

Explaining Black-Box Models through Counterfactuals

Patrick Altmeyer¹, Arie van Deursen¹, and Cynthia C. S. Liem¹

¹Delft University of Technology

ABSTRACT

Machine Learning models like Deep Neural Networks have become so complex and opaque over recent years that they are generally considered Black Boxes. Nonetheless, such models often play a key role in modern automated decision-making systems. Counterfactual Explanations can help human stakeholders make sense of the systems they build, use and endure: they explain how inputs into a system need to change for it to produce different decisions. Explanations that involve realistic and actionable changes can be used for the purpose of Algorithmic Recourse: they offer humans a way to not only understand the behaviour of a system, but also to adjust and react to it. In this article we discuss the usefulness of Counterfactual Explanations for Explainable Machine Learning and demonstrate its implementation in Julia using the `CounterfactualExplanations.jl` package. The package is straightforward to use and designed with a focus on customization and extensibility. We envision it to one day be the go-to place for explaining arbitrary predictive models in Julia through a diverse suite of counterfactual generators.

Keywords

Julia, Explainable Artificial Intelligence, Counterfactual Explanations, Algorithmic Recourse

1. Introduction

This lack of transparency of modern machine learning models like deep neural networks exacerbates a number of other problems typically associated with them: they tend to be unstable [8], encode existing biases [6] and learn representations that are surprising or even counter-intuitive from a human perspective [29]. Nonetheless, they often form the basis for data-driven decision-making systems in real-world applications.

As others have pointed out, this scenario gives rise to an undesirable principal-agent problem involving a group of principals - i.e. human stakeholders - that fail to understand the behaviour of their agent - i.e. the black-box system [5]. The group of principals may include programmers, product managers and other decision-makers who develop and operate the system as well as those individuals ultimately subject to the decisions made by the system. In practice, decisions made by black-box systems are typically left unchallenged since the principals cannot scrutinize them:

“You cannot appeal to (algorithms). They do not listen. Nor do they bend.” [21]

In light of all this, a quickly growing body of literature on Explainable Artificial Intelligence (XAI) has emerged. Counterfactual Explanations (CE) fall into this broader category. They can help hu-

man stakeholders make sense of the systems they develop, use or endure: they explain how inputs into a system need to change for it to produce different decisions. Explainability benefits internal as well as external quality assurance. Explanations that involve realistic and actionable changes can be used for Algorithmic Recourse (AR): they offer the group of principals a way to not only understand their agent’s behaviour but also adjust or react to it.

The availability of open-source software to explain black-box models through counterfactuals is still limited. Most existing implementations are specific to particular methodologies. They are also exclusively built in Python and for Python models. The only existing unifying software approach, for example, is tailored to models built in the two most popular Python libraries for deep learning. The Julia ecosystem has so far lacked an open-source implementation of CE.

Through the work presented here, we aim to close that gap and thereby contribute to broader community efforts towards explainable AI. We envision this package to one day be the go-to place for Counterfactual Explanations in Julia. Thanks to Julia’s unique support for interoperability with foreign programming languages we believe that this library may ultimately also benefit the broader machine learning and data science community.

Our package provides a simple and intuitive interface to generate Counterfactual Explanations for differentiable classification models trained in Julia. It comes with detailed documentation involving various illustrative example datasets, linear and deep learning classifiers and counterfactual generators for binary and multi-class prediction tasks. A carefully designed package architecture allows for a seamless extension of the package functionality through custom generators and models. Through simple examples, we also demonstrate how to use our package to explain models that were built and trained in Python and R, although at the time of writing this feature is still experimental.

The remainder of this article is structured as follows: Section 2 presents related work on Explainable AI as well as a brief overview of the methodological framework underlying CE. Section 3 introduces the Julia package and its high-level architecture. Section 4 then presents a number of basic and advanced usage examples. In Section 5 we demonstrate how the package functionality can be customized and extended. To provide a flavour of its practical use, we use the package to explain models trained on MNIST data in Section 6. Finally, we also discuss the current limitations of our package, as well as its future outlook in Section 7. Section 8 concludes.

2. Background and related work

In this section, we first briefly introduce the broad field of Explainable Artificial Intelligence (XAI), before narrowing it down

to Counterfactual Explanations. We introduce the methodological framework and finally point to existing open-source software.

2.1 Literature on Explainable AI

The field of Explainable AI is still relatively young and made up of a variety of subdomains, definitions, concepts and taxonomies. Covering all of these is beyond the scope of this article, so we will focus only on high-level concepts. The following literature surveys provide more detail: Arrieta et al. (2020) provide a broad overview of XAI [3]; Fan et al. (2020) focus on explainability in the context of deep learning [7]; and finally, Karimi et al. (2020) [12] and Verma et al. (2020) [33] offer detailed reviews of the literature on Counterfactual Explanations and Algorithmic Recourse.¹ Finally, Miller (2019) explicitly discusses the concept of explainability from the perspective of a social scientist [18].

The first broad distinction we want to make here is between **interpretable** and **explainable** AI. These terms are often used interchangeably, but this can lead to confusion. We find the distinction made in [25] useful: interpretable AI involves models that are inherently interpretable and transparent such as general additive models (GAM), decision trees and rule-based models; explainable AI may involve models that are not inherently interpretable but require additional tools to be explainable to humans. Examples of the latter include ensembles, support vector machines and deep neural networks. Some would argue that we best avoid the second category of models altogether and instead focus solely on interpretable AI [25]. While we agree that initial efforts should always be geared towards interpretable models, avoiding black boxes altogether would entail missed opportunities and anyway is probably not very realistic at this point. For that reason, we expect the need for explainable AI to persist in the medium term. Explainable AI can further be broadly divided into **global** and **local** explainability: the former is concerned with explaining the average behaviour of a model, while the latter involves explanations for individual predictions [19]. Tools for global explainability include partial dependence plots (PDP), which involve the computation of marginal effects through Monte Carlo, and global surrogates. A surrogate model is an interpretable model that is trained to explain the predictions of a black-box model.

Counterfactual Explanations fall into the category of local methods: they explain how individual predictions change in response to individual feature perturbations. Among the most popular alternatives to Counterfactual Explanations are local surrogate explainers including local interpretable model-agnostic explanations (LIME) and Shapley additive explanations (SHAP). Since explanations produced by LIME and SHAP typically involve simple feature importance plots, they arguably rely on reasonably interpretable features at the very least. Contrary to Counterfactual Explanations, for example, it is not obvious how to apply LIME and SHAP to visual or audio data. Nonetheless, local surrogate explainers are among the most widely used XAI tools today, potentially because they are easily understood, relatively fast and implemented in popular programming languages. Proponents of surrogate explainers also commonly mention that there is a straightforward way to assess their reliability: a surrogate model that generates predictions in line with those produced by the black-box model is said to have high **fidelity** and therefore considered reliable. As intuitive as this notion may be, it also points to an obvious shortfall of surrogate explainers: even a high-fidelity surrogate model that produces the same predic-

tions as the black-box model 99 per cent of the time is useless and potentially misleading for every 1 out of 100 individual predictions. A recent study has shown that even experienced data scientists tend to put too much trust in explanations produced by LIME and SHAP [15]. Another recent work has shown that both LIME and SHAP can be easily fooled: both methods depend on random input perturbations, a property that can be abused by adverse agents to essentially whitewash strongly biased black-box models [28]. In a related work the same authors find that while gradient-based Counterfactual Explanations can also be manipulated, there is a straightforward way to protect against this in practice [27]. In the context of quality assessment, it is also worth noting that - contrary to surrogate explainers - Counterfactual Explanations always achieve full fidelity by construction: counterfactuals are searched with respect to the black-box classifier, not some proxy for it. That being said, Counterfactual Explanations should also be used with care and research around them is still in its early stages. We shall discuss this in more detail in the following.

2.2 A framework for Counterfactual Explanations

Counterfactual search happens in the feature space²: we are interested in understanding how we need to change individual attributes in order to change the model output to a desired value or label [19]. Typically the underlying methodology is presented in the context of binary classification: $M : \mathcal{X} \mapsto \mathcal{Y}$ where $\mathcal{X} \subset \mathbb{R}^D$ and $\mathcal{Y} = \{0, 1\}$. Further, let $t = 1$ be the target class and let x denote the factual feature vector of some individual sample outside of the target class, so $y = M(x) = 0$. We follow this convention here, though it should be noted that the ideas presented here also carry over to multi-class problems and regression [19].

The counterfactual search objective originally proposed by Wachter et al. (2017) [34] is as follows

$$\min_{x' \in \mathcal{X}} h(x') \quad \text{s. t.} \quad M(x') = t \quad (1)$$

where $h(\cdot)$ quantifies how complex or costly it is to go from the factual x to the counterfactual x' . To simplify things we can restate this constrained objective (Equation 1) as the following unconstrained and differentiable problem:

$$x' = \arg \min_x \ell(M(x'), t) + \lambda h(x') \quad (2)$$

Here ℓ denotes some loss function targeting the deviation between the target label and the predicted label and λ governs the strength of the complexity penalty. Provided we have gradient access for the black-box model M the solution to this problem (Equation 2) can be found through gradient descent. This generic framework lays the foundation for most state-of-the-art approaches to counterfactual search and is also used as the baseline approach in our package. The hyperparameter λ is typically tuned through grid search or in some sense pre-determined by the nature of the problem. Conventional choices for ℓ include margin-based losses like cross-entropy loss and hinge loss. It is worth pointing out that the loss function is typically computed with respect to logits rather than predicted probabilities, a convention that we have chosen to follow.³ Numerous - and in some cases competing - extensions to this simple approach have been developed since Counterfactual Explana-

¹Readers who prefer a text-book approach may also want to consider [19] and [32]

²Or, in the case of Latent Space generators, some latent representation of the feature space.

³While the rationale for this convention is not entirely obvious, implementations of loss functions with respect to logits are often numerically more

tions were first proposed in 2017 (see [33] and [12] for surveys). The various approaches largely differ in how they define the complexity penalty. In the baseline paper [34], for example, $h(\cdot)$ is defined in terms of the Manhattan distance between factual and counterfactual feature values. While this is an intuitive choice, it is too simple to address many of the desirable properties of effective Counterfactual Explanations that have been set out. These desiderata include: **closeness** — the average distance between factual and counterfactual features should be small [34]; **actionability** — the proposed feature perturbation should actually be actionable ([31], [24]); **plausibility** — the counterfactual explanation should be realistic plausible to a human ([10], [26]); **unambiguity** — a human should have no trouble assigning a label to the counterfactual [26]; **sparsity** — the counterfactual explanation should involve as few individual feature changes as possible [26]; **robustness** — the counterfactual explanation should be robust to domain and model shifts [30]; **diversity** — ideally multiple diverse Counterfactual Explanations should be provided [20]; and **causality** — Counterfactual Explanations should respect the structural causal model underlying the data generating process ([14],[13]).

2.3 Existing software

To the best of our knowledge, the package introduced here provides the first implementation of Counterfactual Explanations in Julia and therefore represents a novel contribution to the community. As for other programming languages, we are only aware of one other unifying framework: the recently introduced Python library CARLA [23]. In addition to that, there exists open-source code for some specific approaches to Counterfactual Explanations that have been proposed in recent years. The approach-specific implementations that we have been able to find are generally well-documented, but exclusively in Python. For example, a PyTorch implementation of a greedy generator for Bayesian models proposed in [26] has been released. As another example, the popular InterpretML library includes an implementation of a diverse counterfactual generator proposed by [20].

Generally speaking, software development in the space of XAI has largely focused on various global methods and surrogate explainers: implementations of PDP, LIME and SHAP are available for both Python (e.g. `lime`, `shap`) and R (e.g. `lime`, `iml`, `shapper`, `fastshap`). In the Julia space we have only been able to identify one package that falls into the broader scope of XAI, namely `ShapML.jl` which provides a fast implementation of SHAP.⁴ We also should not fail to mention the comprehensive Interpretable AI infrastructure, which focuses exclusively on interpretable models. Arguably the current availability of tools for explaining black-box models in Julia is limited, but it appears that the community is invested in changing that. The team behind `MLJ.jl`, for example, recruited contributors for a project about both Interpretable and Explainable AI in 2022.⁵ With our work on Counterfactual Explanations we hope to contribute to these efforts. We think that because of its unique transparency the Julia language naturally lends itself towards building a greater degree of trust in Machine Learning and Artificial Intelligence.

stable. For example, the `logitbinarycrossentropy(y, y)` implementation in `Flux.Losses` (used here) is more stable than the mathematically equivalent `binarycrossentropy(y, y)`.

⁴See here: <https://github.com/nredell/ShapML.jl>

⁵For details, see the Google Summer of Code 2022 project proposal: https://julialang.org/jsoc/gsoc/MLJ/#interpretable_machine_learning_in_julia.

3. Introducing: CounterfactualExplanations.jl

Figure 1 provides an overview of the package architecture. It is built around two core modules that are designed to be as extensible as possible through dispatch: 1) `Models` is concerned with making any arbitrary model compatible with the package; 2) `Generators` is used to implement arbitrary counterfactual search algorithms.⁶ The core function of the package `generate_counterfactual` uses an instance of type `<: AbstractFittedModel` produced by the `Models` module and an instance of type `<: AbstractGenerator` produced by the `Generators` module. Relating this to the methodology outlined in Section 2.2, the former instance corresponds to the model M , while the latter defines the rules for the counterfactual search (Equation 2).

3.1 Generators

At the time of writing the following counterfactual generators have been implemented in the package:

- Generic [34]
- Greedy [26]
- DiCE [20]
- Latent Space Search as in REVISE [10] and CLUE [2]
- ClaPROAR [1]
- Gravitational [1]

3.2 Models

The package currently offers native support for models built and trained in Flux. While in general it is assumed that users will use this package to explain their pre-trained models, we provide a simple API call to train the following simple default models:

- Linear Classifier (Logistic Regression)
- Multi-Layer Perceptron (Deep Neural Network)
- Deep Ensemble [16]

Through optional arguments, users can control the neural network architecture and impose regularization through dropout.

As we demonstrate below, it is straightforward to extend the package through custom models, although, at the time of writing, models are required to support auto-differentiation. Support for `torch` models trained in Python and R is possible, but currently still experimental.

3.3 Data and Benchmarking

To allow researchers and practitioners to test and compare counterfactual generators, the package ships with pre-processed synthetic and real-world benchmark datasets from different domains. Real-world datasets include:

- California Housing [22]
- UCI defaultCredit [35]
- Give Me Some Credit [11]
- MNIST [17]

⁶We have made an effort to keep the code base a flexible and extensible as possible, but cannot guarantee at this point that any counterfactual generator can be implemented without further adaptation.

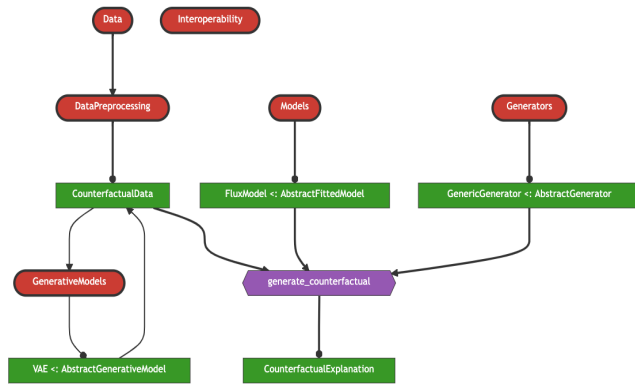


Fig. 1. High-level schematic overview of package architecture. Modules are shown in red, structs in green and functions in blue.

3.4 Plotting

The package also extends common `Plots.jl` methods to facilitate the visualization of results. Calling the generic `plot()` method on an instance of type `<:CounterfactualExplanation`, for example, generates a plot visualizing the entire counterfactual path in the feature space⁷. We will see several examples of this below.

4. Basic Usage

In the following, we begin our exploration of the package functionality with a simple example. We then turn to a more advanced usage example and show how users can impose mutability constraints on features.

4.1 A Simple Generic Generator

Listing 4.1 below provides a complete example demonstrating how the framework presented in Section 2.2 can be implemented in Julia with our package. Using a synthetic data set with linearly separable samples we first define our model and then generate a counterfactual for a randomly selected sample. Figure 2 shows the resulting counterfactual path in the two-dimensional feature space. Features go through iterative perturbations until the desired confidence level is reached as illustrated by the contour in the background, which indicates the classifier's predicted probability that the label is equal to 1.

It may help to go through the relevant parts of the code in some more detail starting from the part involving the model. For illustrative purposes the `Models` module ships with a constructor for a logistic regression model: `LogisticModel(W::Matrix, b::AbstractArray) <: AbstractFittedModel`. This constructor does not fit the regression model but rather takes its underlying parameters as given. In other words, it is generally assumed that the user has already estimated a model. Based on the provided estimates two functions are already implemented that compute logits and probabilities for the model, respectively. Below we will see how users can use `Dispatch` to extend these functions for use with arbitrary models. For now, it is enough to note that those methods define how the

⁷For multi-dimensional input data, standard dimensionality reduction techniques are used to compress the data. In this case, the classifier's decision boundary is approximated through a Nearest Neighbour model. This is still somewhat experimental and will be improved in the future.

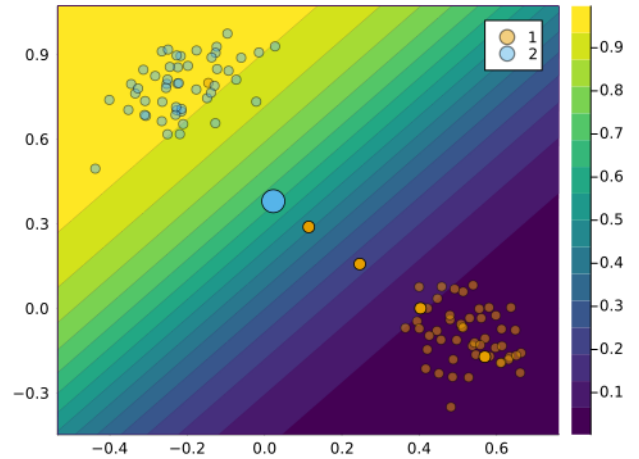


Fig. 2. Counterfactual path using generic counterfactual generator for conventional binary classifier.

model makes its predictions $M(x)$ and hence they form an integral part of the counterfactual search. With the model M defined in the code below we go on to set up the counterfactual search as follows: 1) specify the other class as our `target` label ($t = 1$) in line 7; 2) choose a random sample x from the non-target class in line 10; 3) define the counterfactual generator in line 13; and finally run the counterfactual search in line 15. Generators like the `GenericGenerator` take several optional arguments that govern the strength of the complexity penalty, the step size for gradient descent and the tolerance for convergence among other things. This will be discussed in some more detail when looking at the advanced usage example below.

```

1 # Data and Classifier:
2 counterfactual_data = load_linearly_separable()
3 M = fit_model(counterfactual_data, :Linear)
4
5 # Factual and Target:
6 yhat = predict_label(M, counterfactual_data)
7 target = 2 # target label
8 candidates = findall(vec(yhat) .!= target)
9 chosen = rand(candidates)
10 x = select_factual(counterfactual_data, chosen)
11
12 # Counterfactual search:
13 generator = GenericGenerator()
14 ce = generate_counterfactual(
15     x, target, counterfactual_data, M, generator)

```

In this simple example, the generic generator produces an effective counterfactual: the decision boundary is crossed (i.e. the counterfactual explanation is valid) and upon visual inspection, the counterfactual seems plausible (Figure 2). Still, the example also illustrates that things may well go wrong. Since the underlying model produces high-confidence predictions in regions free of any data - that is regions with high epistemic uncertainty - it is easy to think of scenarios that involve valid but unrealistic counterfactuals. Similarly, any degree of overfitting can be expected to result in more ambiguous Counterfactual Explanations, since it reduces the classifier's sensitivity to regions with high aleatoric uncertainty. Consider, for example, the scenario illustrated in Figure 3, which involves the same logistic classifier, but a massively overfitted ver-

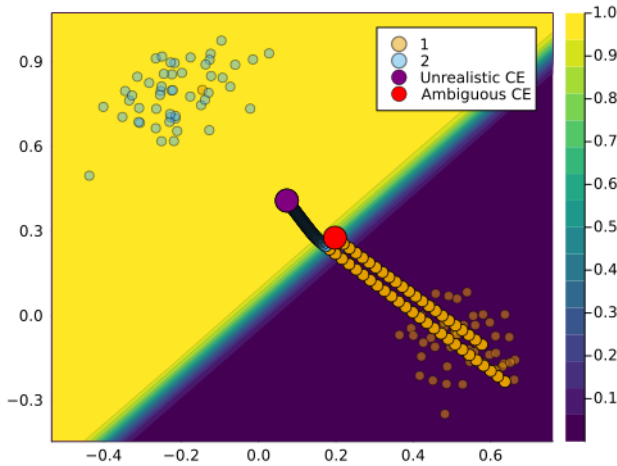


Fig. 3. Unrealistic and ambiguous counterfactuals that may be produced by generic counterfactual search for an overfitted conventional binary classifier.

sion of it. In this case, generic search may yield an entirely ambiguous counterfactual near the decision boundary (purple marker) or an unrealistic counterfactual that has moved well into the target domain, but remains far away from all other samples (red marker).

4.2 More Advanced Generators

The more advanced generators currently implemented in `CounterfactualExplanations.jl` are designed to generate more realistic counterfactuals. In this context, ‘realistic’ is defined in the sense that counterfactuals ought to be generated by the same data-generating process (DGP) that generates the actual data points. To this end, **Latent Space** generators like REVISE [10] use a separate generative model to learn the DGP. We refer to them as Latent Space generators, because they search counterfactuals in the latent embedding learned by the generative model.⁸ The **Greedy** approach [26] instead relies on minimizing predictive uncertainty in order to generate realistic counterfactuals. **CLUE** [2] can be thought of as a combination of these two ideas. The other generator currently implemented, **DiCE** [20], generates multiple counterfactuals at once that are as diverse as possible. This strategy is based on the intuition that a wide variety of diverse explanations may be suitable depending on the practical context. Listing 4.2 below shows a more advanced usage example involving the DiCE generator. Once again it is worth dwelling on this for a moment. The `DiCEGenerator` is instantiated in line 2 with a custom Flux optimizer and decision threshold specified in lines 3 and 4, respectively⁹. The main API call to generate counterfactuals is the same as before, but note that in line 10 we have specified an optional key argument that determines how many counterfactuals are generated. For the DiCE generator it naturally makes sense to generate multiple counterfactuals, but note that this is in principle also possible for all other generators.¹⁰ Figure 4 shows the result-

⁸Currently our implementation relies on a Variational Autoencoder (VAE)

⁹Note that all differentiable generators except the `GreedyGenerator` work with Flux optimizers and accept them as an optional key argument.

¹⁰By default counterfactuals are initialized by adding a small, random perturbation, as this improves adversarial robustness [27]. Therefore, gener-

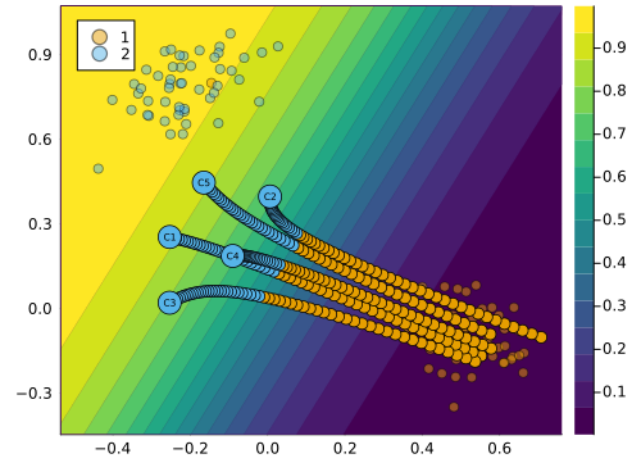


Fig. 4. Counterfactual path using the DiCE generator.

ing output. It was generated by calling the generic plot method directly on the object returned by `generate_counterfactual`.

```

1 # Generator :
2 generator = DiCEGenerator(
3     opt = Flux.Optimise.Descent(0.01),
4     decision_threshold = 0.7,
5     λ = [0.1, 5]
6 )
7 # Counterfactual search:
8 counterfactuals = generate_counterfactual(
9     x, target, counterfactual_data, M, generator;
10    num_counterfactuals=5
11 )

```

4.3 Mutability Constraints

In practice, features usually cannot be perturbed arbitrarily. Suppose, for example, that one of the features used by a bank to predict the creditworthiness of its clients is *gender*. If a counterfactual explanation for the prediction model indicates that female clients should change their gender to improve their creditworthiness, then this is an interesting insight (it reveals gender bias), but it is not usually an actionable transformation in practice. In such cases, we may want to constrain the mutability of features to ensure actionable and realistic recourse. To illustrate how this can be implemented in `CounterfactualExplanations.jl` we will look at the linearly separable toy dataset again.

Mutability of features can be defined in terms of four different options: 1) the feature is mutable in both directions, 2) the feature can only increase (e.g. *age*), 3) the feature can only decrease (e.g. *time left* until your next deadline) and 4) the feature is not mutable (e.g. *skin colour*, *ethnicity*, ...). To specify which category a feature belongs to, users can pass a vector of symbols containing the mutability constraints at the pre-processing stage. For each feature one can choose from these four options: `:both` (mutable in both directions), `:increase` (only up), `:decrease` (only down) and `:none` (immutable). By default, nothing is passed to that key-

ating multiple counterfactuals will yield multiple distinct outcomes even without an explicit diversity constraint.

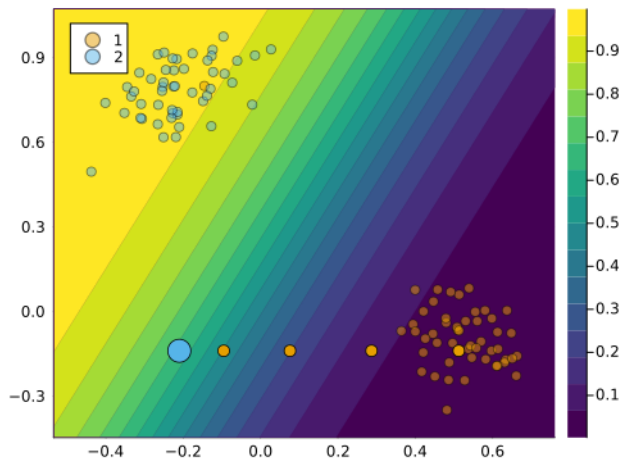


Fig. 5. Counterfactual path with immutable feature.

word argument and it is assumed that all features are mutable in both directions.¹¹

Below we impose that the second feature is immutable.

```
1 counterfactual_data.mutability = [:both, :none]
```

The resulting counterfactual path is shown in Figure 5 below. Since only the first feature can be perturbed, the sample can only move along the horizontal axis.

5. Customization and Extensibility

One of our priorities has been to make `CounterfactualExplanations` customizable and extensible. In the long term, we aim to add support for more default models and counterfactual generators. In the short term, it is designed to allow users to integrate models and generators themselves. Ideally, these community efforts will facilitate our long-term goals.

5.1 Adding Custom Models

At the high level, only two steps are necessary to make any supervised learning model compatible with our package:

Subtyping: the model needs to be declared as a subtype of `AbstractFittedModel`.

Dispatch: the functions `logits` and `probs` need to be extended through custom methods for the model in question.

To demonstrate how this can be done in practice, we will reiterate here how native support for `Flux.jl` [9] deep learning models was enabled.¹² Once again we use synthetic data for an illustrative example. Listing 5.1 below builds a simple model architecture that can be used for a multi-class prediction task. Note how outputs from the final layer are not passed through a softmax activation

function, since the counterfactual loss is evaluated with respect to logits as we discussed earlier. The model is trained with dropout for ten training epochs.

```
1 n_hidden = 32
2 output_dim = length(unique(y))
3 input_dim = 2
4 model = Chain(
5     Dense(input_dim, n_hidden, activation),
6     Dropout(0.1),
7     Dense(n_hidden, output_dim)
8 )
```

Listing 5.1 below implements the two steps that were necessary to make `Flux` models compatible with the package. In line 2 we declare our new struct as a subtype of `AbstractDifferentiableModel`, which itself is an abstract subtype of `AbstractFittedModel`.¹³ Computing logits amounts to just calling the model on inputs. Predicted probabilities for labels can then be computed by passing predicted logits through the softmax function.

```
1 # Step 1)
2 struct MyFluxModel <: AbstractDifferentiableModel
3     model::Any
4     likelihood::Symbol
5 end
6
7 # Step 2)
8 # import functions in order to extend
9 import CounterfactualExplanations.Models: logits
10 import CounterfactualExplanations.Models: probs
11 logits(M::MyFluxModel, X::AbstractArray) =
12     M.model(X)
13 probs(M::MyFluxModel, X::AbstractArray) =
14     softmax(logits(M, X))
15 M = MyFluxModel(model)
```

The API call for actually generating counterfactuals for our new model is the same as before. Figure 6 shows the resulting counterfactual path for a randomly chosen sample. In this case, the contour shows the predicted probability that the input is in the target class ($t = 2$). Generic search yields a valid, realistic and largely unambiguous counterfactual.

5.2 Adding Custom Generators

To illustrate how custom generators can be implemented we will consider a simple example of a generator that extends the functionality of our `GenericGenerator`. We have noted elsewhere that the effectiveness of Counterfactual Explanations depends to some degree on the quality of the fitted model. Another, perhaps trivial, thing to note is that Counterfactual Explanations are not unique: there are potentially many valid counterfactual paths. One idea building on these two observations might be to introduce some form of regularization in the counterfactual search. For example, we could use dropout to randomly switch features on and off in each iteration. Without dwelling further on the merit of this idea, we will now briefly show how this can be implemented.

¹¹Mutability constraints are currently not yet implemented for Latent Space generators.

¹²Flux models are now natively supported by our package and can be instantiated by calling `FluxModel()`.

¹³Note that in line 4 we also provide a field determining the likelihood. This is optional and only used internally to determine which loss function to use in the counterfactual search. If this field is not provided to the model, the loss function needs to be explicitly supplied to the generator.

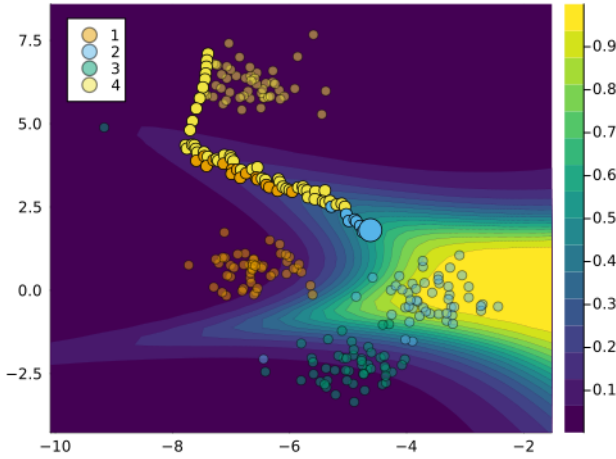


Fig. 6. Counterfactual path using generic counterfactual generator for multi-class classifier.

5.2.1 A Generator with Dropout. Listing 5.2.1 below implements two important steps: 1) create an abstract subtype of the `AbstractGradientBasedGenerator` and 2) create a constructor similar to the `GenericConstructor`, but with one additional field for the probability of dropout.

```
1 # Abstract supertype:
2 abstract type AbstractDropoutGenerator <:
  AbstractGradientBasedGenerator end
3 # Constructor:
4 struct DropoutGenerator <:
  AbstractDropoutGenerator
5     loss::Symbol # loss function
6     complexity::Function # complexity function
7     λ::AbstractFloat # strength of penalty
8     decision_threshold::Union{Nothing, AbstractFloat}
9     opt::Any # optimizer
10    τ::AbstractFloat # tolerance for convergence
11    p_dropout::AbstractFloat # dropout rate
12 end
```

Next, in listing 5.2.1 we define how feature perturbations are generated for our custom dropout generator: in particular, we extend the relevant function through a method that implements the dropout logic.

```
1 using CounterfactualExplanations.Generators
2 function Generators.generate_perturbations(
3     generator::AbstractDropoutGenerator,
4     ce::CounterfactualExplanation
5 )
6     s' = deepcopy(ce.s')
7     new_s' = propose_state(
8         generator, ce
9     )
10    Δs' = new_s' - s' # gradient step
11    # Dropout:
12    set_to_zero = sample(
13        1:length(Δs'),
14        Int(round(generator.p_dropout*length(Δs'))),
15        replace=false
16    )
17    Δs'[set_to_zero] .= 0
18    return Δs'
19 end
```

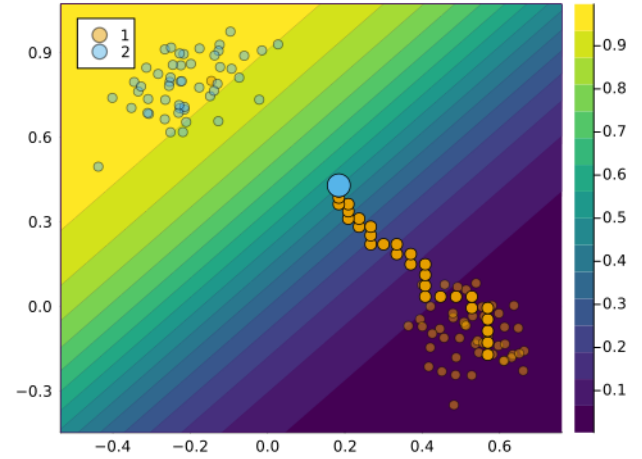


Fig. 7. Counterfactual path for a generator with dropout.

Finally, we proceed to generate counterfactuals in the same way we always do. The resulting counterfactual path is shown in Figure 7.

6. A Real-World Examples

Now that we have explained the basic functionality of `CounterfactualExplanations.jl` through some synthetic examples, it is time to work through examples involving real data.

6.1 Give Me Some Credit

The **Give Me Some Credit** dataset is one of the tabular real-world datasets that ship with the package [11]. It can be used to train a binary classifier to predict whether a borrower is likely to experience financial difficulties in the next two years. In particular, we have an output variable $y \in \{0 = \text{no stress}, 1 = \text{stress}\}$ and a feature matrix X that includes socio-demographic variables like age and income X . A retail bank might use such a classifier to determine if potential borrowers should receive credit or not.

For the classification task we use a Multi-Layer Perceptron with dropout regularization.

Using the Gravitational generator [1] we will generate counterfactuals for ten randomly chosen individuals that would be denied credit based on our pre-trained model. Concerning the mutability of features, we only impose that the age cannot be decreased.

Figure 8 shows the resulting counterfactuals proposed by Wachter in the two-dimensional feature space spanned by the age and income variables. An increase in income and age is recommended for the majority of individuals, which seems plausible: both age and income are typically positively related to creditworthiness.

6.2 MNIST

For our second example, we will look at image data. The MNIST dataset contains 60,000 training samples of handwritten digits in the form of 28x28 pixel grey-scale images [17]. Each image is associated with a label indicating the digit (0-9) that the image represents. The data makes for an interesting case study of Counterfactual Explanations because humans have a good idea of what realistic counterfactuals of digits look like. For example, if you were asked to pick up an eraser and turn the digit in the left panel of Figure 9 into a four (4) you would know exactly what to do: just erase the top part. Schut et al. (2021) [26] leverage this idea to illustrate

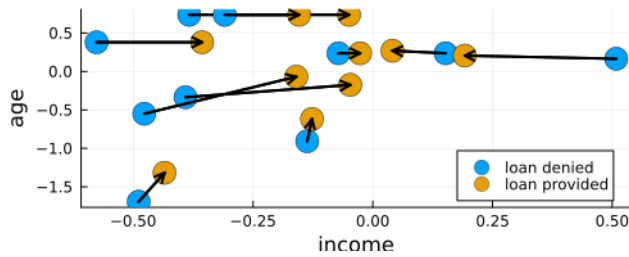


Fig. 8. Give Me Some Credit: counterfactuals for would-be borrowers proposed by the Gravitational Generator.

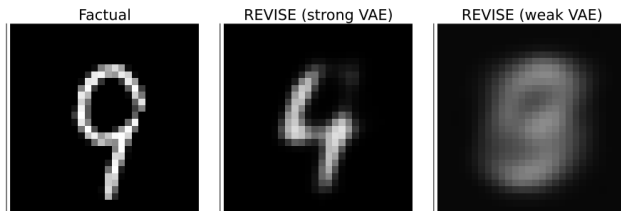


Fig. 9. Counterfactual explanations for MNIST using a Latent Space generator: turning a nine (9) into a four (4).

to the reader that their methodology produces realistic counterfactuals. In what follows we replicate some of their findings. You as the reader are therefore the perfect judge to evaluate the quality of the Counterfactual Explanations presented below.

On the model side, we will use a simple multi-layer perceptron (MLP). Listing 6.2 loads the data and the pre-trained MLP. It also loads two pre-trained Variational Auto-Encoders, which will be used by our counterfactual generator of choice for this task: REVISE.

```
1 counterfactual_data = load_mnist()
2 X, y = unpack_data(counterfactual_data)
3 input_dim, n_obs = size(counterfactual_data.X)
4 M = load_mnist_mlp()
5 vae = load_mnist_vae()
6 vae_weak = load_mnist_vae(;strong=false)
```

The proposed counterfactuals are shown in Figure 9. In the case in which REVISE has access to an expressive VAE (centre), the result looks convincing: the perturbed image does look like it represents a four (4). In terms of explainability, we may conclude that removing the top part of the handwritten nine (9) leads the black-box model to predict that the perturbed image represents a four (4). We should note, however, that the quality of counterfactuals produced by REVISE hinges on the performance of the underlying generative model, as demonstrated by the result on the right. In this case, REVISE uses a weak VAE and the resulting counterfactual is invalid. In light of this, we recommend using Latent Space search with care.

7. Discussion and Outlook

We believe that this package in its current form offers a valuable contribution to ongoing efforts towards XAI in Julia. That being said, there is significant scope for future developments, which we briefly outline in this final section.

7.1 Candidate models and generators

At the time of writing the package supports a handful of default models and generators either natively or through minimal augmentation. In future work, we would like to prioritize the addition of further predictive models and generators. Concerning the former, it would be useful to add native support for any supervised models built in MLJ.jl, an extensive Machine Learning framework for Julia [4]. This may also involve adding support for regression models as well as non-differentiable models. In terms of counterfactual generators, there is a list of recent methodologies that we would like to implement including MINT [13], ROAR [30] and PROBE [?]. Through its composable nature, our package allows for combining different approaches.

7.2 Additional datasets

For benchmarking and testing purposes it will be crucial to add more datasets to our library. We would like to prioritize datasets that have typically been used in the literature on counterfactual explanations including Adult, COMPAS and German Credit [12]. That being said, there is also scope for adding data sources that have so far not been explored much in this context including image, audio, natural language and time-series data.

7.3 Adding Foreign Language Support

The Julia language offers unique support for programming language interoperability. For example, calling R or Python is made remarkably easy through auxiliary packages like RCall.jl, PythonCall.jl and PyCall.jl, respectively. Early experimentation has shown that this functionality can be leveraged to make CounterfactualExplanations.jl compatible with models that were developed in foreign programming languages. At the time of writing this feature is still not mature enough, but we hope to add native support for explaining torch models trained in R or Python in the near future.¹⁴

8. Concluding remarks

The goal of this paper was to illustrate the need for explainability in machine learning and the promise of Counterfactual Explanations in this context. To this end, we introduced CounterfactualExplanation.jl: a package for generating Counterfactual Explanations and Algorithmic Recourse in Julia. Through various synthetic and real-world examples, we have demonstrated the basic usage of the package and shown how it can be easily customized and extended. We envision this package to one day constitute the go-to place for explaining arbitrary predictive models through a diverse suite of counterfactual generators. As a major next step, we would therefore like to interface our library with the popular MLJ.jl package for machine learning in Julia. The package can also serve as a testing ground for new and existing methodological approaches to Counterfactual Explanations and Algorithmic Recourse. We invite the Julia community to contribute to these goals through usage, open challenge and active development.

¹⁴Early experiments with this feature can be found here: <https://www.paltmeyer.com/CounterfactualExplanations.jl/dev/tutorials/interop/>

9. Acknowledgements

Patrick is grateful to his PhD supervisors and co-authors, Cynthia C. S. Liem and Arie van Deursen, for being supportive of his work on open-source developments. The authors are grateful to the Julia community for being welcoming and open and for supporting research contributions like this one.

10. References

- [1] Patrick Altmeyer, Giovan Angela, Aleksander Buszydlak, Karol Dobiczek, Arie van Deursen, and Cynthia Liem. Endogenous Macrodynamics in Algorithmic Recourse. In *First IEEE Conference on Secure and Trustworthy Machine Learning*.
- [2] Javier Antorán, Umang Bhatt, Tameem Adel, Adrian Weller, and José Miguel Hernández-Lobato. Getting a clue: A method for explaining uncertainty estimates. arxiv:2006.06848.
- [3] Alejandro Barredo Arrieta, Natalia Diaz-Rodriguez, Javier Del Ser, Adrien Bernetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, et al. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. 58:82–115.
- [4] Anthony D. Blaom, Franz Kiraly, Thibaut Lienart, Yiannis Simillides, Diego Arenas, and Sebastian J. Vollmer. MLJ: A Julia package for composable machine learning. 5(55):2704. doi:10.21105/joss.02704.
- [5] Christian Borch. Machine learning, knowledge risk, and principal-agent problems in automated trading. page 101852.
- [6] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, pages 77–91. PMLR.
- [7] Fenglei Fan, Jinjun Xiong, and Ge Wang. On interpretability of artificial neural networks. arxiv:2001.02522.
- [8] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. arxiv:1412.6572.
- [9] Mike Innes. Flux: Elegant machine learning with Julia. 3(25):602.
- [10] Shalmali Joshi, Oluwasanmi Koyejo, Warut Vijitbenjaronk, Been Kim, and Joydeep Ghosh. Towards realistic individual recourse and actionable explanations in black-box decision making systems. arxiv:1907.09615.
- [11] Kaggle. Give me some credit, Improve on the state of the art in credit scoring by predicting the probability that somebody will experience financial distress in the next two years.
- [12] Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A survey of algorithmic recourse: Definitions, formulations, solutions, and prospects. arxiv:2010.04050.
- [13] Amir-Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse: From counterfactual explanations to interventions. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 353–362.
- [14] Amir-Hossein Karimi, Julius Von Kügelgen, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse under imperfect causal knowledge: A probabilistic approach. arxiv:2006.06831.
- [15] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. Interpreting interpretability: Understanding data scientists’ use of interpretability tools for machine learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–14.
- [16] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. arxiv:1612.01474.
- [17] Yann LeCun. The MNIST database of handwritten digits.
- [18] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. 267:1–38.
- [19] Christoph Molnar. *Interpretable Machine Learning*. Lulu.com.
- [20] Ramaravind K Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 607–617.
- [21] Cathy O’Neil. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown.
- [22] R Kelley Pace and Ronald Barry. Sparse spatial autoregressions. 33(3):291–297.
- [23] Martin Pawelczyk, Sascha Bielawski, Johannes van den Heuvel, Tobias Richter, and Gjergji Kasneci. Carla: A python library to benchmark algorithmic recourse and counterfactual explanation algorithms. arxiv:2108.00783.
- [24] Rafael Poyiadzi, Kacper Sokol, Raul Santos-Rodriguez, Tijl De Bie, and Peter Flach. FACE: Feasible and actionable counterfactual explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 344–350.
- [25] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. 1(5):206–215.
- [26] Lisa Schut, Oscar Key, Rory Mc Grath, Luca Costabello, Bogdan Sacaleanu, Yarin Gal, et al. Generating Interpretable Counterfactual Explanations By Implicit Minimisation of Epistemic and Aleatoric Uncertainties. In *International Conference on Artificial Intelligence and Statistics*, pages 1756–1764. PMLR.
- [27] Dylan Slack, Anna Hilgard, Himabindu Lakkaraju, and Sameer Singh. Counterfactual explanations can be manipulated. 34.
- [28] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186.
- [29] Bob L Sturm. A simple method to determine if a music information retrieval system is a “horse”. 16(6):1636–1644.
- [30] Sohini Upadhyay, Shalmali Joshi, and Himabindu Lakkaraju. Towards Robust and Reliable Algorithmic Recourse. arxiv:2102.13620.
- [31] Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 10–19.
- [32] Kush R. Varshney. *Trustworthy Machine Learning*. Independently Published.

- [33] Sahil Verma, John Dickerson, and Keegan Hines. Counterfactual explanations for machine learning: A review. arxiv:2010.10596.
- [34] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. 31:841.
- [35] I-Cheng Yeh and Che-hui Lien. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. 36(2):2473–2480.