# Part 1 : Residual Networks



- Read paper : https://arxiv.org/pdf/1512.03385
- Focus on **Fig 3.** and **Sec 3.4**

# Dataset

Butterfly and Moth species classification into 100 classes

# Part 2 : Free Lunch



$P(c \mid I) = 0.95$

Raw Importance

Refined Mask

- Goal : Given a trained CNN, segment the objects without additional data.

# Saliency Visualization

- Given an image $I_0$, a class $c$ and a CNN with a class score function $S_c(I)$ we'd like to rank the pixels of $I_0$ based on their influence on the score $S_c(I_0)$

- Considering a linear score model for some class :
$$S_c(I) = w_c^T I + b_c$$
  - Image $I$ is represented in 1-D form and $w_c$ and $b_c$ are respectively the weight and bias vector of the model.
  - Magnitude of elements of $w$ defines the importance of corresponding pixels of $I$ for the class $c$.

- In CNNs the class $S_c(I)$ is highly non-linear, However we can approximate it with a linear function about $I_0$ using 1ˢᵗ order Taylor expansion:

$$S_c(w) \approx w^T I + b$$
$$w = \frac{\partial S_c}{\partial I}$$

# Revisiting



$I$

CNN

FC

$P(c \mid I) = 0.95$

$S_c$

Raw Importance

$$\frac{\partial S_c}{\partial I}$$

Refined Mask

CV magic :
- Blurring
- OTSU Thresholding
- Dilate/Erode

- Sec 3.2 in [1]

[1] Simonyan et al, *Deep Inside Convolutional Networks : Visualizing Image Classification Models and Saliency Maps, 2013*
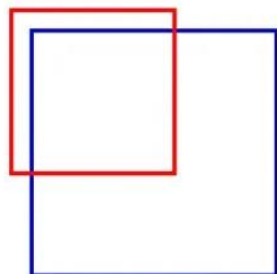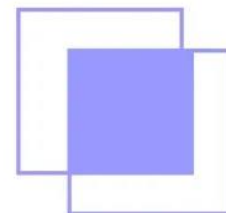
# Alternative : Grad-CAM

- Uses the gradient of a 'target concept' flowing into the final convolutional layer to produce a *coarse localization* map highlighting the important regions in the image for predicting the concept.



$$\alpha_k^c = \frac{1}{Z}\sum_i \sum_j \frac{\partial S_c}{\partial A_{ij}^k}$$

$$L_{grad-cam} = ReLU(\sum_k \alpha_k A^k)$$

$P(c \mid I) = 0.95$

$S_c$

CNN

$A_k$

Feature map

FC

$I$

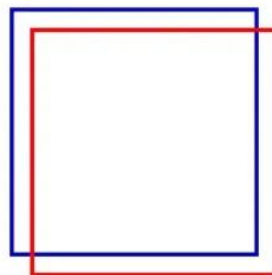# Evaluation: IoU Score



$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

Poor          Good          Excellent

# References

- Simonyan et al, *Deep Inside Convolutional Networks : Visualizing Image Classification Models and Saliency Maps, 2013*

- Selvaraju et al, *Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, 2016*