# Adversarial regression training for visualizing the progression of chronic obstructive pulmonary disease with chest x-rays

Ricardo Bigolin Lanfredi[1][0000−0001−8740−5796], Joyce D. Schroeder[2][0000−0002−7451−4886], Clement Vachet[1][0000−0002−8771−1803], and Tolga Tasdizen[1][0000−0001−6574−0366]

[1] Scientific Computing and Imaging Institute,
University of Utah, Salt Lake City UT 84112, USA
`ricbl@sci.utah.edu`
[2] Department of Radiology and Imaging Sciences,
University of Utah, Salt Lake City UT 84112, USA

**Abstract.** Knowledge of what spatial elements of medical images deep learning methods use as evidence is important for model interpretability, trustiness, and validation. There is a lack of such techniques for models in regression tasks. We propose a method, called visualization for regression with a generative adversarial network (VR-GAN), for formulating adversarial training specifically for datasets containing regression target values characterizing disease severity. We use a conditional generative adversarial network where the generator attempts to learn to shift the output of a regressor through creating disease effect maps that are added to the original images. Meanwhile, the regressor is trained to predict the original regression value for the modified images. A model trained with this technique learns to provide visualization for how the image would appear at different stages of the disease. We analyze our method in a dataset of chest x-rays associated with pulmonary function tests, used for diagnosing chronic obstructive pulmonary disease (COPD). For validation, we compute the difference of two registered x-rays of the same patient at different time points and correlate it to the generated disease effect map. The proposed method outperforms a technique based on classification and provides realistic-looking images, making modifications to images following what radiologists usually observe for this disease. Implementation code is available at `https://github.com/ricbl/vrgan`.

**Keywords:** COPD · chest x-ray · regression interpretation · visual attribution · adversarial training · disease effect · VR-GAN

## 1  Introduction

Methods of visual attribution in deep learning applied to computer vision are useful for understanding what regions of an image models are using [1]. These methods improve a model's interpretability, help in building user trust and validate if the model is using the evidence humans would expect it to use for its task.

They have been mostly used to explain models in classification tasks. A common way of formulating the problem of visual attribution is by asking "What regions of the image are weighted positively in the decision to output this class?". This question is not suitable for regression tasks, for which we propose to ask "What would this image look like if it had this other regression value?".

To answer the proposed question, we draw from conditional generative adversarial networks (GANs) [10] and image-to-image GANs [7], both of which have been shown to model complex non-linear relations between conditional labels, input images, and generated images. We hypothesize that using the regression target value for training a visual attribution model with the proposed novel loss will improve the visualization when compared to a similar formulation that only uses classification labels since it does not lose the information of a continuous regression target value by imposing a set of classes. We name our method visualization for regression with a generative adversarial network (VR-GAN).

We study our model in the context of how a value characterizing chronic obstructive pulmonary disease (COPD) is related to changes in x-ray images. COPD is defined by pulmonary function tests (PFTs) [8]. A PFT measures forced expiratory volume in one second ($FEV_1$), which is the volume of air a patient can exhale in one second, and forced vital capacity ($FVC$), which is the total volume of air a patient can exhale. A patient with $FEV_1/FVC$ ratio lower than 0.7 is diagnosed to have COPD.

Radiology provides a few clues that can be used for raising suspicion of COPD directly from an x-ray [4]. There is a higher chance of COPD when diaphragms are low and flat, corresponding to high lung volumes, and when the lung tissue presents low-attenuation (dark, or lucent), corresponding to emphysema, air trapping or vascular pruning. We show that our model highlights low and flat diaphragm and added lucency. Our method is, to the best of our knowledge, the first data-driven approach to model disease effects of COPD on chest x-rays.

## 1.1   Related work

One way to visualize evidence of a class using deep learning is to perform backpropagation of the outputs of a trained classifier [1]. In [11], for example, a model is trained to predict the presence of 14 diseases in chest x-rays, and class activation maps [15] are used to show what regions of the x-rays have a larger influence on the classifier's decision. However, as shown in [2], these methods suffer from low resolution or from highlighting limited regions of the original images.

In [2], researchers visualize what brain MRIs of patients with mild cognitive impairment would look like if they developed Alzheimer's disease, generating disease effect maps. To solve problems with other visualization methods, they propose an adversarial setup. A generator is trained to modify an input image which fools a discriminator. The modifications the generator outputs are used as visualization of evidence of one class. This setup inspires our method. However, instead of classification labels, we use regression values and a novel loss function.

There have been other works on generating visual attribution for regression. In [13], Seah et al. start by training a GAN on a large dataset of frontal x-rays,

and then train an encoder that maps from an x-ray to its latent space vector. Finally, Seah et al. train a small model for regression that receives the latent vector of the images from a smaller dataset and outputs a value which is used for diagnosing congestive heart failure. To interpret their model, they backpropagate through the small regression model, taking steps in the latent space to reach the threshold of diagnosis, and generate the image associated with the new diagnosis.

The loss function we provide for this task is similar to the cost function provided in [14]. Unlike our formulation, [14] models adversarial attackers and defenders in a game theoretic sense and arrives at an optimal solution for the defenders, using only linear models and applying it to simple features datasets. In [3], Bazrafkan et al. propose a method for training GANs conditioned in a continuous regression value. However, the model has a discriminator in parallel with the regressor, it is not used for visual attribution, and the used loss function is different than the one we propose.

## 2   Method

### 2.1   Problem Definition

We want to generate what an image would look like for different levels of a regression target value, without changing the rest of the content of the image. To formalize this mathematically, we can model an image as $x = f(y, z)$, where $x$ is a dependent variable representing an image, $y$ is an independent variable that determines an aspect of $x$, and $z$ another independent variable representing the rest of the content. In our application, $x$ is an x-ray image, $y$ is the value from a PFT of the same patient taken contemporaneously to the chest x-ray approximating the severity of COPD for that patient, and $z$ represents factors such as patient anatomy unrelated to COPD and position of the body at the moment the x-ray was taken.

We want to construct a model that, given an image $x$ associated with a value $y$ and a content $z$, can generate an image $x'$ conditioned on the same $z$, but on a different value $y'$. By doing this, we can visualize what impacts the change of $y$ to $y'$ has on the image. Similar to [2], we formulate the modified image as

$$x' = \Delta x + x = G(x, y', y) + x, \tag{1}$$

where $G$ is a conditional generator, and $\Delta x$ is a difference map or a disease effect map. By summing $\Delta x$ to $x$, the task of $G$ is made easier, since $G$ only has to model the impact of changing $y$ to $y'$, and the content $z$ should be already in $x$.

### 2.2   Loss function

Fig. 1 shows the loss terms that are defined in this section and how the modules are connected for training a VR-GAN. We start by defining a regressor $R$ that has the task to, given an image $x$, predict its $y$ value. We start building our loss
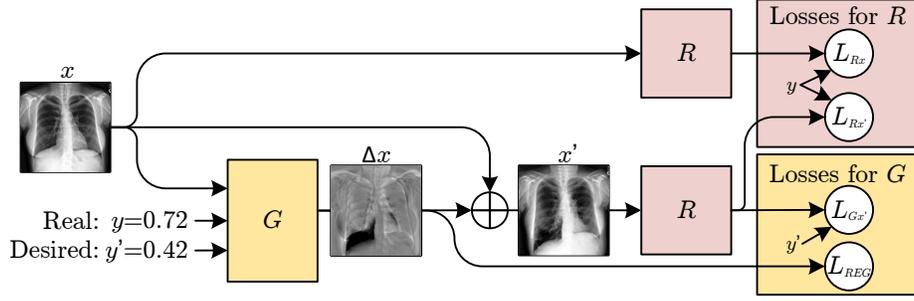
**Fig. 1.** Overall model architecture for training with the proposed adversarial loss. The losses $L_{Rx}$, $L_{Rx'}$ and $L_{Gx'}$ are $L1$ regression losses, and $L_{REG}$ is an $L1$ norm penalty.

function by defining a term, which is used to optimize only the weights of $R$, to assure $R$ can perform the task of regression over the original dataset:

$$L_{Rx} = L_{Rx}(x, y) = \|R(x) - y\|_{L1}. \tag{2}$$

The regressor is used to assess how close to having $y'$ an image $x'$ is. We also define a term, which is used to optimize only the weights of $G$, to make $G$ learn to create a map $\Delta x$ that, when added to the original image, changes the output of $R$ to a certain value $y'$:

$$L_{Gx'} = L_{Gx'}(\Delta x + x, y') = \|R(G(x, y', y) + x) - y'\|_{L1}. \tag{3}$$

Training a model using only these two terms would lead to an $R$ that does not depend on $G$, and, consequently, to a $G$ that can modify the output of $R$ in simple and unexpressive ways, similar to noisy adversarial examples [5]. We define an adversarial term, which is used to optimize the weights of $R$, to make $R$ learn to output the same value as the original image for the modified image:

$$L_{Rx'} = L_{Rx'}(\Delta x + x, y) = \|R(G(x, y', y) + x) - y\|_{L1}. \tag{4}$$

As $G$ learns trivial or unrealistic modifications to images, $R$ learns to ignore them due to Eq. (4), forcing the generator to create more meaningful $\Delta x$. This game between $G$ and $R$ should reach an equilibrium where $G$ produces images that are realistic and induces the desired output from $R$. At this point, $R$ should not be able to find modifications to ignore, being unable to output the original $y$. In our formulation, $R$ replaces a discriminator from a traditional GAN.

We define another loss term to assure that $G$ only generates what is needed to modify the label from $y$ to $y'$ and does not modify regions that would alter $z$. An $L1$ penalty is used over the difference map to enforce sparsity:

$$L_{REG} = L_{REG}(\Delta x) = \|\Delta x\|_{L1}. \tag{5}$$

Intuitively, when the modifications generated by $G$ are unrealistic and ignored by $R$, not having any impact to the term defined in Eq. (3), the norm penalty defined in Eq. (5) should enforce their removal from the disease effect map.

The complete optimization problem is defined as

$$G^* = \underset{G}{\operatorname{argmin}}(\lambda_{Gx'}L_{Gx'} + \lambda_{REG}L_{REG}), R^* = \underset{R}{\operatorname{argmin}}(\lambda_{Rx'}L_{Rx'} + \lambda_{Rx}L_{Rx}), \quad (6)$$

where the $\lambda$'s are hyperparameters. Optimizations are performed alternatingly.

## 3  Experiments

We used a U-Net [12] as $G$. The conditioning inputs $y$ and $y'$, together with their difference, were normalized and concatenated to the U-Net bottleneck layer (Fig. 1). For $R$, we used a Resnet-18 [6], pretrained on ImageNet and with the output changed to a single linear value. We froze the batch normalization parameters in $R$, as in [2]. Since $R$ depends on the supervision from the original regression task, it will only be able to learn to output values in the range of the original $y$. Therefore, during training we sampled $y'$ from the same distribution as $y$. The hyperparameters were chosen as $\lambda_{Gx'} = 0.3$, $\lambda_{REG} = 0.03$, $\lambda_{Rx} = 1.0$, $\lambda_{Rx'} = 0.3$, using validation over the toy dataset presented in Section 3.1. The same set of hyperparameters were used for the x-ray dataset and were not sensitive to change of dataset. Adam [9] was used as the optimizer, with a learning rate of $10^{-4}$. To prevent overfitting, early stopping was used.

We employed the VA-GAN method presented in [2] as a baseline, since it is a classification version of our method[3]. We used $\lambda = 10^2$ and gradient penalty with a factor of 10, as in [2]. Baseline optimizers and models were the same as the ones described for our model.

We compared the results visually to check if they agreed with radiologists' expectations. For quantitative validation, we used the normalized cross-correlation between the generated $\Delta x$ map and the expected $\Delta x$ map, averaged over the test set. Each result is given with its mean and its standard deviation over 5 tests, with training initialized using distinct random seeds.

### 3.1  Toy dataset

To test our model, we generated images of squares, superimposed with a Gaussian filtered white noise and with a resolution of $224 \times 224$. An example is presented in Fig. 2(a). The side length of the square is proportional to a regression target $y$ that follows a Weibull distribution with a shape parameter of 7 and a scale parameter of 0.75. The class threshold for our baseline model [2] was set at 0.7. It was trained to receive images of big squares ($y \geq 0.7$), and output a difference map that made that square small ($y < 0.7$). We generated 10,000 images for training. Since we generated the images, we could evaluate with perfect ground truth for $\Delta x$. We evaluated using input examples where $y \geq 0.7$ and sampling $y' < 0.7$ from the Weibull distribution. This resulted in 5,325 images for validation and 5,424 images for testing.

---

[3] Our implementations of VA-GAN and VR-GAN extends code from `github.com/orobix/Visual-Feature-Attribution-Using-Wasserstein-GANs-Pytorch`
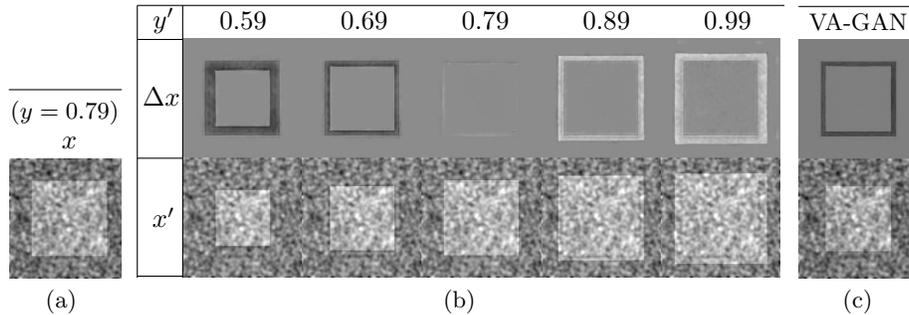
**Fig. 2.** Results on a test example from the toy dataset. **Top:** difference maps $\Delta x$. **Bottom:** images of squares, $x$ or $x'$. **(a)** Original image $x$, with square size $y = 0.79$. **(b)** VR-GAN result for several desired square side lengths $y'$. **(c)** VA-GAN result.

Examples of difference maps and modified versions of the original image for a few levels of the desired square side length are presented in Fig. 2(b). While the baseline presents a fixed modification for an image, shown in Fig. 2(c), and can only generate smaller squares, our method can generate different levels of change for the map, and also generate both bigger and smaller squares. The baseline [2], using VA-GAN, obtained a score of $0.780 \pm 0.007$ for the normalized cross-correlation, while our method, VR-GAN obtained a score of $0.853 \pm 0.014$.

### 3.2   Visualizing the progression of COPD

We gathered a dataset of patients that had a chest x-ray exam and a PFT within 30 days of each other at the University of Utah Hospital from 2012 to 2017. This study was performed under an approved Institutional Review Board process[4] by our institution. Data was transferred from the hospital PACS system to a HIPAA-compliant protected environment. Orthanc[5] was used for data de-identification by removing protected health information. Lung transplants patients were excluded, and only posterioranterior (PA) x-rays were used. PFTs were only associated with their closest x-ray exam and vice-versa. For validation and testing, all subjects with at least one case without COPD (used as original image $x$) and one case with COPD (used as desired modified image $x'$) were selected, using COPD presence as defined by PFTs. This setup was chosen because the trained baseline model can only handle transitions from no disease to disease. For each of these subjects, we used all combinations of pair of cases with distinct diagnoses. The average time between paired exams was 17 months. We used 3,414 images for training, 208 pair of images for validation and 587 pair of images for testing. Images were cropped to a centered square, resized to $256 \times 256$ and randomly cropped to $224 \times 224$. We equalized their histogram and normalized their individual intensity range to [-1,1]. We used $FEV_1/FVC$

---

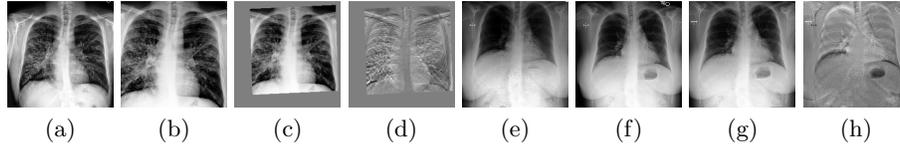[4] IRB_00104019, PI: Schroeder MD

[5] `orthanc-server.com`

**Fig. 3.** Examples of results of the alignment. **(a) and (e):** Reference images without COPD. **(b) and (f):** Images to align, with COPD. **(c) and (g):** Aligned images. **(d) and (h):** Difference between reference and aligned images, used as $\Delta x$ ground truth.
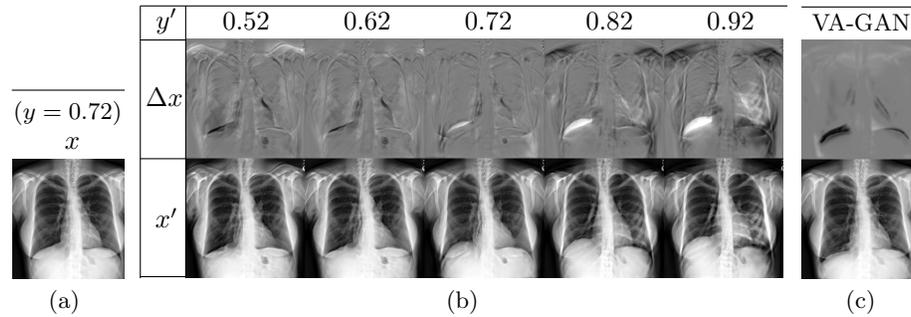


**Fig. 4.** Results on a test example of the COPD dataset. **Top:** disease effect maps $\Delta x$. **Bottom:** chest x-rays, $x$ or $x'$. **(a)** Original image $x$, with $FEV_1/FVC$ ($y$) 0.72. **(b)** VR-GAN results for several desired $FEV_1/FVC$ ($y'$). The lower this value, the more severe the disease. **(c)** VA-GAN results.
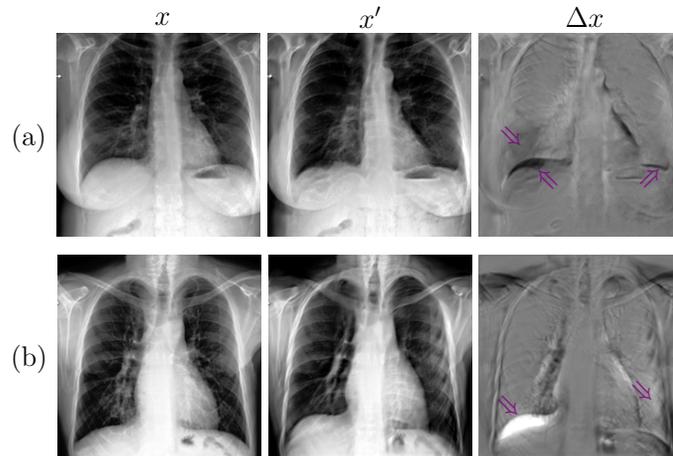


**Fig. 5.** Examples of disease effects that correspond with visual feature changes expected radiologically. In $\Delta x$, gray represents no change, black a decrease and white an increase in image intensity. From left to right: original image $x$, modified image $x'$ and disease effect map $\Delta x$. **(a)** $y = 0.8$ to $y' = 0.4$ (increasing severity). Purple arrows highlight flat and low diaphragm (bottom two arrows) and added lucency (top left arrow). **(b)** $y = 0.37$ to $y' = 0.8$ (decreasing severity). Purple arrows highlight high and curved diaphragm (left arrow) and reduced lucency (right arrow).

as the regression target value $y$. To generate ground truth disease effect maps, we aligned two x-rays of the same patient, using an affine registration employing the pystackreg library[6], and subtracted them. Examples are shown in Fig. 3.

Disease effect maps and modified images for both methods are presented in Fig. 4. Note that our method can modify images to increase and decrease severity by any desired amount in contrast to [2], which can only be trained to modify the classification of images in a single direction. In Fig. 5, we show images generated using VR-GAN which highlight the height and flatness of diaphragm and show changes in the level of lung lucency, features that radiologists use as evidence for COPD on chest x-rays. Small changes in the cardiac contour are consistent with the accommodation of a shift in lung volume. Using normalized cross-correlation, VA-GAN obtained a score of 0.012±0.015, while VR-GAN obtained a score of 0.127±0.017. The low correlation scores may result from imperfect alignments with affine transformations and potential changes between x-ray pairs unrelated to COPD. However, our method still obtained a significantly better score.

## 4   Conclusion

We introduced a visual attribution method for datasets with regression target values and validated it for a toy task and for chest x-rays associated with PFTs, assessing the impact of COPD in the images. We demonstrated significant improvement in the disease effect maps generated by a model trained with adversarial regression when compared to a baseline trained using classification labels. Furthermore, the generated disease effect maps highlighted regions that agree with radiologists' expectations and produced realistic images.

## References

1. Ancona, M., Ceolini, E., Öztireli, C., Gross, M.: Towards better understanding of gradient-based attribution methods for deep neural networks. In: ICLR (2018)
2. Baumgartner, C.F., Koch, L.M., Tezcan, K.C., Ang, J.X.: Visual feature attribution using Wasserstein GANs. In: CVPR (2018)
3. Bazrafkan, S., Corcoran, P.: Versatile auxiliary regressor with generative adversarial network (VAR+GAN). arXiv preprint arXiv:1805.10864 (2018)
4. Foster, W.L., et al.: The emphysemas: radiologic-pathologic correlations. Radio-Graphics **13**(2), 311–328 (1993)
5. Goodfellow, I., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: ICLR (2015)
6. He, K., et al.: Deep residual learning for image recognition. In: CVPR (2016)
7. Isola, P., Zhu, J., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: CVPR (2017)
8. Johnson, J.D., Theurer, W.M.: A stepwise approach to the interpretation of pulmonary function tests. American family physician **89 5**, 359–66 (2014)
9. Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization. In: ICLR (2015)

---

[6] bitbucket.org/glichtner/pystackreg

10. Mirza, M., Osindero, S.: Conditional generative adversarial nets. CoRR **abs/1411.1784** (2014)
11. Rajpurkar, P., Irvin, J., et al.: Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. CoRR **abs/1711.05225** (2017)
12. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional networks for biomedical image segmentation. In: MICCAI (2015)
13. Seah, J.C.Y., et al.: Chest radiographs in congestive heart failure: Visualizing neural network learning. Radiology **290**(2), 514–522 (2019)
14. Tong, L., et al.: Adversarial regression with multiple learners. In: ICML (2018)
15. Zhou, B., Khosla, A., Lapedriza, À., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: CVPR (2016)