

```
using CounterfactualExplanations, Plots, PlotThemes, GraphRecipes
theme(:wong)
default(size=(500, 375))
include("paper/utils.jl")
```

## Introduction

In section Section ??

## Methodological background

Counterfactual search happens in the feature space: 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 ([?]). Typically the underlying methodology is presented in the context of binary classification:  $M : \mathcal{X} \mapsto y$  where  $y \in \{0, 1\}$ . Let  $t = 1$  be the target class and let  $\bar{x}$  denote the factual feature vector of some individual outside of the target class, so  $\bar{y} = M(\bar{x}) = 0$ . We follow this convention here, though it should be noted that the ideas presented here also carry over to multi-class problems and regression ([?]).

### Generic framework

Then the counterfactual search objective originally proposed by [?] is as follows

$$\min_{\underline{x} \in \mathcal{X}} h(\underline{x}) \quad \text{s. t.} \quad M(\underline{x}) = t \quad (1)$$

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

$$\underline{x} = \arg \min_{\underline{x}} \ell(M(\underline{x}), t) + \lambda h(\underline{x}) \quad (2)$$

Here  $\ell$  denotes some loss function targeting the deviation between the target label and the predicted label and  $\lambda$  governs the strength of the complexity penalty. Provided we have gradient access for the black-box model  $M$  the solution to this problem (Equation ??) 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 - **GenericGenerator** - in our package. The hyperparameter  $\lambda$  is typically tuned through grid

search. Conventional choices for  $\ell$  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.<sup>1</sup>

Numerous - and in some cases competing - extensions to this simple approach have been developed since counterfactual explanations were first proposed in 2017 (see [?] and [?] for surveys). The various approaches largely differ in how they define the complexity penalty. In [?], 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 ([?]); **actionability** - the proposed feature perturbation should actually be actionable ([?], [?]); **plausibility** - the counterfactual explanation should be plausible to a human ([?]); **unambiguity** - a human should have no trouble assigning a label to the counterfactual ([?]); **sparsity** - the counterfactual explanation should involve as few individual feature changes as possible ([?]); **robustness** - the counterfactual explanation should be robust to domain and model shifts ([?]); **diversity** - ideally multiple diverse counterfactual explanations should be provided ([?]); and **causality** - counterfactual explanations reflect the structural causal model underlying the data generating process ([?],[?]).

## Counterfactuals for Bayesian models

For what follows it is worth elaborating on the approach proposed in [?]. The authors demonstrate that many of the abovementioned desiderata can be addressed very easily, if the classifier  $M$  is Bayesian. In particular, they show that close, realistic, sparse and unambiguous counterfactuals can be generated by implicitly minimizing the classifier’s predictive uncertainty through a greedy counterfactual search. Formally, they define  $h(\cdot)$  as the predictive entropy of the classifier, which captures both **epistemic** and **aleatoric** uncertainty: the former is high on points far away from the training data while the latter is high in regions of the input space that are inherently noisy. Both are regions we want to steer clear off in our counterfactual search and hence predictive entropy is an intuitive choice for a complexity penalty. The authors further point out that any solution that minimizes cross-entropy loss (Equation ??) also minimizes predictive entropy:  $\arg \min_{\underline{x}} \ell(M(\underline{x}), t) \in \arg \min_{\underline{x}} h(\underline{x})$ . Let  $\widetilde{\mathcal{M}}$  denote the class of binary classifiers that incorporate predictive uncertainty, then the previous observation implies that the optimal solution to counterfactual search (Equation ??) can be restated as follows:

$$\underline{x} = \arg \min_{\underline{x}} \ell(M(\underline{x}), t) \quad , \quad \forall M \in \widetilde{\mathcal{M}} \quad (3)$$

---

<sup>1</sup>While the rationale for this convention is not entirely obvious, implementations of loss functions with respect to logits are often numerically more stable. For example, the `logitbinarycrossentropy`( $\hat{y}$ ,  $y$ ) implementation in `Flux.Losses` (used here) is more stable than the mathematically equivalent `binarycrossentropy`( $\hat{y}$ ,  $y$ ).

We can drop the complexity penalty altogether and still generate effective counterfactual explanations. As we will see below, even a fast and greedy counterfactual search proposed in [?] yields good results in this setting. The approach has been implemented as **GreedyGenerator** in our package and should only be used with classifiers of type  $\widetilde{\mathcal{M}}$ . It is worth noting that the findings in [?] are not mutually exclusive of many of the other methodologies that have been put forward. On the contrary, we believe that they are complementary: the generic counterfactual search proposed in [?], for example, can be shown to produce more plausible counterfactuals in the Bayesian setting. Similarly, there is no obvious reason why recent work on diversity ([?]), robustness ([?]) and causality ([?],[?]) could not be complemented by the findings in [?]. For this reason we are highlighting [?] here and have prioritized it in the development of **CounterfactualExplanations**. While there is no free lunch and  $M \in \widetilde{\mathcal{M}}$  may seem like a hard constraint, recent advances in probabilistic machine learning have shown that the computational cost involved in Bayesian model averaging is lower than we may have thought ([?], [?], [?], [?]).

## Using CounterfactualExplanations

The package is built around two modules that are designed to be as scalable as possible through multiple dispatch: 1) **Models** is concerned with making any arbitrary model compatible with the package; 2) **Generators** is used to implement arbitrary counterfactual search algorithms.<sup>2</sup> The core function of the package `generate_counterfactual` uses an instance of type `T <: FittedModel` produced by the **Models** module (Figure ??) and an instance of type `T <: Generator` produced by the **Generators** module (Figure ??). Relating this back the methodology outlined in Section ??, the former instance corresponds to the model  $M$  while the latter defines the rules for the counterfactual search (Equation ?? and Equation ??). In the following we will demonstrate how to use and extend the package architecture through a few examples.

```
plot(code, method=:tree, fontsize=8, nodeshape=:rect, axis_buffer=0.3)
```

```
p = plot(CounterfactualExplanations.Models.FittedModel, method=:tree, fontsize=8, nodeshape=:rect, axis_buffer=0.3)
savefig(p, "paper/www/models.png")
```

```
p = plot(CounterfactualExplanations.Generator, method=:tree, fontsize=8, nodeshape=:rect, axis_buffer=0.3)
savefig(p, "paper/www/generators.png")
```

---

<sup>2</sup>We have made an effort to keep the code base as flexible and scalable as possible, but cannot guarantee at this point that really any counterfactual generator can be implemented without further adaptation.

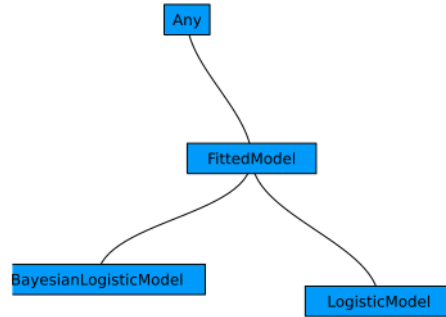


Figure 1: Schematic overview of the `FittedModel` base type and its descendants.

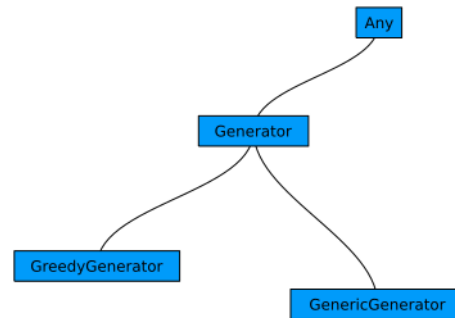


Figure 2: Schematic overview of the `Generator` base type and its descendants.

## Getting started

The code below provides a complete example demonstrating how the framework presented in Section ?? can be implemented in Julia using the `CounterfactualExplanations` package: using a synthetic data set with linearly separable samples we firstly define our model and then generate a counterfactual for a randomly selected sample. Figure ?? 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) <: FittedModel`. 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 multiple dispatch to extend these functions for use with arbitrary models. For now it is enough to note that those methods define how the 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) choose a random sample `x_factual`; 2) compute its factual label `y_factual` as predicted by the model ( $M(\bar{x}) = 0$ ); and 3) specify the other class as our `target` label ( $t = 1$ ) along with a desired level of `confidence` in the final prediction  $M(x) = t$ .

The last two lines of the code below define the counterfactual generator and finally run the counterfactual search. The first three fields of the `GenericGenerator` are reserved for hyperparameters governing the strength of the complexity penalty, the step size for gradient descent and the tolerance for convergence. The fourth field accepts a `Symbol` defining the type of loss function  $\ell$  to be used. Since we are dealing with a binary classification problem logistic binary cross-entropy is an appropriate choice.<sup>3</sup> The fifth and last field can be used to define mutability constraints for the features.

```
# Data:
using CounterfactualExplanations, Random
Random.seed!(1234);
N = 100 # number of data points
using CounterfactualExplanations.Data
x, y = toy_data_linear(N)

# Model:
```

---

<sup>3</sup>As mentioned earlier, the loss function is computed with respect to logits and hence it is important to use logistic binary cross-entropy loss as opposed to just binary cross-entropy.