

Combining Similarity and Adversarial Learning to Generate Visual Explanation: Application to Medical Image Classification

Martin Charachon^{1,2}, C eline Hudelot², Paul-Henry Courn ede², Camille Ruppli¹, Roberto Ardon¹
¹Incepto Medical
²Universit  Paris-Saclay, CentraleSup elec, MICS

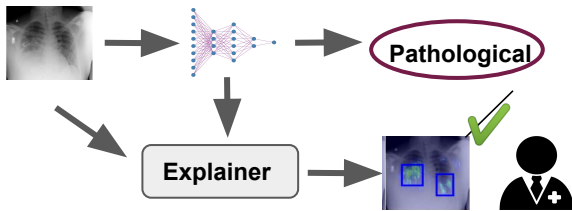
Introduction

Problem: Explain decision of black-box classifiers



Objectives:

What are the **discriminative regions** in the image for the classifier? Are they **interpretable** for humans (clinicians)?



Prior Work

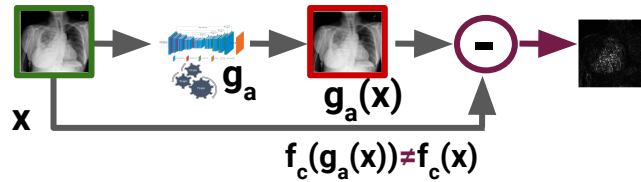
- Saliency maps [1], Activation maps based [2]
- **Perturbation-based** [3, 4, 5]

Contributions

- New definition of **visual explanation** through **adversarial example generation**
- New Optimization workflow combining the **training** of an **adversarial generator** and a **similar generator**
- New **regularization** methods **improving** several explanation methods in terms of **weak localization**

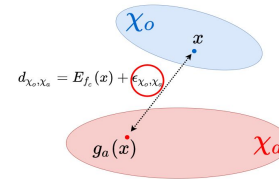
"Naive" Method

Visual explanation: $E_{f_c}(x) = |x - g_a(x)|$



Issues:

- Different space χ_o and χ_a
- Reconstruction errors

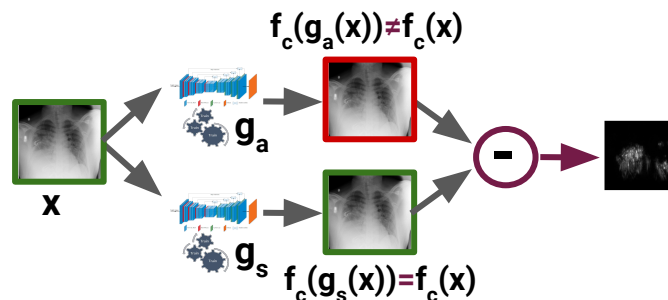
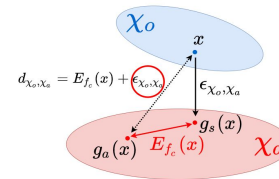


Proposed Method

Approach: Generate an adversarial example $g_a(x) \in \chi_a$
Project x in space $\chi_a \rightarrow g_s(x)$

Visual explanation:

$$E_{f_c}(x) = |g_s(x) - g_a(x)|$$

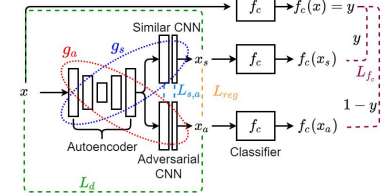


Joint Optimization Problem

Optimization function

$$(\bar{g}_s, \bar{g}_a) = \underset{g_s, g_a}{\operatorname{argmin}} \left\{ \begin{aligned} &E_x \left(L_d(x, g_s(x), g_a(x)) + \right. \\ &L_f(x, g_s(x), g_a(x)) + \\ &L_{reg}(x, g_s(x), g_a(x)) \left. + \right) \\ &+ L_{s,a}(g_s, g_a) \end{aligned} \right\}$$

Optimization framework



Regularization:

N random geometric transformations ψ_i at test time

$$\bar{E}_{f_c}(x) = \frac{1}{N+1} \left[E_{f_c}(x) + \sum_{i=1}^N \psi_i^{-1}(E_{f_c}(\psi_i(x))) \right]$$

Experimental Results

Adversarial and Similar generation

Great similarity: $x, g_s(x)$ and $g_a(x)$
reconstruction errors between $g_s(x)$ and $g_a(x)$

Weak Localization

$$IoU_i = \frac{M_{CT} \cap M_{E_i}}{M_{CT} \cup M_{E_i}} \quad AUC_{Loc} = \sum_i P_i(R_i - R_{i-1})$$

IOU SCORES AT DIFFERENT THRESHOLDS OF BINARIZATION - COMPARISON TO STATE OF THE ART METHODS WITHOUT (TOP) AND WITH (BOTTOM) AUGMENTATIONS

Explanation method	80	85	90	95	98
Gradient [1]	0.203	0.199	0.187	0.152	0.097
GradCAM [2]	0.256	0.252	0.236	0.199	0.117
BBMP [3]	0.277	0.263	0.244	0.199	0.105
Mask Generator [4]	0.233	0.226	0.204	0.154	0.087
"Naive"	0.222	0.219	0.208	0.169	0.103
Ours	0.259	0.264	0.259	0.221	0.137
Ours w/o Aug.	0.177	0.173	0.158	0.118	0.054
Ours w Aug.	0.239	0.230	0.208	0.156	0.087

Explanation method	Total AUC	Partial AUC	Time (s)
Gradient [1]	0.287	0.189	2.04
GradCAM [2]	0.374	0.274	2.83
BBMP [3]	0.397	0.302	5.99
Mask Generator [4]	0.326	0.229	17.14
"Naive"	0.327	0.226	0.99
Ours	0.404	0.308	0.68
Ours w/o Aug.	0.238	0.145	0.10
Ours w Aug.	0.325	0.232	0.75

