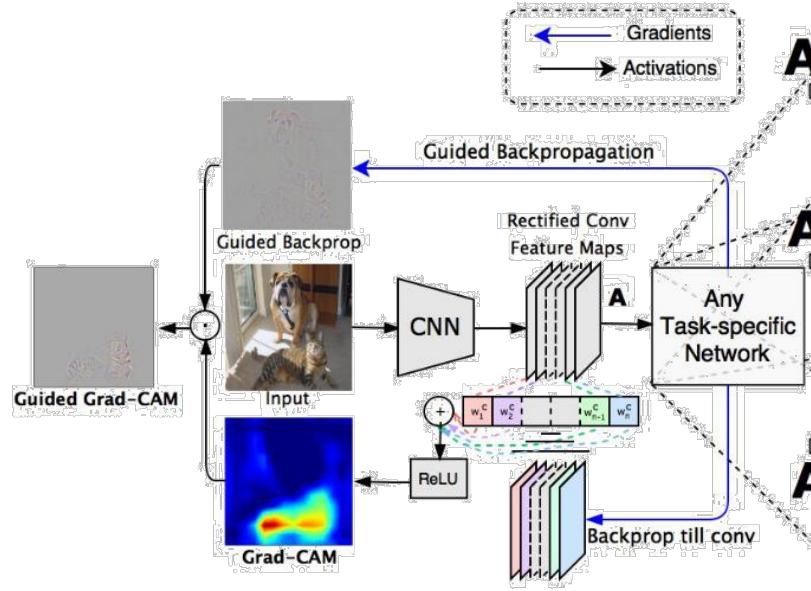
\hat{A}

$$L_{\text{Grad-CAM}}^c = \text{ReLU}(\sum_k \alpha_k^c A^k) \quad (1)$$

$$\alpha^{\text{Cat}} = \begin{pmatrix} \alpha_1^{\text{Cat}} \\ \alpha_2^{\text{Cat}} \end{pmatrix} = \begin{pmatrix} \frac{1}{3} \\ -\frac{2}{3} \end{pmatrix}$$

$$\begin{aligned} L_{\text{Grad-CAM}}^{\text{Cat}} &= \text{ReLU}\left(\frac{1}{3} \cdot \begin{pmatrix} 1 & 0 & 2 \\ 3 & 5 & 0 \\ 1 & 1 & 1 \end{pmatrix} + \left(-\frac{2}{3}\right) \cdot \begin{pmatrix} 0 & 1 & 0 \\ 3 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}\right) \\ &= \text{ReLU}\left(\begin{pmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ -1 & 1 & 0 \\ -\frac{1}{3} & \frac{1}{3} & -\frac{1}{3} \end{pmatrix}\right) \\ &= \begin{pmatrix} \frac{1}{3} & 0 & \frac{2}{3} \\ 0 & 1 & 0 \\ 0 & \frac{1}{3} & 0 \end{pmatrix} \end{aligned}$$

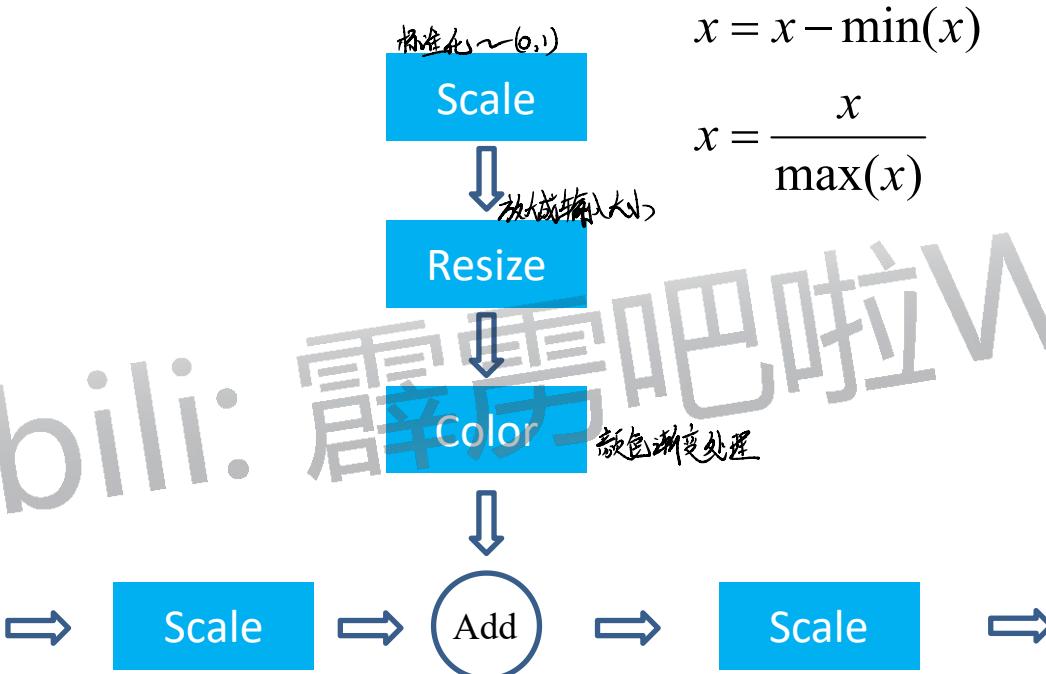
Grad-CAM



热力图如何与原图结合？后处理

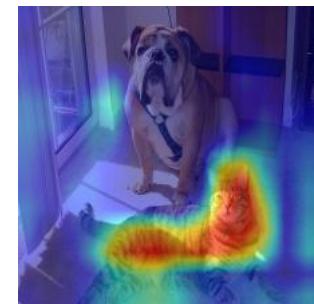
Grad-CAM

Bilibili: 露露吧啦Wz



$$img = \frac{img}{255}$$

$$cam = \frac{cam}{cam_{max}} \cdot 255$$



为什么拿到那个类，做反向传播，而不用交叉熵做反向传播？

①用交叉熵之后，这时计算的梯度包括的意思：输入对损失函数的重要性。

而直接对网络输出值的那个维度值进行求导，得到的是输入对于类别的重要性。

通过 Grad-CAM 与 Gradient ascent 可视化 CNN 图，梯度是啥？

原来说梯度下降或梯度上升，都是为了最优化问题，举例：我通过提高网络输出时，目标类的输出值，也可以优化交叉熵。
但这两个不同目标的梯度，绝不一样，也绝不相似。

总结：即使两个问题等价，但不同的目标函数对应的梯度含义不同。