

# THIS CHANGES TO THAT : COMBINING CAUSAL AND NON-CAUSAL EXPLANATIONS TO GENERATE DISEASE PROGRESSION IN CAPSULE ENDOSCOPY

*Anuja Vats\**   *Ahmed Mohammed\*<sup>†</sup>*   *Marius Pedersen\**   *Nirmalie Wiratunga<sup>‡</sup>*

\* Department of Computer Science, NTNU, Gjøvik, Norway   <sup>†</sup>Sintef Digital, Oslo, Norway   <sup>‡</sup>School of Computing, Robert Gordon University, Aberdeen, Scotland

## ABSTRACT

Due to the unequivocal need for understanding the decision processes of deep learning networks, both modal-dependent and model-agnostic techniques have become very popular. Although both of these ideas provide transparency for automated decision making, most methodologies focus on either using the modal-gradients (model-dependent) or ignoring the model internal states and reasoning with a model’s behavior/outcome (model-agnostic) to instances. In this work, we propose a unified explanation approach that given an instance combines both model-dependent and agnostic explanations to produce an explanation set. The generated explanations are not only consistent in the neighborhood of a sample but can highlight causal relationships between image content and the outcome. We use Wireless Capsule Endoscopy (WCE) domain to illustrate the effectiveness of our explanations. The saliency maps generated by our approach are comparable or better on the softmax information score.

**Index Terms**— Explainable AI, Counterfactual, Semifactual, saliency map, capsule endoscopy

## 1. INTRODUCTION

There has been a rapid integration of deep learning based models in real-world applications, including high risk ones such as healthcare and defence owing to their unparalleled predictive performance [5]. Such real-world deployment and usage of models accompanies with it the moral obligation to make their decision processes transparent. This is necessary not only for accountability of high stake decisions but also for the identification and mitigation of algorithmic or societal bias [14, 7]. This has led research to continue attempts at opening the black boxes, to gain insight in decision making processes [4, 20, 10] while also considering that useful explanations could emerge through model-agnostic explainer methods ([15, 23, 6, 24]). Although both of these approaches are suited to explaining model predictions, dominant explanation approaches today focus on one or the other.

Factual explainers that reason with gradients [19, 10, 20] aim to identify regions or pixels within an image that most

significantly contributed to the prediction and thereafter visualize these attribution weights in saliency maps [4, 20, 19]. For example, given an endoscopic image with an ulcer, a saliency map would highlight the ulcer region in response to a question such as “Why did you make that decision?”. However, although popular saliency methods [19, 20, 10] are fairly easy to implement they also have limitations.

One limitation of gradient-based saliency methods is the use of a baseline image and the sensitivity to the choice of that baseline [13]. Here a baseline helps contrast the query scenario from a “baseline” scenario and typically marks an absence against which the “presence” can be measured, e.g a black image. Since attribution maps are then accumulated over a classical linear path from baseline to the query; the choice of baseline is crucial to the success of the explanation. A second limitation is that the pixel perturbations done to arrive from a baseline to query are typically blind to image content. We argue against such pixel perturbations to create images between the baseline and query image as well as the baseline itself. Because not only are the images in between not natural but are also prone to abnormal gradient behaviours from irrelevant pixels as identified in [10].

Consider our ulcer example from before, the perturbations with respect to clinical biomarkers relating to ulcer abnormality are more meaningful than individual pixels. For example take perturbations that cause “more or less inflammation around a suspected ulcer”, knowing that an ulcer is often accompanied with inflammation is important, as a lack of it might suggest incorrectly to the doctor that the suspected ulcer is just intestinal debris stuck to the surface. Such perturbations are not only more meaningful but every image resulting from them is directly interpretable. This also implies that a more apt baseline would be one that marks absence of the biomarker (here the ulcer) and not complete absence of the signal. A third limitation of such methods is their single-pointwise explanation mode of operation [3, 1], whereby an explanation to a given image is made in isolation of its locality, i.e. without considering its neighborhood (i.e. how explanation changes as the input changes slightly).

Counterfactual reasoning has gained popularity [6] as a locality-aware explainer that is model agnostic (i.e. does not need access to a network’s internal mechanism (gradients, layer activations, etc). Often these provide causally under-

---

Thanks to Research Council of Norway for funding (Project no: 300031). <sup>\*</sup>authors contribute equally.

standable explanations which have been argued to be GDPR compliant [23] and help address questions on fairness, trust and robustness [8]. These explanations generate a counterfactual as an alternative scenario with a desirable outcome that counters the observed (real) outcome. As such they generate explanations through relationships like: “If the ulcer had not been present, this image would not be abnormal.” In other words, it pinpoints how the input must change to flip the outcome. It is clear how such explanations might seem intuitive and interesting [6] to a doctor in our context. In fact, counterfactual thinking is very natural to how humans reason especially in response to negative outcomes in order to prevent them in the future [16].

In vision, explaining an instance with its corresponding counterfactual [8, 2] has become common for highlighting changes that would most easily flip the prediction. In [8], authors perform minimal edits by swapping regions of a query image from a distractor image till a decision flip occurs. However, the choice of a suitable distractor image is crucial for quick convergence but this choice can be unintuitive for some domains such as the medical domain or when little information is available for the dataset. Further, the image resulting from such edits can be unnatural looking at times and therefore lack explainability. For such explanations to be efficient, the changes applied to the image for a different prediction must be minimal and human interpretable [23]. Alipour et al. [2] use the latent space of a pretrained styleGAN for retrieving counterfactual latent codes and is similar to our approach in idea but differs in implementation (their method produces causal explanations only unlike ours, while also employing pretrained attribute detectors in latent space that are largely unavailable for medical domains.) Recently, semi-factuals have been argued to offer advantages similar to counterfactuals [12]. As opposed to counterfactuals that propose explanations as ‘If only’ clause, semi-factuals propose explanation of type ‘even if’ i.e. what changes to the situation would still lead to the same outcome. In our earlier example, a semi-factual image might illustrate the inflammatory changes that occur right before an ulcer starts forming, as this point the doctor will still identify the image as abnormal.

Despite advantages, one of the biggest challenges in using counterfactual and semi-factual explanations (together referred to as contrastive explanations) lies in generating instances that not only expose realistic and progressive visual changes smoothly (as to be directly understandable), but also ensuring progression alignment with the expected class prediction behaviour (congruous change in softmax score) [2]. Addressing this need for aligned progression both in the image and classifier space is precisely the problem we propose to solve in this work. We argue that in favor of human interpretability and algorithmic transparency, explanations that support both the aforementioned modes (causal and non-causal) are better than either one. We demonstrate the effectiveness of our explanations in the domain of WCE with

focus on Ulcerative Colitis (UC). We use the UC biomarkers used by experts in diagnosis such as inflammation and ulcerations as progression attributes to manage counterfactual explanations. The main contributions are:

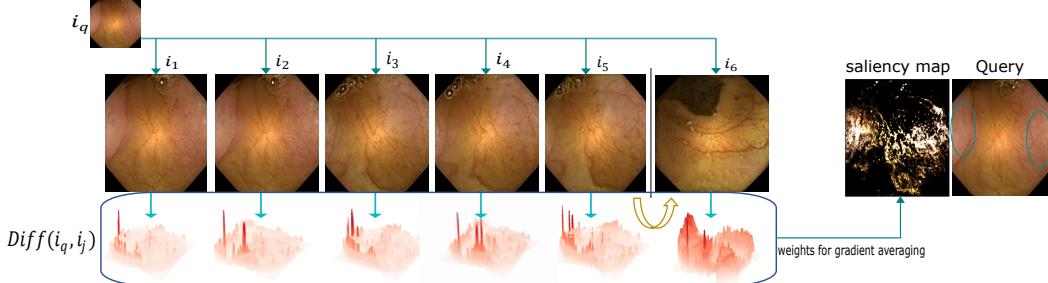
- a unified framework that generates both causal and non-causal explanations for each decision;
- a method to control progression along a specific UC biomarker such that the counterfactual relationships inferred are causal as opposed to being adhoc; and
- a formal algorithm to generate saliency maps that are comparable to (or better) than others on the Softmax Information Curve (SIC) metrics 3.

## 2. METHODOLOGY

Given an attribute of choice (e.g., a UC biomarker like inflammation, vascular pattern etc.) and the query image  $i_a^q$ , the goal is to retrieve two instances that are closest to the decision boundary as semifactual (on the same side) and counterfactual (on the other side), while preserving visual interpretability along a path of images directed by the attribute (Figure 2). Regions of importance is highlighted by a saliency map (can be generated for each image on the path, including the query).

Given a classifier  $\mathcal{C}$  that outputs label  $y \in \{0, 1\}$  through a prediction function  $f : \mathbb{R}^n \rightarrow [0, 1]$  for an image  $x_i \in \mathbb{R}^{512 \times 512}$ , an explanation set is produced,  $\mathcal{X} = \{i_{sm}, i_{cf}, i_{sf}\}$ , along attribute  $a$ . Here  $i_{sm}$  is the saliency map,  $i_{cf}$  is the nearest counterfactual and  $i_{sf}$  the semifactual along  $a$ . We use this to generate an explanation: “image  $i_a^q$  is abnormal with probability  $p$  due to signs/regions highlighted by the saliency map  $i_{sm}$ . The least amount of abnormality required for the prediction to be abnormal is seen in  $i_{sf}$  (semifactual). However, if the abnormal signs change to as in  $i_{cf}$  (counterfactual) the image would no longer be classified as abnormal”. Importantly the changes along the single attribute,  $a$ , is also directly visually interpretable by a user (e.g., a doctor).

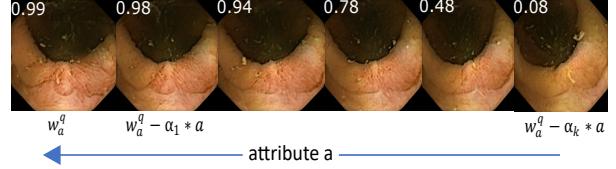
**Attribute discovery in latent space:** We use StyleGAN2 [11] and train it on WCE images (discussed in Dataset and Training details sec). StyleGAN2 uses a mapping network between a latent variable and the network generator,  $G$ , which transforms the latent variable to an intermediate  $d$ -dimensional space,  $W$ , of latent vectors,  $w \in \mathbb{R}^d$ , where style attributes are known to be more amenable to control. We use SeFA [17] for the unsupervised discovery of attributes in the intermediate  $W$  space. In the natural image domain, pretrained attribute detectors can be utilized for labeling these attributes however for our case of pathological and anatomical variations of the colon such attribute detectors are not available a priori. We perform clustering on images using TSNE [21] for isolating attributes relevant to pathological changes. This is done by planting seed images before



**Fig. 1.** Figure shows  $i_a$  and the corresponding directional derivatives. The derivatives expose the semantic similarity between the query and its neighbors. We use this similarity to weigh in the contribution of each neighbor towards the saliency map.

clustering that had been identified by a doctor as good representatives of UC pathological changes. Upon clustering we sampled the attributes closest to seed images and have used these as explanation attributes.

**Generating the explanation set  $\mathcal{X}$ :** Once relevant attributes are identified, to explain a query  $i_a^q$  with latent  $w_a^q \in \mathcal{R}^d$  such that  $i_a^q = G(w_a^q)$  along attribute  $a$ , a set of  $k$  local images  $i_a$  is created,  $i_a = \{i_a^1, i_a^2, \dots, i_a^k\}$  from latents  $w_a = \{w_a^1, w_a^2, \dots, w_a^k\}$  where  $w_a^j = w_a^q - \alpha_j * a$  and  $\alpha_j$  varies linearly in  $[A, B]$  and  $a \in \mathbb{R}^{512}$  is the aforementioned attribute vector. In Figure 3, attribute  $a$  corresponds to (reddish) inflammatory regions and set  $i_a$  can be understood as images with decrease in severity of such inflammation as  $\alpha$  progresses from  $A = 0$  to  $B = 30$ .  $i_{cf}, i_{sf}$  in  $\mathcal{X}$  are retrieved based on the classifier output for  $i_a$  such that  $i_{cf} = argmax(\sigma(\mathcal{C}(i_a^j))) \forall \sigma(\mathcal{C}(i_a^j)) < 0.5$  and  $i_{sf} = argmin(\sigma(\mathcal{C}(i_a^j))) \forall \sigma(\mathcal{C}(i_a^j)) > 0.5$  where  $\sigma$  is the softmax function. For the saliency map, to avoid the spuriousness observed in previous literature, we use the latent space to curate a neighborhood such that every image in the neighborhood of a query varies only along the chosen attribute. In other words, the pixel changes that occur in this neighborhood are neither uniform nor content blind [10, 20], but



**Fig. 3.** Images in  $i_a$  along attribute  $a$ . Top left corner shows softmax score. Notice how apart from effected region (for attribute  $a$ ), other regions in the image undergo only minimal changes. As a result, the generated explanations are consistent in the locality of a query.

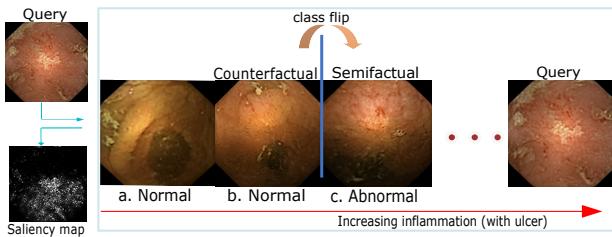
targeted towards those pixels that most strongly affect the attribute/biomarker. We use directional derivatives in  $i_a$  along attribute  $a$  for identifying these regions and weight them based on semantic similarity with  $i_a^q$  to generate the saliency map. The directional derivative  $Diff(i_a^q, i_a^j)$  between the query and  $i_a^j = G(w_a^j) \forall \{w_a^j\}_{j=1}^k$  is given by:

$$Diff(i_a^q, i_a^j) = \left| \frac{G(w_a^q) - G(w_a^j)}{1} \right| \quad (1)$$

The directional derivatives  $Diff(i_a^q, i_a^j)$  over  $i_a^q$  and  $i_a$  exposes pixels with consistent change in the direction of increasing/decreasing attribute (see Figure. 1), in other words it is a measure of semantic similarity to the image being explained. We use these derivatives to measure the contribution of each image in  $i_a$ . A formal algorithm is described in algorithm 1.

#### Dataset and Training Details:

The dataset consists of approximately 200k unlabeled WCE images. The majority of images come from WCE examinations of 10 patients with varying UC activity, as well as other pathologies with PillCam Colon 2 Capsule, Medtronic. The images are 576x576 in resolution with varying degree of bowel cleanliness. In addition to this we use PS-DeVCEM dataset[22] with 80k images of the same capsule modality. Remaining images come from the OSF-Kvasir Dataset [18] with 3478 images from seven classes taken with the capsule modality Olympus EC-S10. We use StyleGAN2 without progressive growing and work exclusively on the original inter-



**Fig. 2.** The approach explains a query image along the ulcer attribute path together with a semifactual and counterfactual along the same path. Here the query exhibits an abnormality with inflammation. Even with inflammation reduced down to as in (c) the prediction would still be abnormal (semifactual). However, if only the visual signs change from (c) to as in (b), the prediction would be normal (counterfactual).

---

**Algorithm 1:** Saliency map generation

<https://www.overleaf.com/project/632c20113b1a4fb9b68a2cfb>

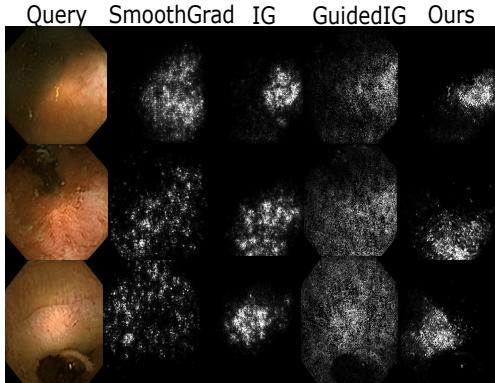
**Input:** Classifier  
 $C, \text{query } i_a^q; i_a = \{i_a^1, i_a^2 \dots i_a^k, i_a^q\};$   
**Output:**  $i_{sm}$   
**foreach**  $i_a^q$  **do**  
 predict output class probabilities for  $i_a$   
 $output \leftarrow C(i_a)$   
 backpropagate and collect gradients wrt  $i_a$   
 $[grad_a^1, grad_a^2 \dots grad_a^k] \leftarrow i_a.grad()$   
 directional derivatives along attribute a  
**foreach**  $i \in (i_a \setminus \{i_a^q\})$  **do**  
 $Diff(i_q, i) \leftarrow \left| \frac{G(w_a^q) - G(w_a^i)}{1} \right|$   
**end**  
 $S(i_q, a) \leftarrow \frac{\sum_{j=1}^k grad_a^j \cdot Diff(i_q, j)}{k}$   
 $i_{sm} \leftarrow meanThresholding(S(i_q, a))$   
**end**


---

mediate latent space  $W$  and not the extended space  $W^+$ . The model was trained on TwinTitan RTX for 30 days.<sup>1</sup>

### 3. RESULTS

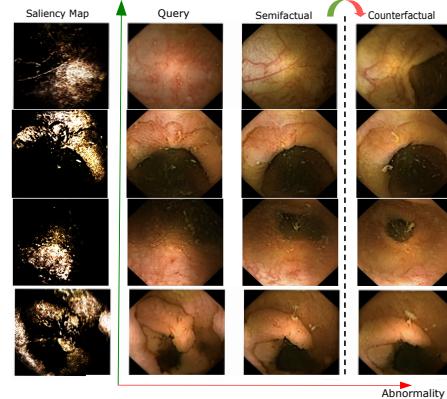
**Qualitative Comparison:** GuidedIG [10] produces noisy saliency maps (as only pixels with low partial derivatives are moved towards their original intensity at each step to avoid high gradient regions and thus abnormal behavior), but if the pixels affecting the decision are not localized but spread globally across the image, (as in WCE), the resulting saliency map can appear to be noisier. Similarly, while SmoothGrad[19] captures the right regions, the saliency maps are overall noisy. Integrated Gradients [20] correlates very closely with our maps. Figure 5 shows  $\mathcal{X}$  for various query images.



**Fig. 4.** Qualitative comparison of saliency maps between our approach and other approaches. Integrated Gradients (IG) [20], Guided integrated gradients [10], SmoothGrad [19]

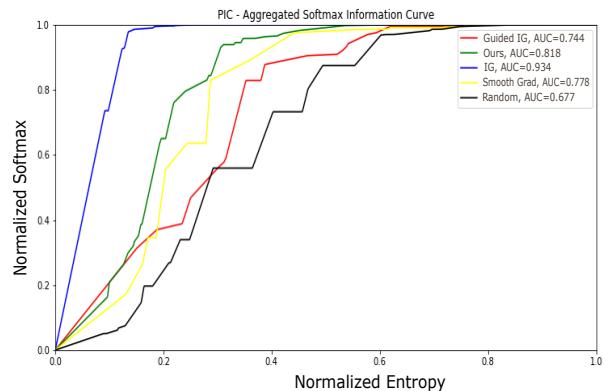
**Quantitative Comparison:** We use Softmax Information Curve (SIC AUC) [9] for quantitative comparison. SIC AUC

<sup>1</sup>Github : <https://github.com/anuja13/ContrastiveExplanations>



**Fig. 5.** Figure shows  $\mathcal{X}$  generated with this approach on different query images (column 3). Best viewed in color.

measures the softmax score of a model against salient regions indicated by the saliency map. Figure 6 shows the SIC AUC for different approaches averaged over 50 images. Integrated gradients achieve the best score followed by our approach. We suspect this to be due to the SIC score's preference for smallest regions of effect (as in IG) instead of identifying all contributing regions (as in ours).



**Fig. 6.** Quantitative comparison of saliency maps : Median Softmax Information curves

### 4. CONCLUSION

In this work, we propose a framework for generating causal as well as non-causal explanations for any image classifier. Our model is network agnostic and supports not only visual insight into model decisions, but offers end users the opportunity to visualize alternate scenarios relevant to the current situation, for example in prognosis of UC as shown in this work. To the best of our knowledge, this is one of the first works to propose a single framework for generating both causal as well as non-causal explanation for deep learning based models.

## References

- [1] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim. Sanity checks for saliency maps. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [2] K. Alipour, A. Lahiri, E. Adeli, B. Salimi, and M. Pazzani. Explaining image classifiers using contrastive counterfactuals in generative latent spaces. *arXiv preprint arXiv:2206.05257*, 2022.
- [3] D. Alvarez-Melis and T. S. Jaakkola. On the robustness of interpretability methods. 2018. cite arxiv:1806.08049Comment: presented at 2018 ICML Workshop on Human Interpretability in Machine Learning (WHI 2018), Stockholm, Sweden.
- [4] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11(61):1803–1831, 2010.
- [5] S. Benjamens, P. Dhunnoo, and B. Meskó. The state of artificial intelligence-based fda-approved medical devices and algorithms: an online database. *NPJ digital medicine*, 3(1):1–8, 2020.
- [6] R. M. Byrne. Counterfactuals in explainable artificial intelligence (xai): Evidence from human reasoning. In *IJCAI*, pages 6276–6282, 2019.
- [7] S. Chodosh. Courts use algorithms to help determine sentencing, but random people get the same results. *Popular Science*, 2018.
- [8] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee. Counterfactual visual explanations. In K. Chaudhuri and R. Salakhutdinov, editors, *36th ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 2376–2384. PMLR, 09–15 Jun 2019.
- [9] A. Kapishnikov, T. Bolukbasi, F. Viégas, and M. Terry. Xrai: Better attributions through regions. In *ICCV*, pages 4948–4957, 2019.
- [10] A. Kapishnikov, S. Venugopalan, B. Avci, B. Wedin, M. Terry, and T. Bolukbasi. Guided integrated gradients: An adaptive path method for removing noise. In *CVPR*, pages 5050–5058, 2021.
- [11] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020.
- [12] E. M. Kenny and M. T. Keane. On generating plausible counterfactual and semi-factual explanations for deep learning. *Proc. AAAI*, 35(13):11575–11585, May 2021.
- [13] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim. The (un)reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.
- [14] Z. Obermeyer, B. Powers, C. Vogeli, and S. Mullanathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.
- [15] M. T. Ribeiro, S. Singh, and C. Guestrin. Why should i trust you? explaining the predictions of any classifier. In *22nd ACM SIGKDD*, pages 1135–1144, 2016.
- [16] N. J. Roese. Counterfactual thinking. *Psychological bulletin*, 121(1):133, 1997.
- [17] Y. Shen and B. Zhou. Closed-form factorization of latent semantics in gans. In *CVPR*, pages 1532–1540, 2021.
- [18] P. H. Smedsrød, V. Thambawita, S. A. Hicks, H. Gjestang, O. O. Nedrejord, E. Næss, H. Borgli, D. Jha, T. J. D. Berstad, S. L. Eskeland, et al. Kvasir-capsule, a video capsule endoscopy dataset. *Scientific Data*, 8(1):1–10, 2021.
- [19] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [20] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- [21] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- [22] M. A. . P. M. Vats Anuja. From labels to priors in capsule endoscopy: a prior guided approach for improving generalization with few labels. *Scientific Reports*, 12:15708, 2022.
- [23] S. Wachter, B. Mittelstadt, and C. Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.
- [24] N. Wiratunga, A. Wijekoon, I. Nkisi-Orji, K. Martin, C. Palihawadana, and D. Corsar. Discern: Discovering counterfactual explanations using relevance features from neighbourhoods. In *33rd ICTAI*, pages 1466–1473. IEEE, 2021.