

# Integrated Gradients

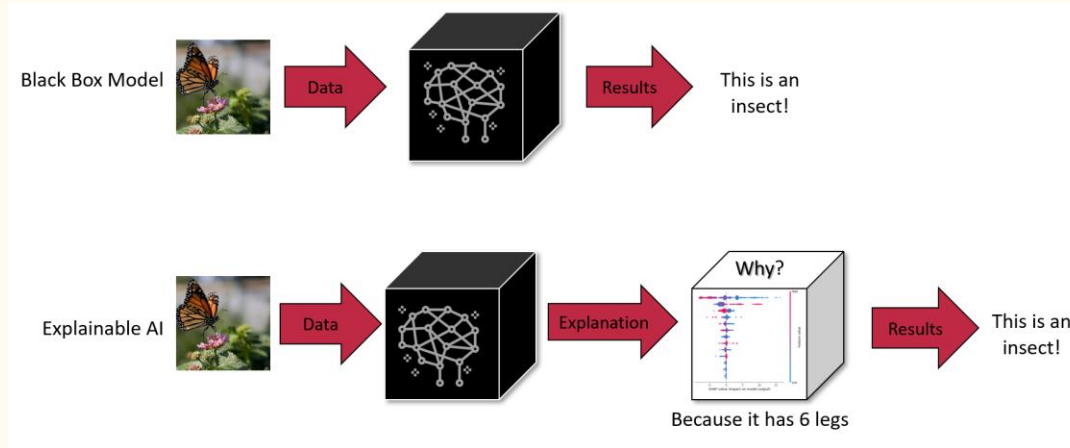
---

Henrique Jesus M12359  
Prof. Dr. Hugo Proença

Visão Computacional

# Introduction

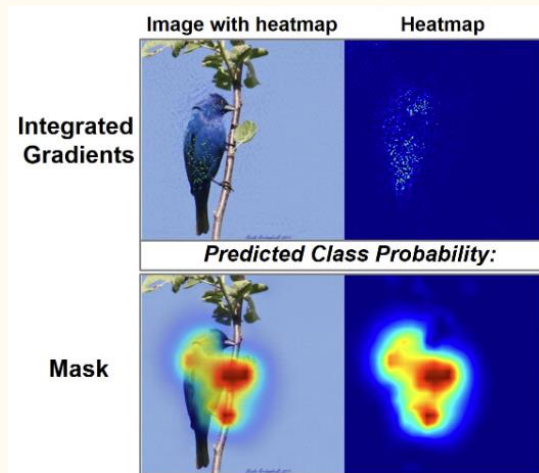
- Explainable AI (XAI)
- Black Box



Source: <https://images.squarespace-cdn.com/content/v1/5bce4071ab1a620db382773e/1626294924656-L8ZFLDNLKYK7C33DHX4I9/Picture1.png>

# How to make a model explainable?

- Paper "Axiomatic Attribution for Deep Networks"
- Integrated Gradients



Source: <https://media.arxiv-vanity.com/render-output/7235504/IntroductionFinal2.jpg>

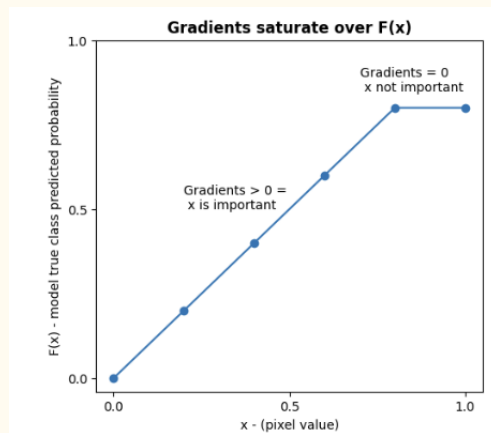
# Gradients

- The direction of steepest increase of the function at that point.
- Determine the rate of change of a function in relation to its variables.

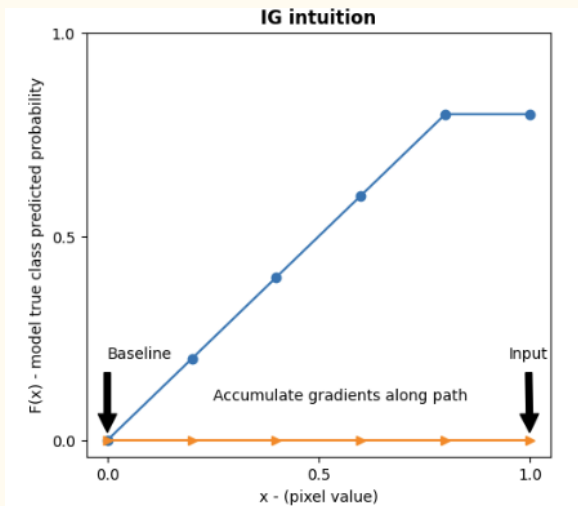
$$f(x, y) = x^2 + y^2$$

$$\frac{\partial f}{\partial x} = 2x \quad \frac{\partial f}{\partial y} = 2y$$

$$\text{grad}(f(1, 2)) = (2(1), 2(2)) = (2, 4)$$



# Integrated Gradients



$$\text{IntegratedGradients}_i(x) ::= (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha(x - x'))}{\partial x_i} d\alpha$$

where:

$i$  = feature (individual pixel)

$x$  = input (image tensor)

$x'$  = baseline (image tensor)

$\alpha$  = interpolation constant

$$\text{IntegratedGrads}_i^{\text{approx}}(x) ::= (x_i - x'_i) \times \sum_{k=1}^m \frac{\partial F(x' + \frac{k}{m}(x - x'))}{\partial x_i} \times \frac{1}{m}$$

where:

$k$  = interpolation constant step

$m$  = number of steps in the Riemann sum approximation of the integral

$$\text{IntegratedGrads}_i^{\text{approx}}(x) ::= \overbrace{(x_i - x'_i)}^{5.} \times \sum_{k=1}^m \overbrace{\frac{\partial F(x' + \frac{k}{m}(x - x'))}{\partial x_i}}^{\overbrace{2.}^{\overbrace{1.}^3}} \times \overbrace{\frac{1}{m}}^{4.}$$

1. Generate alphas  $\alpha$
2. Generate interpolated images  $= (x' + \frac{k}{m} \times (x - x'))$
3. Compute gradients between model  $F$  output predictions with respect to input features  $= \frac{\partial F(\text{interpolated path inputs})}{\partial x_i}$
4. Integral approximation through averaging gradients  $= \sum_{k=1}^m \text{gradients} \times \frac{1}{m}$
5. Scale integrated gradients with respect to original image  $= (x_i - x'_i) \times \text{integrated gradients}$ .