

Fairness and Robustness of Contrasting Explanations^{*}

André Artelt¹ and Barbara Hammer¹

CITEC - Cognitive Interaction Technology
Bielefeld University, 33619 Bielefeld, Germany
{aartelt,bhammer}@techfak.uni-bielefeld.de

Abstract. Fairness and explainability are two important and closely related requirements of decision making systems. While ensuring and evaluating fairness as well as explainability of decision making systems has been extensively studied independently, only little effort has been investigated into studying fairness of explanations on their own - i.e. the explanations it self should be fair. In this work we formally and empirically study individual fairness and robustness of contrasting explanations - in particular we consider counterfactual explanations as a prominent instance of contrasting explanations. Furthermore, we propose to use plausible counterfactuals instead of closest counterfactuals for improving the individual fairness of counterfactual explanations.

Keywords: XAI · Contrastive Explanations · Counterfactual Explanations · Fairness · Robustness

1 Introduction

Fairness and transparency are fundamental building blocks of ethical artificial intelligence (AI) an machine learning (ML) based decision making systems. In particular, the increasing use of automated decision making systems have strengthened the demand for trustworthy systems. The criticality of transparency was also recognized by policy makers which resulted in legal regulations like the EU’s GDPR [29] that grants the user a right to an explanation. Therefore, the research community focused a lot on the question how to realize explainability and transparency of AI and ML systems [15, 18, 34, 37]. Nowadays, there exist diverse methods for explaining ML models [18, 27]. One specific family of methods are model-agnostic methods [18, 31]. Model-agnostic methods are flexible in the sense that they are not tailored to a particular model or representation. This makes model-agnostic methods (in theory) applicable to many different types of ML models. In particular, “truly” model-agnostic methods do not need access to the training data or model internals. It is sufficient to have an interface for passing data points to the model and obtaining the output of the model - the underlying model it self is viewed as a black-box.

^{*} We gratefully acknowledge funding from the VW-Foundation for the project *IMPACT* funded in the frame of the funding line *AI and its Implications for Future Society*.

Examples of model-agnostic methods are feature interaction methods [16], feature importance methods [14], partial dependency plots [42] and local methods that approximates the model locally by an explainable model (e.g. a decision tree) [17, 32]. These methods explain the models by using features as vocabulary.

A different class of model-agnostic explanations are example-based explanations where a prediction or behavior is explained by a (set of) data points [1]. Instances of example-based explanations are prototypes & criticisms [20] and influential instances [21]. Another instance of example-based explanations are counterfactual explanations [39]. A counterfactual explanation is a change of the original input that leads to a different (specific) prediction/behavior of the decision making system - *what has to be different in order to change the prediction of the system?* Such an explanation is considered to be fairly intuitive, human-friendly and useful because it tells people what to do in order to achieve a desired outcome [27, 39]. Further, there exists strong evidence that explanations by humans are often counterfactual in nature [10]. We will focus on these types of explanations in this work.

Despite the recent success stories of AI and ML systems, ML systems were also involved in prominent failures like predictive policing [5] and loan approval [40]. Many of those failures, where the systems exhibit unethical behaviour, deal with predictions that are made based of sensitive attributes like race, gender, etc. Using such sensitive attributes for making decisions is considered to be unethical and thus unacceptable. Building “fair” systems that respect our notion of fairness (ethical correct behaviour) requires a formal definition of fairness - such a formalism can then be used for verifying and enforcing fairness of ML and AI based systems.

As a consequence, a number of several (formal) fairness criteria has been proposed [11, 26]. A large group of criteria are concerned with the dependency between a sensitive attribute/feature and the response variable (prediction/behaviour of the system) - these criteria belong to group-based fairness criteria because they do not care about individuals but focus on whole groups of individuals only. While they all share the idea that the prediction of the system should “not depend” on the sensitive attribute, they differ in the definition of “independency”. Prominent examples of these kinds of fairness criteria are [11] demographic parity, equalized odds and predictive rate parity. However, it was shown that there does not exist a perfect criteria that can not be exploited (e.g. finding a setting in which the particular criteria is satisfied but an obvious unfairness still exists), and many criteria even contradict each other - it is impossible for some sets of criteria to be all satisfied at the same time.

Another very intuitive formalization of fairness is individual fairness [13]. The idea behind individual fairness is to “treat similar individuals similar” - which is considered to be very intuitive and similar to the concept of individual fairness from other scientific disciplines [8, 13]. Despite its appealing simple intuition, a major problem of individual fairness is to properly formalize the notion of individuality - i.e. given two individuals we have to be able to compute a score

that tells us how similar these individuals are. In this work we will focus on individual fairness.

Related work While fairness and robustness are known to be closely related to each other [12, 28], robustness of explanations has only been recently started to be investigated. For instance it was shown recently that explanation methods are also vulnerable to adversarial attacks [4, 19] - i.e. an explanation (e.g. a saliency map or feature importances) can be (arbitrarily) changed by applying small perturbations to the original sample which is going to be explained. This instability of explanations also holds true for counterfactual explanations [23] and the necessity of local stability of explanations is widely accepted [3, 4, 19, 23, 24]. However, computing stable and robust counterfactual explanations is still an open-research problem [38].

Fairness and explanations are also closely related to each other: For instance, counterfactual explanations can be used for detecting bias and unfairness in decision making systems [35]. Given the complex meaning and definition of fairness, the authors of [36] propose to use explanations methods for explaining the (un-) fairness of a model to a lay person. Explanations methods like counterfactual explanations can also be used for defining a fairness criteria like it was done in case of counterfactual fairness [22] in which a decision making system is considered to be fair if changing the sensitive attribute while holding everything else that is not causally dependent on the sensitive attribute (under a causal model) constant must not change the prediction of the decision making system.

Missing stability and robustness can lead to unfair explanations and thus compromise the trustworthiness of the decision making system [3, 4].

Our contributions In this work we formally and empirically study the robustness and individual fairness of counterfactual explanations and propose to use plausible instead of closest counterfactual explanations because we find evidence that the latter yields a better individual fairness.

The remainder of this work is structured as follows: First, in section 2 we briefly review counterfactual explanations and individual fairness. Next, we formally define our notion of individual fairness of contrasting explanations in section 3.1 and formally study fairness and robustness of counterfactual explanations in section 3.3 - we also propose to use plausible instead of closest counterfactual explanations for improving the individual fairness of the explanations. We empirically evaluate the individual fairness of closest and plausible counterfactual explanations in section 4. Finally, our work closes with a summary and outlook in section 5.

Note that all proofs and derivations, as well as additional plots of the experiments, can be found in the appendices A and B.

2 Foundations

2.1 Counterfactual Explanations

Counterfactual explanations (often just called counterfactuals) contrast samples by counterparts with minimum change of the appearance but different class label [27, 39] and can be formalized as follows:

Definition 1 (Counterfactual explanation [39]). *Assume a prediction function $h : \mathbb{R}^d \rightarrow \mathcal{Y}$ is given. Computing a counterfactual $\mathbf{x}' \in \mathbb{R}^d$ for a given input $\mathbf{x} \in \mathbb{R}^d$ is phrased as an optimization problem:*

$$\arg \min_{\mathbf{x}' \in \mathbb{R}^d} \ell(h(\mathbf{x}'), y') + C \cdot \theta(\mathbf{x}', \mathbf{x}) \quad (1)$$

where $\ell(\cdot)$ denotes a suitable loss function, y' the requested prediction, and $\theta(\cdot)$ a penalty term for deviations of \mathbf{x}' from the original input \mathbf{x} . $C > 0$ denotes the regularization strength.

In this work we assume that data come from a real-vector space and continuous optimization is possible. In this context [6], two common regularizations are the weighted Manhattan distance and the generalized L2 distance. Depending on the model and the choice of $\ell(\cdot)$ and $\theta(\cdot)$, the final optimization problem might be differentiable or not. If it is differentiable, we can use a gradient-based optimization algorithm like conjugate gradients, gradient descent or (L-)BFGS. Otherwise, we have to use a black-box optimization algorithm for continuous optimization like Downhill-Simplex method.

While the formalization of the optimization problem Eq. (1) is model agnostic (i.e. it does not make any assumptions on the model h), it can be beneficial to rewrite the optimization problem Eq. (1) in constraint form [6]:

$$\arg \min_{\mathbf{x}' \in \mathbb{R}^d} \theta(\mathbf{x}', \mathbf{x}) \quad (2a)$$

$$\text{s.t. } h(\mathbf{x}') = y' \quad (2b)$$

The authors of [6] have shown that the constraint optimization problem Eq. (2) can be turned (or efficiently approximated) into convex programs for many standard machine learning models like GLM, QDA, LVQ, etc.. Since convex programs can be solved quite efficiently [9], the constraint form [6] becomes superior over the original black-box modelling [39] if we have access to the underlying model h .

Counterfactuals stated in its simplest form, like in Definition 1 (also called closest counterfactuals), are very similar to adversarial examples, since there are no guarantees that the resulting counterfactual is plausible and feasible in the data domain. As a consequence, the absence of such constraints often leads to counterfactual explanations that are not plausible [7, 25, 30]. To overcome this problem, the several approaches [7, 25] propose to allow only those samples that lie on the data manifold - e.g. by enforcing a lower threshold $\delta > 0$ for their

probability/density. In particular, the authors of [7] build upon the constraint form Eq. (2) and propose the following extension of Eq. (2) for computing plausible counterfactuals:

$$\arg \min_{\mathbf{x}' \in \mathbb{R}^d} \theta(\mathbf{x}', \mathbf{x}) \quad (3a)$$

$$\text{s.t. } h(\mathbf{x}') = y' \quad (3b)$$

$$\hat{p}_y(\mathbf{x}') \geq \delta \quad (3c)$$

where $\hat{p}_y(\cdot)$ denotes a class dependent density estimator. Because the true density is usually not known, they further propose to replace the density constraint Eq. (3c) with an approximation of a Gaussian mixture model (GMM) which then can be written a set of convex quadratic constraints and hence nicely fits into the work of [6] for using convex programming for computing counterfactual explanations.

2.2 Individual Fairness

Individual fairness requires to “treat similar individuals similar” [13]. Transferring this idea to the ML world where we have a make predictions $h : \mathcal{X} \rightarrow \mathcal{Y}$, we can formalize individual fairness as follows:

$$d(\mathbf{x}_1, \mathbf{x}_2) \leq \epsilon \implies \Delta(h(\mathbf{x}_1), h(\mathbf{x}_2)) \leq \epsilon \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X} \quad (4)$$

where $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ denotes a similarity measure on the individuals \mathcal{X} , $\epsilon > 0$ denotes a threshold up to which we consider two individuals as similar and $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ denotes a similarity measure on the predictions of the ML systems h . A critical choice, which highly depends on the specific use-case, are the similarity measures $d(\cdot)$ and $\Delta(\cdot)$. A very simple possible default choice is to use the p-norm, which can be enhanced with some kind of feature weighting, - i.e. a real valued vector space as a representation of the individuals is assumed.

In this work we use a p-norm for measuring similarity of individuals as well as the similarity of two predictions.

3 Fairness and Robustness of Counterfactual Explanations

In the sub sequel we formally study and define individual fairness and robustness of counterfactual explanations for general as well as specific prediction functions $h : \mathcal{X} \rightarrow \mathcal{Y}$.

3.1 Individual Fairness of Counterfactual explanations

We aim for a formalization of individual fairness of counterfactual explanations. Inspired by the intuition (and formalization) of individual fairness in section 2.2, we propose the following definition that formalizes the intuition that counterfactual explanations of similar individuals should be similar:

Definition 2 (Individual fairness of counterfactual explanations). *Let $h : \mathcal{X} \rightarrow \mathcal{Y}$ be a prediction function and $(\mathbf{x}_{orig}, y_{orig}) \in \mathcal{X} \times \mathcal{Y}$ with $h(\mathbf{x}_{orig}) = y_{orig}$ be a sample prediction that has to be explained. For this purpose, let $(\mathbf{x}', y') \in \mathcal{X} \times \mathcal{Y}$ be a counterfactual explanations of $(\mathbf{x}_{orig}, y_{orig})$.*

Let $\mathbf{x} \sim \mathcal{B}_\epsilon(\mathbf{x}_{orig})$ be a randomly perturbed sample of \mathbf{x}_{orig} with $h(\mathbf{x}) = y_{orig}$ that is ϵ close to \mathbf{x}_{orig} - i.e. $d(\mathbf{x}_{orig}, \mathbf{x}) \leq \epsilon$ for some suitable metric $d(\cdot)$.

Let $(\mathbf{x}'', y'') \in \mathcal{X} \times \mathcal{Y}$ be a counterfactual explanations of this perturbed sample (\mathbf{x}, y_{orig}) with the same target label y' .

We define the individual fairness of the explanation (\mathbf{x}', y') as the expected distance between the counterfactual explanations of the original sample \mathbf{x}_{orig} and a perturbed sample \mathbf{x} :

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{B}_\epsilon(\mathbf{x}_{orig})} [d(\mathbf{x}', \mathbf{x}'')] \quad (5)$$

Remark 1. By replacing the expectation in Eq. (5) with the sample mean gives us a method for empirically comparing the individual fairness (according to Definition 2) of different counterfactual explanations.

Note that the fairness criteria Eq. (5) is to be minimized - i.e. smaller values correspond to a better individual fairness.

3.2 Perturbations

While there are infinitely many possible ways of perturbing a given input \mathbf{x} - i.e. choosing $\mathcal{B}_\epsilon(\cdot)$ in Definition 2 -, we focus on two specific perturbations in this work: Perturbation by Gaussian noise and a perturbation by masking features. Perturbing a given input \mathbf{x} with Gaussian noise means to add a small amount of normally distributed noise $\boldsymbol{\delta}$ to \mathbf{x} :

$$\mathbf{x} = \mathbf{x} + \boldsymbol{\delta} \quad \text{where } \boldsymbol{\delta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \quad (6)$$

where the size and shape of the perturbation can be controlled by the covariance matrix $\boldsymbol{\Sigma}$ - in this work we use a diagonal matrix and often choose $\boldsymbol{\Sigma} = \mathbb{I}$. However, note that in this perturbation we can not guarantee that the perturbed sample is ϵ close because $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ can yield arbitrarily large values although this is kind of unlikely.

While a perturbation by Gaussian noise Eq. (6) potentially changes every feature, feature masking allows a much more precise way of perturbing a given input:

$$\mathbf{x} = \mathbf{x} \odot \mathbf{m} \quad \text{where } \mathbf{m} \in \{0, 1\}^d \quad (7)$$

where \odot denotes the element wise multiplication and the size of the perturbation can be controlled by the number of masked features. Also note that the number of 0s (number of masked features) as well as their position (feature id) can vary - in this work: given a fixed number of masked features, we select the masked features randomly.

3.3 Robustness of Counterfactual Explanations

In this section we formally study robustness and fairness of different prediction functions h . First, we give a very general bound on the robustness of closest counterfactual explanations in Theorem 1.

Theorem 1 (General bound on closest counterfactuals of perturbed samples). *Let $h : \mathcal{X} \rightarrow \mathcal{Y}$ be a prediction function and let $(\mathbf{x}_{orig}, y_{orig}) \in \mathcal{X} \times \mathcal{Y}$ be a sample for which we are given a closest counterfactual (see Eq. (2)) $(\mathbf{x}', y') \in \mathcal{X} \times \mathcal{Y}$. Let $\mathbf{x} \sim \mathcal{B}_\epsilon(\mathbf{x}_{orig})$ be a perturbed version of \mathbf{x}_{orig} such that $\|\mathbf{x}_{orig} - \mathbf{x}\|_2 \leq \epsilon$ - we denote the corresponding closest counterfactual of \mathbf{x} with the same target prediction y' as \mathbf{x}'' .*

We can then bound the difference between the two counterfactuals \mathbf{x}' and \mathbf{x}'' as follows:

$$\|\mathbf{x}' - \mathbf{x}''\|_2 \leq 2\epsilon + 2\|\mathbf{x}_{orig} - \mathbf{x}'\|_2 \quad (8)$$

In case of a binary linear classifier, we can refine the bound from Theorem 1 as stated in Corollary 1.

Corollary 1 (Bound on closest counterfactuals of perturbed samples for binary linear classifiers). *In case of a binary linear classifier - i.e. $h(\mathbf{x}) = \text{sgd}(\mathbf{w}^\top \mathbf{x})$ -, we can refine the bound from Theorem 1 as follows:*

$$\|\mathbf{x}' - \mathbf{x}''\|_2 \leq 2\epsilon + 2|\mathbf{w}^\top \mathbf{x}_{orig}| \quad (9)$$

Remark 2. Note that while the general bound Eq. (8) in Theorem 1 depends on the closest counterfactual of the unperturbed sample \mathbf{x}_{orig} , the bound Eq. (9) in Corollary 1 only depends on the original sample \mathbf{x}_{orig} , the model parameter \mathbf{w} and the perturbation bound ϵ . Both bounds (Theorem 1 and Corollary 1) are rather loose because they do not make any assumption on the perturbation \mathcal{B}_ϵ except that it must be bounded by ϵ - in addition the bound in Theorem 1 does not even make any assumption on h at all.

Making additional assumptions allow us to come up with more precise (and potentially more useful) statements as shown in the next theorem Theorem 2. In Theorem 2 (and the consequential corollaries Corollary 2 and Corollary 3) we study the individual fairness of a linear binary classifier under Gaussian noise. It turns out, that the individual fairness Definition 2 of closest counterfactual explanations of a binary linear classifier under Gaussian noise depends on the dimension d only - i.e. the larger the dimension of the input space, the larger the individual unfairness.

Theorem 2 (Individual fairness of closest counterfactuals of a linear binary classifier under Gaussian noise). *Let $h : \mathbb{R}^d \rightarrow \{-1, 1\}$ be a binary linear classifier - i.e. $h(\mathbf{x}) = \text{sgd}(\mathbf{w}^\top \mathbf{x})$. The individual fairness of closest counterfactuals (see Definition 2) under Gaussian noise Eq. (6) - with an arbitrary diagonal covariance $\Sigma = \text{diag}(\sigma_i^2)$ - at an arbitrary (correctly classified) sample $(\mathbf{x}_{orig}, y_{orig}) \in \mathbb{R}^d \times \{-1, 1\}$ can be stated as follows:*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}_{orig}, \Sigma)} [d(\mathbf{x}', \mathbf{x}'')] = \text{trace}(\Sigma) - \mathbf{w}^\top \Sigma \mathbf{w} \quad (10)$$

where we assume the squared Euclidean distance as a distance metric $d(\cdot)$ for measuring the distance between two counterfactuals.

Corollary 2. *If we assume the identity matrix \mathbb{I} as a covariance matrix Σ of the Gaussian noise in Theorem 2, Eq. (10) simplifies as follows:*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}_{\text{orig}}, \mathbb{I})} [d(\mathbf{x}', \mathbf{x}'')] = d - 1 \quad (11)$$

Corollary 3. *Theorem 2 and Corollary 2 imply the following upper bound on the probability that the individual unfairness is larger than some $\delta > 0$:*

$$\mathbb{P}\left(d(\mathbf{x}', \mathbf{x}'') \geq \delta\right) \leq \frac{d-1}{\delta} \quad (12)$$

Remark 3. We can interpret Theorem 2 (and in particular the consequential correlaries Corollary 2 and Corollary 3) as the “curse of dimensionality for individual fairness of closest counterfactual explanations” because the larger the dimension d of the data space, the larger the individual unfairness.

We can still make some statements on the individual fairness Definition 2 of closest counterfactual explanations, when using bounded uniform noise instead of Gaussian noise, as stated in Theorem 3 and Corollary 4.

Theorem 3 (Individual fairness of closest counterfactuals of a linear binary classifier under bounded uniform noise). *Let $h : \mathbb{R}^d \rightarrow \{-1, 1\}$ be a binary linear classifier - i.e. $h(\mathbf{x}) = \text{sgd}(\mathbf{w}^\top \mathbf{x})$. The individual fairness of closest counterfactuals (see Definition 2) under a bounded uniform noise at an arbitrary (correctly classified) sample $(\mathbf{x}_{\text{orig}}, y_{\text{orig}}) \in \mathbb{R}^d \times \{-1, 1\}$ can be stated follows:*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{U}(\mathbf{x}_{\text{orig}} \pm \epsilon \mathbf{1})} [d(\mathbf{x}', \mathbf{x}'')] = \frac{\epsilon^2(d-1)}{3} \quad (13)$$

Corollary 4. *Theorem 3 implies the following upper bound on the probability that the individual unfairness is larger than some $\delta > 0$:*

$$\mathbb{P}\left(d(\mathbf{x}', \mathbf{x}'') \geq \delta\right) \leq \frac{\delta \epsilon^2(d-1)}{3} \quad (14)$$

Since our formal statements so far suggest a potential presence of unfairness even for simple models like binary linear classifiers - we will empirically confirm this in the experiments (see section 4) -, we propose to add some kind of regularization for improving the individual fairness Definition 2 of counterfactual explanations. We propose to use plausible instead of closest counterfactual because we think that the problem of individual unfairness of closest counterfactuals comes from the fact that in case of a “wiggly” decision boundary small perturbations of the input cause a completely different closest counterfactual (similar to adversarial attacks). Under the assumption that the set of plausible counterfactuals is less “wiggly” we would expect to observe better individual fairness when considering plausible instead of closest counterfactuals. We empirically evaluate this hypothesis in section 4.

4 Experiments

We empirically evaluate the individual fairness (Definition 2) of closest and plausible counterfactual explanations. For this purpose, we compute closest and plausible counterfactuals of perturbed data points for a diverse set of classifiers and data sets.

Data sets We use the three standard data sets:

- The “Breast Cancer Wisconsin (Diagnostic) Data Set” [41] whereby we add a PCA dimensionality reduction to 5 dimensions to the model.
- The “Wine data set” [33].
- The “Optical Recognition of Handwritten Digits Data Set” [2] whereby we add a PCA dimensionality reduction to 40 dimensions to the model.

Because PCA preprocessing is an affine transformation, we can integrate the transformation into the convex programs and therefore still compute counterfactuals in the original data space [6].

Models We use the following diverse set of models: softmax regression, generalized learning vector quantization (GLVQ) and decision tree classifier. We use the same hyperparameters across all data sets - for all vector quantization models we use 3 prototypes per class and use 7 as the maximum depth of decision tree classifiers.

Setup We report the results of the following experiments over a 4-fold cross validation: We fit all models on the training data (depending on the data set this might involve a PCA as a preprocessing) and compute a closest and plausible counterfactual explanations of all samples from the test set that are classified correctly by the model - whereby we compute counterfactuals of the original as well as the perturbed sample. We use two different types of perturbations: Gaussian noise Eq. (6) with $\Sigma = \mathbb{I}$ and feature masking Eq. (7) for one up to half of the total number of features. In case of a multi-class problem, we chose a random target label that is different from the original label. We compute and report the distance between the counterfactuals of the original sample and the perturbed sample Eq. (5) - we do this separately for closest and plausible counterfactuals. Furthermore, we use MOSEK¹ as a solver for all mathematical programs. The complete implementation of the experiments is available on GitHub².

Results The results of using Gaussian noise Eq. (6) for perturbing the samples are shown in Table 1. The results on the digit data set for increasingly masking more and more features Eq. (7) are shown in Fig. 1 - plots for the other data sets are given in appendix B.

We observe that in all cases the plausible counterfactual explanations are less affected by perturbations than the closest counterfactuals - thus we consider them

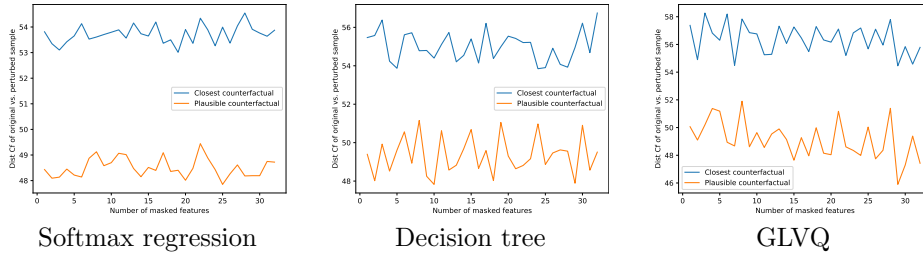
¹ We gratefully acknowledge an academic license provided by MOSEK ApS.

² <https://github.com/andreArtelt/FairnessRobustnessContrastingExplanations>

<i>Data set</i>	Wine		Breast cancer		Handwritten digits	
<i>Method</i>	Closest	Plausible	Closest	Plausible	Closest	Plausible
Softmax regression	10.16	1.87	24.04	22.48	53.71	48.78
Decision tree	9.25	2.42	24.05	23.11	56.56	49.40
GLVQ	9.95	1.74	23.34	21.42	57.66	49.46

Table 1: Comparing the median absolute distance between counterfactual of original sample and perturbed sample (using Gaussian noise Eq. (6)) - closest and plausible counterfactual explanations. Smaller values are better - best values are **highlighted**.

Fig. 1: Handwritten digits data set: Median absolute distance between counterfactual of original sample and perturbed sample (using feature masking Eq. (7)) for closest and plausible counterfactual explanations - for different number of masked features. Smaller values are better.



to be better under individual fairness (Definition 2). The size of the differences depends a lot on the combination of model and data set. However, in all cases the difference is significant. In case of increasingly masking features, we observe that although the distance between counterfactuals of original and perturbed sample is subject to some variance, the difference between closest and plausible counterfactual is always significant - even when masking up to 50% of all features.

5 Discussion and Conclusion

In this work we argued that not only the fairness of decision making systems is important but also fairness of explanations is important. We studied the robustness of contrasting explanations - in particular we focused on counterfactual explanations. Besides deriving robustness bounds, we also focused on individual fairness of contrasting explanations - we studied formally and empirically the individual fairness of counterfactual explanations. In addition, we proposed to use plausible instead of closest counterfactuals for increasing the individual fairness of counterfactual explanations - we empirically evaluated and compared the individual fairness of closest vs. plausible counterfactual explanations. We found

evidence that plausible counterfactuals provide better individual fairness than closest counterfactual explanations.

In future work we plan to further study formal fairness and robustness guarantees and bounds of more models and different perturbations. We also would like to investigate other approaches and methodologies for computing plausible counterfactual explanations - the work [7] we used in this work is only one possible approach for computing plausible counterfactuals, other approaches exist as well [25, 30]. Finally, we are highly interested in studying the problem of individual fairness of (contrasting) explanations from a psychological perspective - i.e. investigating how people actually experience individual fairness of contrasting explanations and whether this experience is successfully captured/modelled by our proposed formalizations and methods.

References

1. Aamodt, A., Plaza, E.: Case-based reasoning: Foundational issues, methodological variations, and system approaches. *AI communications* (1994)
2. Alpaydin, E., Kaynak, C.: Optical recognition of handwritten digits data set. <https://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits> (1998)
3. Alvarez-Melis, D., Jaakkola, T.S.: Towards robust interpretability with self-explaining neural networks. In: Bengio, S., Wallach, H.M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*. pp. 7786–7795 (2018), <https://proceedings.neurips.cc/paper/2018/hash/3e9f0fc9b2f89e043bc6233994dfcf76-Abstract.html>
4. Anders, C.J., Pasliev, P., Dombrowski, A., Müller, K., Kessel, P.: Fairwashing explanations with off-manifold detergent. In: *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event. Proceedings of Machine Learning Research*, vol. 119, pp. 314–323. PMLR (2020), <http://proceedings.mlr.press/v119/anders20a.html>
5. Angwin, J., Larson, J., Mattu, S., Kirchner, L.: Machine bias - theres software used across the country to predict future criminals. and its biased against blacks. (2016), <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
6. Artelt, A., Hammer, B.: On the computation of counterfactual explanations - A survey. *CoRR* **abs/1911.07749** (2019), <http://arxiv.org/abs/1911.07749>
7. Artelt, A., Hammer, B.: Convex density constraints for computing plausible counterfactual explanations. *29th International Conference on Artificial Neural Networks (ICANN)* (2020)
8. Binns, R.: On the apparent conflict between individual and group fairness. In: Hildebrandt, M., Castillo, C., Celis, E., Ruggieri, S., Taylor, L., Zanfir-Fortuna, G. (eds.) *FAT* '20: Conference on Fairness, Accountability, and Transparency*, Barcelona, Spain, January 27-30, 2020. pp. 514–524. ACM (2020). <https://doi.org/10.1145/3351095.3372864>, <https://doi.org/10.1145/3351095.3372864>

9. Boyd, S., Vandenberghe, L.: *Convex Optimization*. Cambridge University Press, New York, NY, USA (2004)
10. Byrne, R.M.J.: Counterfactuals in explainable artificial intelligence (xai): Evidence from human reasoning. In: *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. pp. 6276–6282. International Joint Conferences on Artificial Intelligence Organization (7 2019). <https://doi.org/10.24963/ijcai.2019/876>, <https://doi.org/10.24963/ijcai.2019/876>
11. Caton, S., Haas, C.: Fairness in machine learning: A survey. *CoRR* **abs/2010.04053** (2020), <https://arxiv.org/abs/2010.04053>
12. Chang, H., Nguyen, T.D., Murakonda, S.K., Kazemi, E., Shokri, R.: On adversarial bias and the robustness of fair machine learning. *CoRR* **abs/2006.08669** (2020), <https://arxiv.org/abs/2006.08669>
13. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.S.: Fairness through awareness. In: Goldwasser, S. (ed.) *Innovations in Theoretical Computer Science 2012*, Cambridge, MA, USA, January 8-10, 2012. pp. 214–226. ACM (2012). <https://doi.org/10.1145/2090236.2090255>, <https://doi.org/10.1145/2090236.2090255>
14. Fisher, A., Rudin, C., Dominici, F.: All Models are Wrong but many are Useful: Variable Importance for Black-Box, Proprietary, or Misspecified Prediction Models, using Model Class Reliance. *arXiv e-prints* arXiv:1801.01489 (Jan 2018)
15. Gilpin, L.H., Bau, D., Yuan, B.Z., Bajwa, A., Specter, M., Kagal, L.: Explaining explanations: An overview of interpretability of machine learning. In: *5th IEEE International Conference on Data Science and Advanced Analytics, DSAA 2018, Turin, Italy, October 1-3, 2018*. pp. 80–89 (2018). <https://doi.org/10.1109/DSAA.2018.00018>, <https://doi.org/10.1109/DSAA.2018.00018>
16. Greenwell, B.M., Boehmke, B.C., McCarthy, A.J.: A simple and effective model-based variable importance measure. *CoRR* **abs/1805.04755** (2018), <http://arxiv.org/abs/1805.04755>
17. Guidotti, R., Monreale, A., Ruggieri, S., Pedreschi, D., Turini, F., Giannotti, F.: Local rule-based explanations of black box decision systems. *CoRR* **abs/1805.10820** (2018), <http://arxiv.org/abs/1805.10820>
18. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D.: A survey of methods for explaining black box models. *ACM Comput. Surv.* **51**(5), 93:1–93:42 (Aug 2018). <https://doi.org/10.1145/3236009>, <http://doi.acm.org/10.1145/3236009>
19. Heo, J., Joo, S., Moon, T.: Fooling neural network interpretations via adversarial model manipulation. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E.B., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. pp. 2921–2932 (2019), <https://proceedings.neurips.cc/paper/2019/hash/7fea637fd6d02b8f0adf6f7dc36aed93-Abstract.html>
20. Kim, B., Koyejo, O., Khanna, R.: Examples are not enough, learn to criticize! criticism for interpretability. In: *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*. pp. 2280–2288 (2016)
21. Koh, P.W., Liang, P.: Understanding black-box predictions via influence functions. In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. pp. 1885–1894 (2017), <http://proceedings.mlr.press/v70/koh17a.html>

22. Kusner, M.J., Loftus, J.R., Russell, C., Silva, R.: Counterfactual fairness. In: Guyon, I., von Luxburg, U., Bengio, S., Wallach, H.M., Fergus, R., Vishwanathan, S.V.N., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, December 4-9, 2017, Long Beach, CA, USA. pp. 4066–4076 (2017), <https://proceedings.neurips.cc/paper/2017/hash/a486cd07e4ac3d270571622f4f316ec5-Abstract.html>
23. Laugel, T., Lesot, M., Marsala, C., Detyniecki, M.: Issues with post-hoc counterfactual explanations: a discussion. *CoRR* **abs/1906.04774** (2019), <http://arxiv.org/abs/1906.04774>
24. Laugel, T., Renard, X., Lesot, M., Marsala, C., Detyniecki, M.: Defining locality for surrogates in post-hoc interpretability. *CoRR* **abs/1806.07498** (2018), <http://arxiv.org/abs/1806.07498>
25. Looveren, A., Klaise, J.: Interpretable counterfactual explanations guided by prototypes. *CoRR* **abs/1907.02584** (2019)
26. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. *CoRR* **abs/1908.09635** (2019), <http://arxiv.org/abs/1908.09635>
27. Molnar, C.: *Interpretable Machine Learning* (2019), <https://christophm.github.io/interpretable-ml-book/>
28. Nanda, V., Dooley, S., Singla, S., Feizi, S., Dicker-son, J.P.: Fairness through robustness: investigating robustness disparity in deep learning. In: *FAccT’21: Conference on Fairness, Accountability, and Transparency*, Virtual Event, Canada, March 310, 2021. ACM (2021). <https://doi.org/10.1145/3442188.3445910>, <https://doi.org/10.1145/3442188.3445910>
29. parliament, E., council: Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (general data protection regulation). <https://eur-lex.europa.eu/eli/reg/2016/679/oj> (2016)
30. Poyiadzi, R., Sokol, K., Santos-Rodriguez, R., Bie, T.D., Flach, P.A.: FACE: feasible and actionable counterfactual explanations. *CoRR* **abs/1909.09369** (2019), <http://arxiv.org/abs/1909.09369>
31. Ribeiro, M.T., Singh, S., Guestrin, C.: Model-agnostic interpretability of machine learning. In: *ICML Workshop on Human Interpretability in Machine Learning (WHI)* (2016)
32. Ribeiro, M.T., Singh, S., Guestrin, C.: ”why should i trust you?”: Explaining the predictions of any classifier. In: *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. pp. 1135–1144. KDD ’16, ACM, New York, NY, USA (2016). <https://doi.org/10.1145/2939672.2939778>, <http://doi.acm.org/10.1145/2939672.2939778>
33. S. Aeberhard, D.C., de Vel, O.: Comparison of classifiers in high dimensional settings. Tech. Rep. no. 92-02 (1992)
34. Samek, W., Wiegand, T., Müller, K.: Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *CoRR* **abs/1708.08296** (2017), <http://arxiv.org/abs/1708.08296>
35. Sokol, K., Flach, P.A.: Counterfactual explanations of machine learning predictions: Opportunities and challenges for AI safety. In: Espinoza, H., hÉigeartaigh, S.Ó., Huang, X., Hernández-Orallo, J., Castillo-Effen, M. (eds.) *Workshop on Artificial Intelligence Safety 2019 co-located with the Thirty-Third AAAI Conference on Artificial Intelligence 2019 (AAAI-19)*, Honolulu, Hawaii, January

- 27, 2019. CEUR Workshop Proceedings, vol. 2301. CEUR-WS.org (2019), http://ceur-ws.org/Vol-2301/paper_20.pdf
36. Stevens, A., Deruyck, P., Veldhoven, Z.V., Vanthienen, J.: Explainability and fairness in machine learning: Improve fair end-to-end lending for kiva. In: 2020 IEEE Symposium Series on Computational Intelligence, SSCI 2020, Canberra, Australia, December 1-4, 2020. pp. 1241–1248. IEEE (2020). <https://doi.org/10.1109/SSCI47803.2020.9308371>, <https://doi.org/10.1109/SSCI47803.2020.9308371>
37. Tjoa, E., Guan, C.: A survey on explainable artificial intelligence (XAI): towards medical XAI. CoRR **abs/1907.07374** (2019), <http://arxiv.org/abs/1907.07374>
38. Verma, S., Dickerson, J., Hines, K.: Counterfactual explanations for machine learning: A review (2020)
39. Wachter, S., Mittelstadt, B.D., Russell, C.: Counterfactual explanations without opening the black box: Automated decisions and the GDPR. CoRR **abs/1711.00399** (2017), <http://arxiv.org/abs/1711.00399>
40. Waddell, K.: How algorithms can bring down minorities’ credit scores. The Atlantic (2016)
41. William H. Wolberg, W. Nick Street, O.L.M.: Breast cancer wisconsin (diagnostic) data set. [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)) (1995)
42. Zhao, Q., Hastie, T.: Causal interpretations of black-box models. Journal of Business & Economic Statistics **0**(ja), 1–19 (2019). <https://doi.org/10.1080/07350015.2019.1624293>, <https://doi.org/10.1080/07350015.2019.1624293>

A Proofs and Derivations

1. *Proof (Theorem 1).* Since the perturbation is bounded by $\epsilon > 0$, it holds that:

$$\|\mathbf{x}_{\text{orig}} - \mathbf{x}\|_2 \leq \epsilon \quad (15)$$

Furthermore, if the closest counterfactual \mathbf{x}' of \mathbf{x}_{orig} is different from the closest counterfactual \mathbf{x}'' of \mathbf{x} (perturbed \mathbf{x}_{orig}), it must hold that:

$$\|\mathbf{x} - \mathbf{x}''\|_2 \leq \|\mathbf{x} - \mathbf{x}'\|_2 \quad (16)$$

Because of the triangle inequality we know that the following holds:

$$\|\mathbf{x}'' - \mathbf{x}'\|_2 \leq \|\mathbf{x} - \mathbf{x}''\|_2 + \|\mathbf{x} - \mathbf{x}'\|_2 \quad (17)$$

Plugging Eq. (16) into Eq. (17) yields:

$$\begin{aligned} \|\mathbf{x}' - \mathbf{x}''\|_2 &\leq \|\mathbf{x} - \mathbf{x}''\|_2 + \|\mathbf{x} - \mathbf{x}'\|_2 \\ &\leq \|\mathbf{x} - \mathbf{x}'\|_2 + \|\mathbf{x} - \mathbf{x}'\|_2 \\ &= 2\|\mathbf{x} - \mathbf{x}'\|_2 \end{aligned} \quad (18)$$

By making use of the triangle inequality and Eq. (15), we find that:

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}'\|_2 &\leq \|\mathbf{x}_{\text{orig}} - \mathbf{x}\|_2 + \|\mathbf{x}_{\text{orig}} - \mathbf{x}'\|_2 \\ &\leq \epsilon + \|\mathbf{x}_{\text{orig}} - \mathbf{x}'\|_2 \end{aligned} \quad (19)$$

Plugging Eq. (19) into Eq. (18) yields the desired bound Eq. (8):

$$\begin{aligned}\|\mathbf{x}' - \mathbf{x}''\|_2 &\leq 2\|\mathbf{x} - \mathbf{x}'\|_2 \\ &= 2\epsilon + 2\|\mathbf{x}_{\text{orig}} - \mathbf{x}'\|_2\end{aligned}\quad (20)$$

□

2. *Proof (Corollary 1).* First, we prove that the closest counterfactual explanations \mathbf{x}' of a sample \mathbf{x}_{orig} under a binary linear classifier $h(\mathbf{x}) = \text{sgd}(\mathbf{w}^\top \mathbf{x})$ (we assume w.l.o.g. $\|\mathbf{w}\|_2 = 1$) can be explicitly stated as follows:

$$\mathbf{x}' = \mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}}) \mathbf{w} \quad (21)$$

Computing the closest counterfactual of some \mathbf{x}_{orig} under a binary linear classifier can be formalized as the following optimization problem:

$$\min_{\mathbf{x}' \in \mathbb{R}^d} \|\mathbf{x}_{\text{orig}} - \mathbf{x}'\|_2^2 \quad (22a)$$

$$\text{s.t. } \mathbf{w}^\top \mathbf{x}' = 0 \quad (22b)$$

Note that the constraint Eq. (22b) “replaces/approximates” the constraint $h(\mathbf{x}') = y'$ Eq. (2b). The constraint Eq. (22b) requires that the solution \mathbf{x}' lies directly on the decision boundary. We assume that points on the decision boundary are classified as y' - while this approach is debatable, it offers an easy solution to the original problem because otherwise we would have to project onto an open set which is “difficult” (once we are on the decision boundary we could add an infinitesimally small constant to the solution for crossing the decision boundary if this is really necessary).

We solve Eq. (22) by using the method of Lagrangian multipliers. Since Eq. (22) is a convex optimization problem, we only have globally optimal solutions. The Lagrangian of Eq. (22) is given as follows:

$$\mathcal{L}(\mathbf{x}', \lambda) = \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}' + \mathbf{x}'^\top \mathbf{x}' - \lambda \mathbf{w}^\top \mathbf{x}' \quad (23)$$

The gradient of the Lagrangian Eq. (23) with respect to \mathbf{x}' can be written as follows:

$$\begin{aligned}\nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}', \lambda) &= \nabla_{\mathbf{x}'} \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - \nabla_{\mathbf{x}'} 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}' + \nabla_{\mathbf{x}'} \mathbf{x}'^\top \mathbf{x}' - \nabla_{\mathbf{x}'} \lambda \mathbf{w}^\top \mathbf{x}' \\ &= -2\mathbf{x}_{\text{orig}} + 2\mathbf{x}' - \lambda \mathbf{w}\end{aligned}\quad (24)$$

The optimality condition requires the gradient Eq. (24) being equal to zero:

$$\begin{aligned}\nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}', \lambda) &= \mathbf{0} \\ \Leftrightarrow -2\mathbf{x}_{\text{orig}} + 2\mathbf{x}' - \lambda \mathbf{w} &= \mathbf{0} \\ \Leftrightarrow \mathbf{x}' &= \mathbf{x}_{\text{orig}} + \frac{\lambda}{2} \mathbf{w}\end{aligned}\quad (25)$$

Plugging Eq. (25) back into the Lagrangian Eq. (23) yields the Lagrangian dual:

$$\begin{aligned}
\mathcal{L}_D(\lambda) &= \min_{\mathbf{x}' \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}', \lambda) \\
\mathcal{L}\left(\mathbf{x}' = \mathbf{x}_{\text{orig}} + \frac{\lambda}{2}\mathbf{w}, \lambda\right) &= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \left(\mathbf{x}_{\text{orig}} + \frac{\lambda}{2}\mathbf{w}\right) + \\
&\quad \left(\mathbf{x}_{\text{orig}} + \frac{\lambda}{2}\mathbf{w}\right)^\top \left(\mathbf{x}_{\text{orig}} + \frac{\lambda}{2}\mathbf{w}\right) - \lambda \mathbf{w}^\top \left(\mathbf{x}_{\text{orig}} + \frac{\lambda}{2}\mathbf{w}\right) \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - \lambda \mathbf{x}_{\text{orig}}^\top \mathbf{w} + \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} + \\
&\quad \lambda \mathbf{x}_{\text{orig}}^\top \mathbf{w} + \frac{\lambda^2}{4} \mathbf{w}^\top \mathbf{w} - \lambda \mathbf{x}_{\text{orig}}^\top \mathbf{w} - \frac{\lambda^2}{2} \mathbf{w}^\top \mathbf{w} \\
&= \frac{\lambda^2}{4} \mathbf{w}^\top \mathbf{w} - \lambda \mathbf{x}_{\text{orig}}^\top \mathbf{w} - \frac{\lambda^2}{2} \mathbf{w}^\top \mathbf{w} \\
&= -\frac{\lambda^2}{4} \mathbf{w}^\top \mathbf{w} - \lambda \mathbf{x}_{\text{orig}}^\top \mathbf{w}
\end{aligned} \tag{26}$$

The gradient of the Lagrangian dual Eq. (26) can be written as follows:

$$\begin{aligned}
\frac{\partial}{\partial \lambda} \mathcal{L}_D(\lambda) &= -\frac{\partial}{\partial \lambda} \frac{\lambda^2}{4} \mathbf{w}^\top \mathbf{w} - \frac{\partial}{\partial \lambda} \lambda \mathbf{x}_{\text{orig}}^\top \mathbf{w} \\
&= -\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w} - \mathbf{x}_{\text{orig}}^\top \mathbf{w}
\end{aligned} \tag{27}$$

Next, the optimality condition requires that the gradient Eq. (27) is equal to zero:

$$\begin{aligned}
\frac{\partial}{\partial \lambda} \mathcal{L}_D(\lambda) &= 0 \\
\Leftrightarrow -\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w} - \mathbf{x}_{\text{orig}}^\top \mathbf{w} &= 0 \\
\Leftrightarrow \lambda &= -2\mathbf{x}_{\text{orig}}^\top \mathbf{w} = -2\mathbf{w}^\top \mathbf{x}_{\text{orig}}
\end{aligned} \tag{28}$$

where we made use of $\|\mathbf{w}\|_2 = 1$.

Finally, we obtain the solution of the original problem Eq. (22) by plugging the solution of the dual problem Eq. (28) into Eq. (25):

$$\begin{aligned}
\mathbf{x}' &= \mathbf{x}_{\text{orig}} + \frac{\lambda}{2}\mathbf{w} \\
&= \mathbf{x}_{\text{orig}} + \frac{-2\mathbf{w}^\top \mathbf{x}_{\text{orig}}}{2} \mathbf{w} \\
&= \mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}}) \mathbf{w}
\end{aligned} \tag{29}$$

which concludes this sub-proof.

Plugging Eq. (21) into the bound Eq. (8) from Theorem 1, and again assuming w.l.o.g. that $\|\mathbf{w}\|_2 = 1$, yields the desired bound Eq. (9):

$$\begin{aligned}
\|\mathbf{x}' - \mathbf{x}''\|_2 &\leq 2\epsilon + 2\|\mathbf{x}_{\text{orig}} - \mathbf{x}'\|_2 \\
&= 2\epsilon + 2\|\mathbf{x}_{\text{orig}} - (\mathbf{x}_{\text{orig}} - \mathbf{w}^\top \mathbf{x}_{\text{orig}} \mathbf{w})\|_2 \\
&= 2\epsilon + 2\|\mathbf{w}^\top \mathbf{x}_{\text{orig}} \mathbf{w}\|_2 \\
&= 2\epsilon + 2|\mathbf{w}^\top \mathbf{x}_{\text{orig}}|
\end{aligned} \tag{30}$$

□

3. *Proof (Theorem 2).* From the proof of Corollary 1 we now that the closest counterfactual explanation \mathbf{x}' of a sample \mathbf{x} under a linear binary classifier $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ can be stated explicitly Eq. (21):

$$\mathbf{x}' = \mathbf{x} - (\mathbf{w}^\top \mathbf{x}) \mathbf{w} \tag{31}$$

Applying the analytic solution Eq. (31) to the squared Euclidean distance between the closest counterfactual \mathbf{x}' of the original sample \mathbf{x}_{orig} and the closest counterfactual \mathbf{x}'' of the corresponding perturbed sample \mathbf{x} yields:

$$\begin{aligned}
d(\mathbf{x}', \mathbf{x}'') &= (\mathbf{x}' - \mathbf{x}'')^\top (\mathbf{x}' - \mathbf{x}'') \\
&= \mathbf{x}'^\top \mathbf{x}' - 2\mathbf{x}'^\top \mathbf{x}'' + \mathbf{x}''^\top \mathbf{x}'' \\
&= (\mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}}) \mathbf{w})^\top (\mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}}) \mathbf{w}) - \\
&\quad 2(\mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}}) \mathbf{w})^\top (\mathbf{x} - (\mathbf{w}^\top \mathbf{x}) \mathbf{w}) + \\
&\quad (\mathbf{x} - (\mathbf{w}^\top \mathbf{x}) \mathbf{w})^\top (\mathbf{x} - (\mathbf{w}^\top \mathbf{x}) \mathbf{w}) \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x} - 2(\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + 2\mathbf{x}_{\text{orig}}^\top (\mathbf{w}^\top \mathbf{x}) \mathbf{w} + \\
&\quad \mathbf{x}^\top \mathbf{x} + 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w}) (\mathbf{x}^\top \mathbf{w}) - 2(\mathbf{x}^\top \mathbf{w})^2 + (\mathbf{x}_{\text{orig}}^\top \mathbf{w})^2 - \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w})^\top (\mathbf{x}^\top \mathbf{w}) + (\mathbf{x}^\top \mathbf{w})^2 \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + 2(\mathbf{w}^\top \mathbf{x}) (\mathbf{x}_{\text{orig}}^\top \mathbf{w}) + \\
&\quad \mathbf{x}^\top \mathbf{x} - (\mathbf{x}^\top \mathbf{w})^2
\end{aligned} \tag{32}$$

Taking the expectation of Eq. (32) over an arbitrary density p yields:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim p} [d(\mathbf{x}', \mathbf{x}'')] &= \mathbb{E}_{\mathbf{x} \sim p} \left[\mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \right. \\
&\quad \left. 2(\mathbf{w}^\top \mathbf{x}) (\mathbf{x}_{\text{orig}}^\top \mathbf{w}) + \mathbf{x}^\top \mathbf{x} - (\mathbf{x}^\top \mathbf{w})^2 \right] \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - \mathbb{E} [2\mathbf{x}_{\text{orig}}^\top \mathbf{x}] - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w}) \mathbb{E} [\mathbf{w}^\top \mathbf{x}] + \mathbb{E} [\mathbf{x}^\top \mathbf{x}] - \mathbb{E} [(\mathbf{x}^\top \mathbf{w})^2]
\end{aligned} \tag{33}$$

Working out the specific expectations from Eq. (33) and under a Gaussian distribution $\mathbf{x} \sim \mathcal{N}(\mathbf{x}_{\text{orig}}, \mathbf{\Sigma})$ with $\mathbf{\Sigma} = \text{diag}(\sigma_i^2)$ - i.e. $(\mathbf{x})_i$ are uncorrelated -

yields:

$$\begin{aligned}
\mathbb{E}\left[2\mathbf{x}_{\text{orig}}^\top \mathbf{x}\right] &= \mathbb{E}\left[2\sum_i (\mathbf{x}_{\text{orig}})_i (\mathbf{x})_i\right] \\
&= 2\sum_i (\mathbf{x}_{\text{orig}})_i \mathbb{E}[(\mathbf{x})_i] \\
&= 2\sum_i (\mathbf{x}_{\text{orig}})_i (\mathbf{x}_{\text{orig}})_i \\
&= 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}}
\end{aligned} \tag{34}$$

$$\begin{aligned}
\mathbb{E}\left[\mathbf{w}^\top \mathbf{x}\right] &= \mathbb{E}\left[\sum_i (\mathbf{w})_i (\mathbf{x})_i\right] \\
&= \sum_i (\mathbf{w})_i \mathbb{E}[(\mathbf{x})_i] \\
&= \sum_i (\mathbf{w})_i (\mathbf{x}_{\text{orig}})_i \\
&= \mathbf{w}^\top \mathbf{x}_{\text{orig}}
\end{aligned} \tag{35}$$

$$\begin{aligned}
\mathbb{E}\left[\mathbf{x}^\top \mathbf{x}\right] &= \mathbb{E}\left[\sum_i (\mathbf{x})_i^2\right] \\
&= \sum_i \mathbb{E}[(\mathbf{x})_i^2] \\
&= \sum_i \left(\mathbb{E}[(\mathbf{x})_i]^2 + \text{Var}[(\mathbf{x})_i]\right) \\
&= \sum_i \left((\mathbf{x}_{\text{orig}})_i^2 + \sigma_i^2\right) \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} + \text{trace}(\mathbf{\Sigma})
\end{aligned} \tag{36}$$

$$\begin{aligned}
\mathbb{E}\left[(\mathbf{x}^\top \mathbf{w})^2\right] &= \mathbb{E}[\mathbf{x}^\top \mathbf{w}]^2 + \text{Var}[\mathbf{x}^\top \mathbf{w}] \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \text{Var}\left[\sum_i (\mathbf{w})_i (\mathbf{x})_i\right] \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \sum_i \text{Var}[(\mathbf{w})_i (\mathbf{x})_i] \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \sum_i (\mathbf{w})_i^2 \sigma_i^2 \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \mathbf{w}^\top \mathbf{\Sigma} \mathbf{w}
\end{aligned} \tag{37}$$

where we made use of the assumption that $\|\mathbf{w}\|_2 = 1$.

Substituting Eq. (34), Eq. (35), Eq. (36), Eq. (37) in Eq. (33) yields:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}_{\text{orig}}, \mathbf{\Sigma})} [\text{d}(\mathbf{x}', \mathbf{x}'')] &= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - \mathbb{E} [2\mathbf{x}_{\text{orig}}^\top \mathbf{x}] - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w}) \mathbb{E} [\mathbf{w}^\top \mathbf{x}] + \mathbb{E} [\mathbf{x}^\top \mathbf{x}] - \mathbb{E} [(\mathbf{x}^\top \mathbf{w})^2] \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w})^2 + \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} + \text{trace}(\mathbf{\Sigma}) - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 - \mathbf{w}^\top \mathbf{\Sigma} \mathbf{w} \\
&= \text{trace}(\mathbf{\Sigma}) - \mathbf{w}^\top \mathbf{\Sigma} \mathbf{w}
\end{aligned} \tag{38}$$

which concludes the proof. \square

4. *Proof (Corollary 2).* Substituting \mathbb{I} for $\mathbf{\Sigma}$ in Eq. (10) from Theorem 2 yields the claimed expectation:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}_{\text{orig}}, \mathbf{\Sigma})} [\text{d}(\mathbf{x}', \mathbf{x}'')] &= \text{trace}(\mathbf{\Sigma}) - \mathbf{w}^\top \mathbf{\Sigma} \mathbf{w} \\
&= \text{trace}(\mathbb{I}) - \mathbf{w}^\top \mathbb{I} \mathbf{w} \\
&= d - 1
\end{aligned} \tag{39}$$

\square

5. *Proof (Corollary 3).* Plugging the expectation from Corollary 2 into Markov's inequality yields the claimed bound:

$$\begin{aligned}
\mathbb{P}(\text{d}(\mathbf{x}', \mathbf{x}'') \geq \delta) &\leq \frac{\mathbb{E}[\text{d}(\mathbf{x}', \mathbf{x}'')]}{\delta} \\
&= \frac{d - 1}{\delta}
\end{aligned} \tag{40}$$

\square

6. *Proof (Theorem 3).* From the proof of Theorem 2 we know that the expectation over an arbitrary density p of the distance between the closest counterfactual of the original sample and the perturbed sample can be written as follows:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim p} [\text{d}(\mathbf{x}', \mathbf{x}'')] &= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - \mathbb{E} [2\mathbf{x}_{\text{orig}}^\top \mathbf{x}] - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w}) \mathbb{E} [\mathbf{w}^\top \mathbf{x}] + \mathbb{E} [\mathbf{x}^\top \mathbf{x}] - \mathbb{E} [(\mathbf{x}^\top \mathbf{w})^2]
\end{aligned} \tag{41}$$

Next, working out the specific expectations from Eq. (41) under a bounded uniform noise $\mathbf{x} \sim \mathcal{U}(\mathbf{x}_{\text{orig}} \pm \epsilon \mathbf{1})$ - i.e. $(\mathbf{x})_i$ are uncorrelated - yields:

$$\begin{aligned}
\mathbb{E}\left[2\mathbf{x}_{\text{orig}}^\top \mathbf{x}\right] &= \mathbb{E}\left[2 \sum_i (\mathbf{x}_{\text{orig}})_i (\mathbf{x})_i\right] \\
&= 2 \sum_i (\mathbf{x}_{\text{orig}})_i \mathbb{E}[(\mathbf{x})_i] \\
&= 2 \sum_i (\mathbf{x}_{\text{orig}})_i \frac{1}{2} \left((\mathbf{x}_{\text{orig}})_i - \epsilon + (\mathbf{x}_{\text{orig}})_i + \epsilon \right) \\
&= 2 \sum_i (\mathbf{x}_{\text{orig}})_i (\mathbf{x}_{\text{orig}})_i \\
&= 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}}
\end{aligned} \tag{42}$$

$$\begin{aligned}
\mathbb{E}\left[\mathbf{w}^\top \mathbf{x}\right] &= \mathbb{E}\left[\sum_i (\mathbf{w})_i (\mathbf{x})_i\right] \\
&= \sum_i (\mathbf{w})_i \mathbb{E}[(\mathbf{x})_i] \\
&= \sum_i (\mathbf{w})_i (\mathbf{x}_{\text{orig}})_i \\
&= \mathbf{w}^\top \mathbf{x}_{\text{orig}}
\end{aligned} \tag{43}$$

$$\begin{aligned}
\mathbb{E}\left[\mathbf{x}^\top \mathbf{x}\right] &= \mathbb{E}\left[\sum_i (\mathbf{x})_i^2\right] \\
&= \sum_i \mathbb{E}[(\mathbf{x})_i^2] \\
&= \sum_i \left(\mathbb{E}[(\mathbf{x})_i]^2 + \text{Var}[(\mathbf{x})_i] \right) \\
&= \sum_i \left((\mathbf{x}_{\text{orig}})_i^2 + \frac{1}{12} \left((\mathbf{x}_{\text{orig}})_i + \epsilon - (\mathbf{x}_{\text{orig}})_i + \epsilon \right)^2 \right) \\
&= \sum_i \left((\mathbf{x}_{\text{orig}})_i^2 + \frac{4\epsilon^2}{12} \right) \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} + \frac{d\epsilon^2}{3}
\end{aligned} \tag{44}$$

$$\begin{aligned}
\mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] &= \mathbb{E}[\mathbf{x}^\top \mathbf{w}]^2 + \text{Var}[\mathbf{x}^\top \mathbf{w}] \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \text{Var}\left[\sum_i (\mathbf{w})_i (\mathbf{x})_i\right] \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \sum_i \text{Var}[(\mathbf{w})_i (\mathbf{x})_i] \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \sum_i (\mathbf{w})_i^2 \frac{\epsilon^2}{3} \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \mathbf{w}^\top \mathbf{w} \frac{\epsilon^2}{3} \\
&= (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \frac{\epsilon^2}{3}
\end{aligned} \tag{45}$$

Substituting Eq. (42), Eq. (43), Eq. (44), Eq. (45) in Eq. (41) yields:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim \mathcal{U}(\mathbf{x}_{\text{orig}} \pm \epsilon \mathbf{1})}[\text{d}(\mathbf{x}', \mathbf{x}'')] &= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - \mathbb{E}[2\mathbf{x}_{\text{orig}}^\top \mathbf{x}] - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w}) \mathbb{E}[\mathbf{w}^\top \mathbf{x}] + \mathbb{E}[\mathbf{x}^\top \mathbf{x}] - \mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] \\
&= \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - 2\mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 + \\
&\quad 2(\mathbf{x}_{\text{orig}}^\top \mathbf{w})^2 + \mathbf{x}_{\text{orig}}^\top \mathbf{x}_{\text{orig}} + \frac{d\epsilon^2}{3} - (\mathbf{w}^\top \mathbf{x}_{\text{orig}})^2 - \frac{\epsilon^2}{3} \\
&= \frac{\epsilon^2(d-1)}{3}
\end{aligned} \tag{46}$$

which concludes the proof. \square

7. *Proof (Corollary 4).* Plugging the expectation from Theorem 3 into Markov's inequality yields the claimed bound:

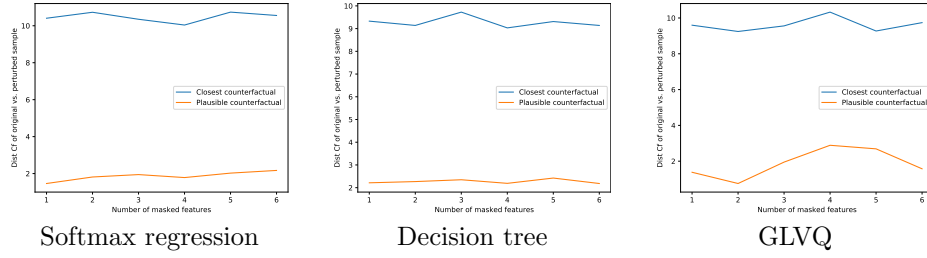
$$\begin{aligned}
\mathbb{P}(\text{d}(\mathbf{x}', \mathbf{x}'') \geq \delta) &\leq \frac{\mathbb{E}[\text{d}(\mathbf{x}', \mathbf{x}'')]}{\delta} \\
&= \frac{\delta\epsilon^2(d-1)}{3}
\end{aligned} \tag{47}$$

\square

B Additional Plots

B.1 Wine data set

Fig. 2: Wine data set: Median absolute distance between counterfactual of original sample and perturbed sample (using feature masking Eq. (7)) for closest and plausible counterfactual explanations - for different number of masked features. Smaller values are better.



B.2 Breast cancer data set

Fig. 3: Breast cancer data set: Median absolute distance between counterfactual of original sample and perturbed sample (using feature masking Eq. (7)) for closest and plausible counterfactual explanations - for different number of masked features. Smaller values are better.

