

LANGUAGE MODELS CAN HELP TO LEARN HIGH-PERFORMING COST FUNCTIONS FOR RECOURSE

Anonymous authors

Paper under double-blind review

ABSTRACT

Algorithmic recourse is a specialised variant of counterfactual explanation, concerned with offering actionable recommendations to individuals who have received adverse outcomes from automated systems. Most recourse algorithms assume access to a cost function, which quantifies the effort involved in following recommendations. Such functions are useful for filtering down recourse options to those which are most actionable. In this study, we explore the use of large language models (LLMs) to help label data for training recourse cost functions, while preserving important factors such as transparency, fairness, and performance. We find that LLMs do generally align with human judgements of cost and can label data for the training of effective cost functions, moreover they can be fine-tuned with simple prompt engineering to maximise performance and improve current recourse algorithms in practice. Previously, recourse cost definitions have mainly relied on heuristics and missed the complexities of feature dependencies and fairness attributes, which has drastically limited their usefulness. Our results show that it is possible to train a high-performing, interpretable cost function by consulting an LLM via careful prompt engineering. Furthermore, these cost functions can be customised to add or remove biases as befitting the domain and problem. Overall, this study suggests a simple, accessible method for accurately quantifying notions of cost, effort, or distance between data points that correlate with human intuition, with possible applications throughout the explainable AI field.

1 INTRODUCTION

Algorithmic recourse has emerged as one of the most impactful areas of explainable AI (Karimi et al., 2022). The field focuses on generating actionable counterfactual recommendations to users who were treated unfavorably by automated systems, with the canonical example being a rejected bank loan application, and what actions could be taken by a user to have it accepted in future (Ustun et al., 2019). In such a scenario, a cost function is needed to quantify how much effort a recourse recommendation would take, so that algorithms can consider this during optimisation. Separately, it is worth noting that the field has branched out to consider positive outcomes with gain functions and semifactual recourse (Kenny & Huang, 2024). In either case, these functions must align with human domain knowledge and intuition, so they can inform appropriate recourse selection. In this paper, we focus on cost functions, and show how large language models (LLMs) can be used to largely automate their design while maintaining desirable aspects such as transparency and fairness.

Typically in recourse, a cost function is assumed a priori, often as some variant of an L_p norm on the feature space (Keane et al., 2021). For example, an L_0 norm assigns higher cost to recourse recommendations that change more features, although this ignores other factors such as how much they are changed. A (weighted) L_1 or L_2 norm can incorporate magnitude information, but not pairwise or higher-order interactions between features. These can be added in an ad hoc manner, but are challenging to formalise and combine. As an alternative, we examine if the issue of recourse cost can be addressed in a flexible and scalable way by tapping into the tacit domain knowledge of LLMs. We show how, with the right prompting, LLMs can be consulted to compare the costs of pairs of recourses, creating a labelled dataset for training either neural network cost functions or transparent tree-based ones (Kanamori et al., 2022; Bewley & Lecue, 2022). Our results suggest that future research into cost functions may benefit from the use of LLMs.

2 COST FUNCTION DESIDERATA

We begin by considering what constitutes a high-performing cost function for recourse applications. Ideally, a cost function should satisfy many intricate criteria which basic L_p norms cannot, such as variable feature weighting and dependencies. Here, we outline our desiderata grounded in prior literature, which will form the basis for subsequent evaluation.

1. *Feature Cost.* A cost function should have different weighting considerations for each feature in the data. For example, adding an additional credit card is generally easier than increasing your down payment (Rawal & Lakkaraju, 2020).
2. *Relative Cost.* A cost function should weigh the cost of a given change differently at different points in the distribution, if appropriate. For example, going from the 55-60th percentile in an exam score may be easier than going from the 90-95th (Ustun et al., 2019).
3. *Dependent Cost.* It must be possible to represent relevant dependencies between two or more features. For example, applying for college funding is usually easier if you are native to a country rather than an immigrant (Karimi et al., 2022).
4. *Fair Cost.* Cost functions should take into account any fairness properties relevant to a given domain and application (Von Kügelgen et al., 2022). In this paper, we define fairness as the cost function not varying its output if demographic information is mutated.

These desiderata have been extensively discussed in the literature cited above. We do not claim this to be an exhaustive list, but a reasonable starting point.¹

3 METHOD

This section outlines our four-step framework for learning cost functions. First, synthetic recourse examples are generated by randomly perturbing a set of data points subject to actionability constraints. Second, pairs of recourse examples are selected at random for comparison. Third, an LLM is queried to provide ratings (i.e. labels) for these comparisons. Finally, the resultant dataset is used to train a cost function. In this process, we assume access to a capable chatbot LLM which may be queried at liberty, and that the data domain is tabular in nature.

3.1 GENERATING SYNTHETIC RECOURSES

Let $\mathcal{D} = \{x_i\}_{i=1}^N \subset \mathbb{R}^d$ denote a given dataset, where each x_i represents a d -dimensional feature vector. For our purposes, we benefit from \mathcal{D} being as diverse as possible. We define a stochastic perturbation function $\phi : \mathbb{R}^d \times \mathcal{A} \rightarrow \Delta(\mathbb{R}^d)$, where \mathcal{A} denotes a set of actionability constraints (see Appendix A for details). The number of features to be perturbed is problem-specific and will determine the cost function’s capabilities in deployment. Here, we randomly select this number from a truncated geometric distribution, which favors perturbations of one feature to focus on sparsity, which is desired in recourse (Keane et al., 2021; Karimi et al., 2022). See Appendix B for details.

For each data point x_i and perturbed feature $f \in \{1, \dots, d\}$, we apply the following perturbation:

$$x'_i[f] = \begin{cases} \sim \text{Uniform}(\text{categories}_f) & \text{if } f \text{ is categorical} \\ x_i[f] + \epsilon : \epsilon \sim \mathcal{E}_f & \text{if } f \text{ is continuous,} \end{cases} \quad (1)$$

where $\text{Uniform}(\text{categories}_f)$ is a uniform distribution over categories and \mathcal{E}_f is a finite set of perturbations for a continuous feature (positive/negative multiples of the standard deviation across \mathcal{D}).

This process generates a set of recourse examples $\mathcal{R} = \{(x_i, x'_i)\}_{i=1}^N$, where each x_i represents an original instance and x'_i is the corresponding synthetic perturbation. We use a finite set of perturbation magnitudes for numerical features because it allows a direct comparison between exactly the same change at different parts of a feature distribution. This helps to learn relative differences in

¹Another possible desideratum is *Individual Cost*, whereby even if two individuals have identical feature values, they could still have different ideal recourse recommendations based on their preferences (Nauta et al., 2023; Rawal & Lakkaraju, 2020). However, this is largely a separate human computer interaction (HCI) question, and we are instead focused on the training of the cost function itself.

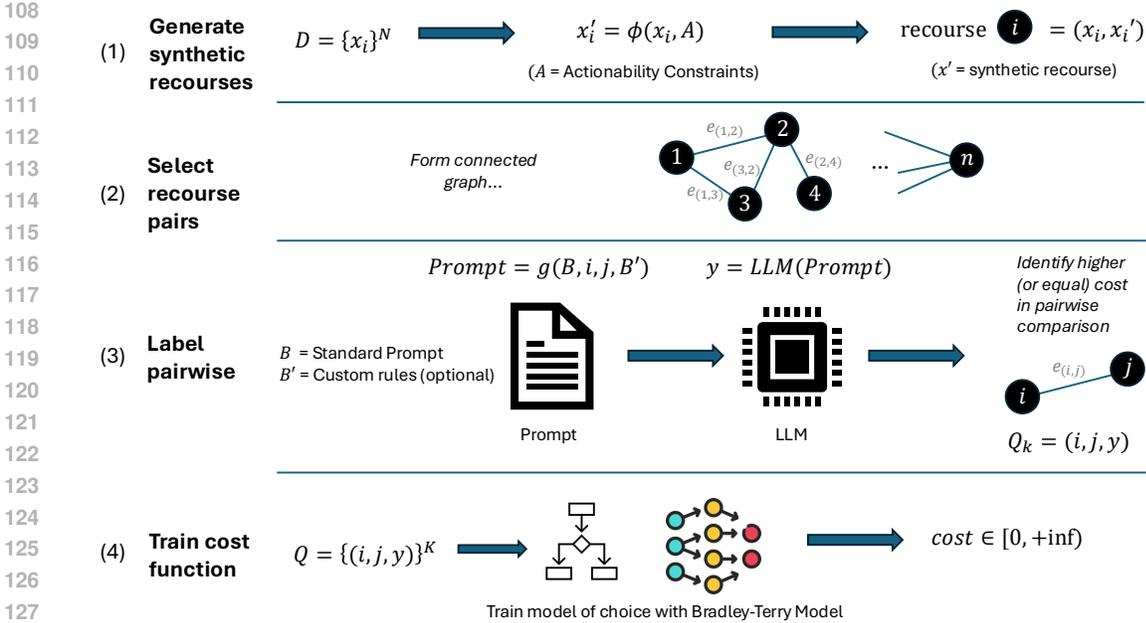


Figure 1: Method Schematic: (1) Each instance in a given dataset is perturbed within actionability constraints to simulate a recourse situation. (2) Pairs of these recourses are selected for comparison in such a way as to form a connected graph (where a path exists between all pairs of recourses). (3) Each edge in the graph is then labeled with an LLM which judges which of the two corresponding recourses takes a higher cost to achieve (or optionally an additional “equal cost” option if specified). (4) The dataset of comparisons is used to train the cost function, in our case either a transparent tree model or an MLP.

cost for the same change, thereby addressing the relative cost criterion (i.e., Desideratum 2). The parameters can be tuned to suit the specific requirements of the problem domain.

3.2 SELECTING RECOURSE PAIRS

Next, we select a set of $K \leq N^2$ pairs of recourse examples from \mathcal{R} which will be presented to an LLM for cost comparison. This process can be understood as connecting the recourses into an undirected graph structure. In forming this graph, we enforce that each recourse must have a minimum of K_{\min} edges, and that the graph as a whole forms a single connected component (where a path exists between all pairs of recourses). We find that this improves the performance of the final cost function, as it allows the costs for all recourses to be estimated on the same scale. To enforce the relative cost criterion, we prioritise edges between recourses which perturb the same continuous feature at two different parts of the distribution by exactly the same amount. This has the effect of forcing the LLM to reason about the difference in cost between e.g. increasing salary from 30-35k versus 50-55k. We also add edges to enforce comparisons of the same feature changes for different feature dependencies, e.g. two recourses which have the same increase in loan amount, but different credit ratings, which can be used to enforce the relative cost criterion of Desideratum 3 (see Section 4 later). The total additional edges from this enforcement is set to 10% of the total data for both, adding 20% extra data on average. Aside from these considerations, we find that the algorithm used to construct the graph of recourse pairs is relatively unimportant. In practice, any algorithm forming a connected graph subject to the K_{\min} constraint seems to work well. We used a random spanning tree algorithm in all experiments.

3.3 PAIRWISE LLM LABELLING

For the standard prompt structure \mathcal{B} (see Appendix D), we begin by instructing the LLM that it is a helpful assistant to a data scientist which labels data. It is then told the task of comparing two

162 individuals and their respective feature changes. We then enumerate the features, as well as their
 163 descriptions. The LLM is then asked to reason about which of the two given recourses requires
 164 more effort for the individual to achieve (i.e. cost), and finally to respond with a label of 1 (first
 165 requires more effort), or 0 (second requires more effort). Optionally, we also permit a third category
 166 of 0.5, indicating a judgement that equal effort is required, which is a useful de-biasing signal in con-
 167 texts where features represent sensitive demographic attributes. **The prompt then gives a high-level**
 168 **overview of the desiderata in Section 2.** In addition, the LLM is instructed to use chain-of-thought
 169 to increase performance and reduce social biases (Kamruzzaman & Kim, 2024). In other experi-
 170 ments, we also fine-tune the prompt more with a set of desired cost function parameters, denoted
 171 by \mathcal{B}' . For the full prompts, see Appendix D. The output of this stage is a set of K comparisons
 172 $\mathcal{Q} = \{(i, j, y)\}_{k=1}^K$, where i and $j \neq i$ are indices of a pair of recourse examples from \mathcal{R} and
 173 $y \in \{0, 0.5, 1\}$ denotes the LLM’s effort/cost judgement.

174 3.4 TRAINING THE COST FUNCTION

176 Finally, we use the dataset of LLM comparisons \mathcal{Q} to train a cost function $C : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$.
 177 Inspired by Rawal & Lakkaraju (2020), as well as the dominant approach to learning reward models
 178 from pairwise comparisons (Kwon et al., 2023), we train cost functions using the Bradley-Terry
 179 model. That is, given a cost function C and a pair of recourses (x_i, x'_i) and (x_j, x'_j) , we define the
 180 predicted probability that recourse i has higher cost than recourse j as

$$181 \hat{y}_C(i, j) = \frac{1}{1 + \exp(C(x_j, x'_j) - C(x_i, x'_i))}. \quad (2)$$

184 Our cost function training objective is to minimise the binary cross-entropy between these predicted
 185 comparison probabilities and the labels provided by the LLM across all training examples:

$$186 \arg \min_{C \in \mathcal{M}} \left[- \sum_{(i, j, y) \in \mathcal{Q}} y \log(\hat{y}_C(i, j)) + (1 - y) \log(1 - \hat{y}_C(i, j)) \right], \quad (3)$$

189 where \mathcal{M} is a chosen model class. Since this loss is differentiable, we can define \mathcal{M} as the class of
 190 MLP neural networks and train by stochastic gradient descent. As an alternative, we also consider
 191 the class of axis-aligned decision trees up to a maximum leaf count L_{\max} , which offers greater trans-
 192 parency. To train a non-differentiable tree with the pairwise Bradley-Terry loss, we use a bespoke
 193 algorithm developed by (Bewley & Lecue, 2022) (and refined in (Bewley et al., 2022)).

194 We one-hot encode categorical features (or binary encode ones with only 2 categories) and con-
 195 catenate the original data point x , the perturbed recourse point x' and the feature-wise difference
 196 $x' - x$ into a single vector $[x, x', x' - x] \in \mathbb{R}^{3d}$. In practice, we found that this simple feature aug-
 197 mentation step significantly improved the models’ ability to learn costs. As a final post-processing
 198 step, we shift the outputs of trained models to ≥ 0 on all training data. This has no impact on the
 199 Bradley-Terry loss, but produces the expected behaviour for a non-negative cost function to only
 200 output non-negative values.

202 4 EVALUATION

204 In our evaluation, we seek to understand how to train effective cost functions utilising LLMs.
 205 Throughout, we focus on three datasets, the Home Equity Line of Credit (HELOC) dataset (Mstz,
 206 2024) for predicting whether someone will repay their account, the Adult Census dataset (Becker &
 207 Kohavi, 1996), for predicting if an individual earns higher than 50k per year, and the German Credit
 208 dataset (Hofmann, 1994), for classifying a client’s credit risk. All are binary classification tasks,
 209 and we considered the first 800 instances from each dataset for training/testing of the cost func-
 210 tion. All categorical features were modeled as binary 0/1 options, except German Credit which has
 211 multi-categorical features one-hot encoded. After creating the dataset of pairwise comparisons, and
 212 adding the additional links described in Section 3, we had 22,000 pairwise training examples on av-
 213 erage, which was divided into 80/20% training/testing, respectively, for the cost functions. Currently
 214 GPT-4o represents state-of-the-art performance on many benchmarks (OpenAI, 2024), and indeed
 215 it is shown to be fairer than prior models Bowen III et al. (2024), so we used it in all our tests. As
 our data is mostly synthetic, we do not expect GPT-4o’s known memorization of the datasets to be
 an issue (Bordt et al., 2024).

4.1 COMPARING HUMAN AND LLM JUDGEMENT OF COST

A natural first question is whether or not LLMs can provide a judgement of cost that aligns with human intuition. Hence, our evaluation started with a study to compare pairwise choices of cost between GPT-4o and humans. Participants were shown two individuals, a proposed change (i.e. their recourse), and asked to select which of these would require more “effort” (i.e. cost). We limited the options to a forced choice between Recourse 1 and Recourse 2, with no equal effort option, which we primarily reserved for situations involving demographic fairness (which this study did not involve). With a distribution of responses from humans in hand, this was then compared to GPT-4o’s responses on the same questions using our prompt template in Appendix D. Note that because the LLM would arbitrarily choose Option 1 when its chain of thought communicated it was unsure, we allowed it a third option to identify this, and then replaced these data points with random answers. This allowed a more accurate comparison to humans, who tend to choose at random when unsure (Gigerenzer & Goldstein, 1996).

The materials covered all three datasets, with six questions for each. These six questions were split into three sets of two representing the first three parts of the desiderata, respectively.² Participants were also asked to choose how “close” they felt the two were, so we could compare their uncertainty with the LLM. See Figure 6 for an example and the supplementary material for the full survey.

We randomly recruited thirty industry data scientists for the purposes of the study. The participants were not compensated; all volunteered to participate. In total, 20 of the participants were male, 10 were female, all were aged 18+, and there was a mix of native/non-native English speakers.³ The study obtained IRB approval.

The metric of interest was how the distributions of responses from humans matches that of the LLM. The test used was the Chi-square test of independence. A second metric was whether or not the most common response from the LLM and humans was identical, represented as mode: Y (they were equal), or mode: N (they were not equal). Lastly, we asked humans to quantify how far apart they felt the two options were, so we could quantify their certainty compared to the LLM.

The LLM was unsure of the answer 13.8% of the time, and this was replaced with random responses to simulate human uncertainty. Correlating LLM uncertainty to humans, we observe a strong positive correlation (Person’s $r=0.5$; $p < 0.04$) in Figure 7, indicating that users and the LLM had approximately the same level of uncertainty across the same questions. Figure 2 shows more results. Overall, there is a tendency of the LLM to accurately align with the human labellers, with 15/18 of questions having statistically similar distributions (i.e. $p > 0.05$). When considering the most common responses (i.e. mode: Y), 15/18 different questions are also in agreement, two of

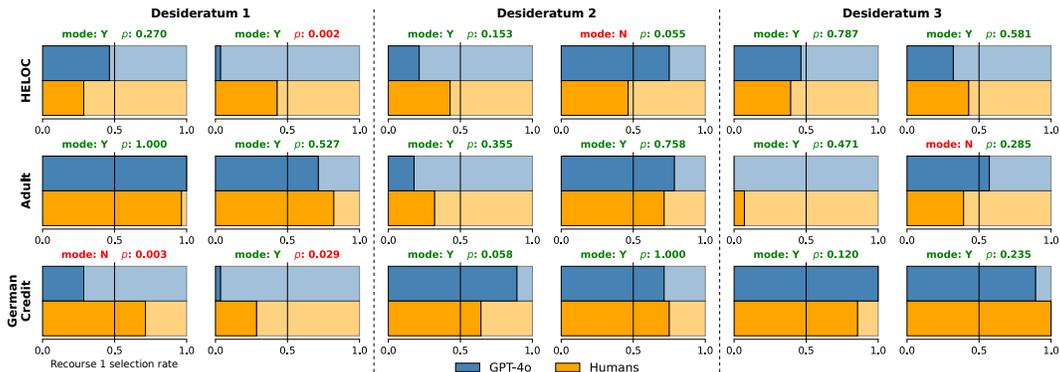


Figure 2: Human Study Results: 15/18 of the questions had the same modal response (i.e. the mode was Y), and 15/18 statistically similar distributions. Overall, 17/18 had one or the other, and were mostly aligned. Note we are trying to show the p -value is greater than 0.05, because we do not want to reject the null that humans and LLMs are aligned in cost judgement.

²Note that we did not evaluate Desideratum 4 with humans due to ethical concerns.

³Although the human sample is somewhat biased, results show they are aligned with the LLM, which increasingly simulate population user responses in surveys (De Bona et al., 2024).

270 these encompassing the questions without statistically similar distributions. Together, this can be in-
 271 terpreted to suggest that the LLM is in alignment for 17/18 of the questions. These results highlight
 272 that LLMs largely agree with human judgment of cost.
 273

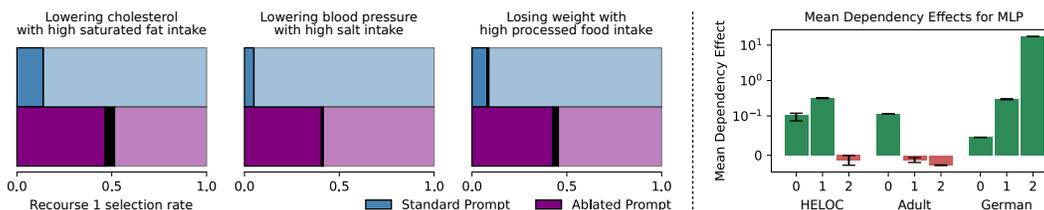
274 4.2 TRAINING THE COST FUNCTIONS

275 We used two types of prompts to label data for the cost functions, the *standard prompt*, and the
 276 *custom prompt*. The **standard prompt is identical to what was used in the user study and uses only**
 277 **a high-level description of the desiderata to instruct the LLM, whilst the custom prompt attempts**
 278 **to fine-tune the resultant cost function with a ground truth we defined in the prompt (i.e., B' in**
 279 **Figure 1).** The point of this custom prompt is to see if we can e.g. re-order feature importance,
 280 manipulate the spectrum of cost for numerical features, add dependencies, and fairness attributes,
 281 see Section H for details on the ground truth chosen. **This is important because (for example) the**
 282 **definition of fairness varies (Mehrabi et al., 2021), so we need to fine-tune different aspects of the**
 283 **cost function in practice.** The choice of ground truth is largely irrelevant, we are simply seeing
 284 if it can be worked into the final cost function via the prompt. We trained either an MLP model
 285 or a tree for 50,000 batches of size 32. So, in total, there are 2 models we are testing across 3
 286 datasets with 2 prompt types. We chose these models because trees help with transparency required
 287 in financial applications (Bewley et al.), and MLPs are differentiable, which is often required in
 288 recourse algorithms (Wachter et al., 2017).
 289

290 4.3 DEPENDENCY TEST

291 Perhaps the primary advantage of using LLMs to learn cost functions is that they have the potential
 292 to naturally model causal feature dependencies, which is the most intractable part of hand-designing
 293 a cost function. In this test, we examine the ability of LLMs to naturally label this with our standard
 294 prompt (i.e., no dependencies are mentioned in the prompt). We consider both synthetic and real
 295 data in this process. **Synthetic data is considered because there is a risk that the LLM can only reason**
 296 **about causal dependencies on well known recourse datasets used for counterfactual generation, as**
 297 **the generated counterfactuals may be in the LLM training data.**
 298
 299

300 **Synthetic Data** The generation of the synthetic data is detailed in Appendix K. In short, we crafted
 301 a novel dataset of known scientific dependencies in a medical domain, where data privacy laws
 302 should give additional reassurances that no such dataset was used to fine-tune the LLM, or sub-
 303 sequently used in recourse research papers. The dependencies were (1) that it is harder to lower
 304 cholesterol levels with a high saturated fat intake, (2) that it is harder to lower blood pressure
 305 with a high dietary salt intake, and (3) that it is harder to lose weight if consuming a large amount of
 306 heavily processed food. The ground truth was always Recourse 2, and we allowed the LLM to chose
 307 Recourse 1 or 2 as the one of higher cost, or 0 (i.e., uncertain). We compared our standard prompt
 308 to an ablated version without the desiderata, to understand more how this helps. The most important
 309 difference between these is that the ablated prompt has no explicit instruction to consider dependen-
 310



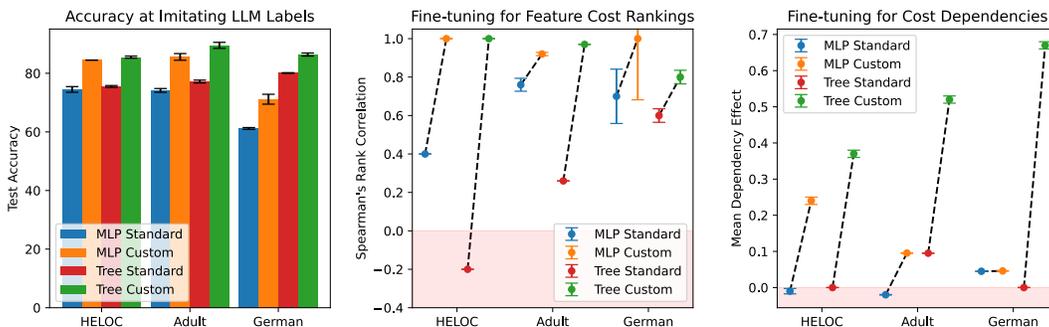
311
 312
 313
 314
 315
 316
 317
 318 **Figure 3: (left) Synthetic Data:** Comparing the standard prompt (with the desiderata included
 319 in the prompt) and the ablated version (without the desiderata), the LLM was 90% + accurate at
 320 labelling the three known scientific causal dependencies, but only with the desiderata inserted into
 321 the prompt. Note the black areas in the data indicate the probability of the LLM being uncertain.
 322 **(right) Real Data:** The trained MLP cost functions successfully learned 6/9 of the ground truth
 323 dependencies suggested by Claude Sonnet 3.5., showing a general trend that the LLM can generally
 identify suitable dependencies in its labeling which are subsequently learned by the cost functions.

324
325
326
327
328
329
330
cies when evaluating cost of recourses. The results are shown in Figure 3(left), where the standard prompt correctly identified all three dependencies with a mean accuracy of 91%, compared to the same prompt with the desiderata ablated which was not significantly better than random guessing. Overall, this shows how we can trust the LLM to label reasonable causal dependencies in novel domains, but only if we (1) use chain-of-thought prompting and (2) the desiderata⁴, which includes instructions to the LLM to explicitly look for dependencies.

331
332
333
334
335
336
337
338
339
340
341
342
343
Real Data. We prompted Claude Sonnet 3.5 to list the most important feature dependencies in each dataset (to help avoid leakage with GPT-4o), and repeated this 10 times to pick out three which were listed the most for our ground truths, see Appendix D and G. We iterated all the testing data with each cost function variation, and manually adjusted the data to subtract the cost of the less costly recourse option from the higher, hence, a positive score shows that the dependency is present in the cost function. In Figure 3 and Figure 4, we refer to this as the “*Mean Dependency Effect*”, where positive scores indicate the dependency has been learned to match the ground truth. The present results can be seen in Figure 3(right). Overall, 6/9 of all dependencies were modeled in accordance with Claude’s ground truth in the MLP cost function, showing a generally positive ability to learn appropriate dependencies. In contrast, the tree models only learned one of these with the other eight showing 0 cost. The reason for this is likely that the tree would require most splits to learn the necessary dependency, but the MLP forms a smoother interpretation of the labels and learned the dependencies more easily.

344 4.4 FINE-TUNING EXPERIMENTS

345
346
347
348
349
350
351
352
353
354
355
356
Going forward, we consider a suite of experiments which try to fine-tune the prompt to achieve different results in the cost function. This is because judgement of cost often needs to be tuned to certain context. For example, during an economic downturn, a bank might need to adjust its lending criteria and weight features differently. In addition, the definition of fairness varies substantially between contexts (Mehrabi et al., 2021), so this also needs to be finetuned occasionally. We design a fine-tuning experiment using our custom prompting scheme B' originally described in Figure 1. The prompting scheme adds a high-level description B' to the original prompt B indicating (1) how costly each feature should be to mutate in order, (2) how numerical features should change in cost at different parts of the distribution, (3) any dependencies we want to exaggerate or create, and (4) any fairness attributes desired. Not all these need to be specified during fine-tuning of cost functions, but we test all here for a complete experiment. When perturbing features in the subsequent



368
369
370
371
372
373
374
375
376
377
Figure 4: Accuracy and Desideratum 1 and 3 Fine-Tuning: (left) Accuracy of the cost functions at imitating the LLM’s pairwise labels on test data. (middle) Ability of the cost functions to be fine-tuned to correlate with ground truth feature cost rankings specified in the custom prompt. A score of 1 illustrates the rankings are perfectly learned by the cost function. (right) Ability of the cost functions to be fine-tuned to weight dependencies specified in the custom prompt. Notably, the Adult MLP flips from a negative to positive dependency effect, showing we can reverse the cost of certain dependencies if desired. Standard error is shown. Red background indicated negative correlation or dependency effect for middle and right plots, respectively.

⁴In additional unreported experiments, we found that the ability of the LLM to reason successfully about causal dependencies requires chain-of-thought prompting.

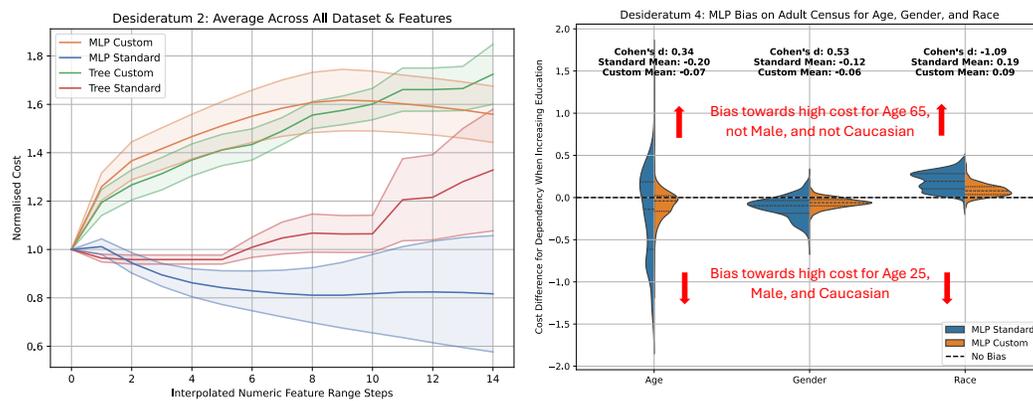


Figure 5: Desideratum 2/4 Fine-Tuning Results: (left) As the numeric features rise in value for the custom models, so does their respective cost relative to the standard ones. (right) Instructing the LLM to not unfairly discriminate between demographic information increased its fairness when suggesting recourse which increased education level. Overall, these results show how it is possible to fine-tune the cost function on Desideratum 2 and 4.

tests, categorical features were either flipped for binary or randomly changed for multi-categorical features, numerical features were perturbed upwards a standard deviation.

Desideratum 1 Here the ground truth specified in the custom prompt was a specific rank ordering of how costly each feature should be to mutate. Each testing datum had each feature perturbed to test its cost, the results for each feature were averaged and reported across four random seeds. Each feature was rank ordered in a list and compared against the ground truth defined in the custom prompt with Spearman’s rho ρ . Results can be seen in Figure 4(middle), where the custom prompt is compared to the standard one (which did not specify what the most costly features should be). Overall, the custom prompt-based cost functions successfully moved towards the new features rank orderings as instructed in B' , with Heloc learning them perfectly, illustrating that it is possible to realign the relative importance of features if desired.

Desideratum 2 Here the ground truth specified in the custom prompt was that each numerical feature should be harder to mutate the higher it gets in value. Each testing datum had each numerical feature perturbed upwards to test its cost at 16 evenly spaced intervals. Each feature across all datasets were averaged and again shown across four random seeds in Figure 5(left). The average Spearman’s ρ for the custom models across all features and datasets was 0.41, compared to 0.04 on the standard models, showing that the numerical features have a gradual trend of increasing their cost the higher the mutation starts, which aligns with the original ground truth schematic B' given to the LLM. This illustrates that it is possible to fine-tune the relative cost of numerical features across their spectrum if desired.

Desideratum 3 Here the ground truth specified in the custom prompt was to purposefully enforce the worst performing previously tested dependencies in each dataset in Section ??(right). The point is to see if we can correct them to be a positive mean cost dependency. As before, each testing datum had each dependency tested the same as Section 4.3, all were averaged and again shown across four random seeds in Figure 4. Notably, the negative cost associated with the dependency in Adult Census and Heloc flipped to be positive, showing it is possible to fine-tune this if desired. Moreover, the tree models all went from no/little cost associated with each dependency, to a positive one. Lastly, the strength of the positive cost in Heloc and German Credit for the MLP models increased, showing that by adding the dependency directly to the prompt, we can strengthen the dependency cost. This illustrates that it is possible to fine-tune dependencies if desired, simply by instructing the LLM to explicitly consider this dependency.

Desideratum 4 Here the ground truth specified in the custom prompt was that the LLM should never use demographic information when considering the cost of other mutations, so here we tested

if mutating education upwards differed in cost between demographics. Each testing datum in Adult Census and had its education perturbed upwards while considering the datum being male/female, white/not-white, and aged 25/65. Figure 5 shows the results were the custom prompting with these fairness constraints was significantly less biased than the standard prompt alternative in all three demographic features. Specifically, Cohen’s d was 0.34, 0.53, and -1.09 for age, gender, and race, respectively, showing small to large effect sizes. This illustrates that it is possible to make the cost function fairer simply by adding this constraint to the prompt.

4.5 COST FUNCTION FIDELITY

It’s important to understand how accurate the decision tree and MLP cost functions are at imitating the LLM’s reasoning, since we are trying to distill the LLM’s knowledge into small cost function, which can be judged based on how accurately it predicts pairwise comparisons the LLM labeled. Note there is noise in the LLM labels due to its inherent temperature settings, so 100% accuracy would be unwarranted, and indeed some noise has been shown to improve preference learning (Laidlaw & Russell, 2021). Models were trained for 50,000 batches of size 32, and evaluated on the labels of the remaining pairwise comparisons labeled by the LLM described in Section 4. The results can be seen in Figure 4(left). Overall, the custom models always achieved higher accuracy, because there was less noise in their labeling process due to the specific constraints in the prompt. Tree models on average also did better, but this is mostly due to their ability to classify equal cost, which the MLP could not, as it has a non-discrete function output. German Credit performed worse on average also due to the sparser one-hot encoding feature space.⁵ Overall, this illustrates that the cost functions have learned to imitate the original LLM labeling well.

4.6 CASE STUDY

Here we showcase how our cost functions can improve current recourse algorithms by Keane & Smyth (2020) and Wachter et al. (2017), the prior being a data driven approach with an L_1 cost function and the latter a gradient-based method using a median-absolute deviation (MAD) cost function. A simple MLP classifier was trained on Adult Census and achieved 82% accuracy on the training and testing data. **Note we repeated this evaluation on Heloc and German Credit in Appendix I.** The data was standard normalised for a fair comparison between features when checking distances using each method’s default cost function, this was then compared to our custom MLP cost function which was plugged into each method. Adult census was used in all tests with the standard prompting scheme. For full implementation details of each method, the data, and the architectures see Appendix I. In total, we evaluated on the same 6000 instances for each method, and for each of these we attempted recourse generation if they were negative predictions by the model initially. Keane & Smyth (2020) generated 1027 successful recourses, whilst Wachter et al. (2017) was an average of 2322 between methods.

	Male	Age	Native-US	Married	Education	Hours-Work	Private Work	Caucasian
Keane and Smyth (2020) - Data Driven								
L_1	0	554	1	29	266	177	0	0
Ours	0	380	1	30	275	341	0	0
Wachter et al. (2017) - SGD Driven								
MAD	20	53	78	332	1594	4	8	1
Ours	1	28	3	313	1750	174	84	0

Table 1: Case Study Results: Each number represents the number of times each method recommended mutating that feature for recourse. In Keane & Smyth (2020), our method recommended mutating age less and hours-worked more as the main trade-off. In Wachter et al. (2017), our method recommended mutating education and hours-work in comparison to MAD which favored features such as male, and native-US, which are generally considered less actionable.

⁵These two models were also compared against a baseline LLM GPT-4o which was instructed to label every recourse option with a numerical cost value from 0-1, but the results were not competitive and not reported.

486 Table 1 shows the results where we counted how often each method recommended mutating a partic-
487 ular feature to achieve recourse. Broadly speaking, from a feature weighting perspective, the
488 method by Keane & Smyth (2020) prioritized age as a feature, whilst our method reduced this by
489 focusing instead on hours-work. For the method by Wachter et al. (2017) the default MAD cost
490 function suggested many questionable recourses such as changing gender 20x more than our cost
491 function, age almost 2x times, and even race once. However, perhaps the most inactionable feature
492 (native-us) was suggested 78 times, compared to ours which was just 3 times.

494 5 RELATED WORK

496 In the counterfactual literature, early work used the median absolute deviation as a distance func-
497 tion (Wachter et al., 2017), which has some desirable properties such as robustness to outliers, but
498 can’t deal with categorical features or actionability constraints. Early work in this area by Ustun
499 et al. (2019) proposed total and maximum-log percentile shift measures, which can address *relative*
500 *cost*, but not the other desiderata constraints. Other researchers such as Karimi et al. (2020) pro-
501 posed a weighted combination of L_p norms across features, which does deal with *feature cost*, but
502 again misses the other constraints in Section 2. Other recent work continued the use of L_p norms
503 (Karimi et al., 2020; Ramakrishnan et al., 2020), while others investigated HCI questions (Tominaga
504 et al., 2024). In a more recent trend, work has begun to focus on individualised cost (De Toni et al.,
505 2022; Yetukuri et al.; Nguyen et al., 2024), with some also focusing on the Bradley-Terry loss for
506 pairwise comparisons specifically (Rawal & Lakkaraju, 2024), albeit without LLM assistance. In
507 comparison to these works, we are concerned with how to automate the training of high-performing
508 cost functions at scale with LLMs, which follows all the core desiderata constraints in Section 2 laid
509 out in the literature.

510 LLMs have recently been applied to various tasks (Han et al., 2024; Hollmann et al., 2024; Heggel-
511 mann et al., 2023; Borisov et al., 2022), but, here we are focused on utilising their latent knowledge
512 for labeling data for training cost functions, which has not been explored before. Perhaps the most
513 similar work to ours was suggested by Rawal & Lakkaraju (2020). Specifically, they learned a pref-
514 erence function using the Bradley-Terry Loss, pairwise comparisons, and MAP estimates (Hunter,
515 2004; Caron & Doucet, 2012). However, their approach would require human labellers, and doesn’t
516 take *relative* or *dependent* feature cost into account. To help automate similar processes in related
517 areas, recent work has utilised LLMs as judges or evaluators to produce pairwise preferences for
518 learning reward models. Most popularly they are used in RLHF for aligning language models with
519 human preferences (Ouyang et al., 2022), but the ability of this to help with cost functions has not
520 been evaluated until the present work.

521 There is a literature on evaluating how well LLMs correlate with human judgement, but it is difficult
522 to interpret because as much work has shown positive results (Liu et al., 2023; Chiang et al., 2024),
523 as negative (Bavaresco et al., 2024; Koo et al., 2023). Some of this work has highlighted how the
524 discrepancy of results is likely due to a narrow focus on tasks (Bavaresco et al., 2024), suggesting
525 that LLMs may need to be evaluated on very specific use cases to uncover credible ones. Bearing
526 this in mind, the ability of LLMs to correlate with human judgement of cost has not been explored
527 previously, which we addressed in the present paper with our human study.

528 6 CONCLUSION

530 The problem of algorithmic recourse, and counterfactual explanation more broadly, has grown in
531 importance the past several years as AI is increasingly used for high-stakes decisions (Keane et al.,
532 2021; Karimi et al., 2022; Gajcin & Dusparic, 2024; Kothari et al., 2024). However, one of the
533 core unresolved issues plaguing research in the area has been the lack of appropriate cost functions,
534 which has limited the practical value of recourse recommendations. In this paper, we first explored
535 LLM’s natural ability to align with human judgments of cost, showing that they do largely correlate.
536 We then showed that the cost functions can be fine-tuned to fit a variety of use cases. Lastly, we
537 also demonstrated the practical outcomes of using these cost functions in two real-world algorithms.
538 In future work, it would be interesting to investigate the ability of LLMs to learn gain functions for
539 semifactual recourse (Kenny & Huang, 2024), as opposed to counterfactual recourse, which likely
involves other considerations.

REFERENCES

- 540
541
542 Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Al-
543 bert Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, et al. Llms instead
544 of human judges? a large scale empirical study across 20 nlp evaluation tasks. *arXiv preprint*
545 *arXiv:2406.18403*, 2024.
- 546 Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI:
547 <https://doi.org/10.24432/C5XW20>.
- 548
549 Tom Bewley and Freddy Lecue. Interpretable preference-based reinforcement learning with tree-
550 structured reward functions. In *Proceedings of the 21st International Conference on Autonomous*
551 *Agents and Multiagent Systems*, pp. 118–126, 2022.
- 552 Tom Bewley, Salim I Amoukou, Saumitra Mishra, Daniele Magazzeni, and Manuela Veloso. Coun-
553 terfactual metarules for local and global recourse. In *Forty-first International Conference on*
554 *Machine Learning*.
- 555
556 Tom Bewley, Jonathan Lawry, Arthur Richards, Rachel Craddock, and Ian Henderson. Reward
557 learning with trees: Methods and evaluation. *arXiv preprint arXiv:2210.01007*, 2022.
- 558
559 Sebastian Bordt, Harsha Nori, Vanessa Rodrigues, Besmira Nushi, and Rich Caruana. Elephants
560 never forget: Memorization and learning of tabular data in large language models. *arXiv preprint*
561 *arXiv:2404.06209*, 2024.
- 562 Vadim Borisov, Kathrin Seßler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language
563 models are realistic tabular data generators. *arXiv preprint arXiv:2210.06280*, 2022.
- 564
565 Donald E Bowen III, S McKay Price, Luke CD Stein, and Ke Yang. Measuring and mitigating racial
566 bias in large language model mortgage underwriting. *Available at SSRN 4812158*, 2024.
- 567
568 Francois Caron and Arnaud Doucet. Efficient bayesian inference for generalized bradley–terry mod-
569 els. *Journal of Computational and Graphical Statistics*, 21(1):174–196, 2012.
- 570
571 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
572 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena:
573 An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*,
574 2024.
- 575
576 Francesco Bombassei De Bona, Gabriele Dominici, Tim Miller, Marc Langheinrich, and Martin
577 Gjoreski. Evaluating explanations through llms: Beyond traditional user studies. *arXiv preprint*
arXiv:2410.17781, 2024.
- 578
579 Giovanni De Toni, Paolo Viappiani, Stefano Teso, Bruno Lepri, and Andrea Passerini. Personalized
580 algorithmic recourse with preference elicitation. *arXiv preprint arXiv:2205.13743*, 2022.
- 581
582 Jasmina Gajcin and Ivana Dusparic. Redefining counterfactual explanations for reinforcement learn-
583 ing: Overview, challenges and opportunities. *ACM Computing Surveys*, 56(9):1–33, 2024.
- 584
585 Gerd Gigerenzer and Daniel G Goldstein. Reasoning the fast and frugal way: models of bounded
586 rationality. *Psychological review*, 103(4):650, 1996.
- 587
588 Sungwon Han, Jinsung Yoon, Sercan O Arik, and Tomas Pfister. Large language models can au-
589 tomatically engineer features for few-shot tabular learning. *arXiv preprint arXiv:2404.09491*,
590 2024.
- 591
592 Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David
593 Sontag. Tabllm: Few-shot classification of tabular data with large language models. In *Internat-*
ional Conference on Artificial Intelligence and Statistics, pp. 5549–5581. PMLR, 2023.
- Hans Hofmann. Statlog (German Credit Data). UCI Machine Learning Repository, 1994. DOI:
<https://doi.org/10.24432/C5NC77>.

- 594 Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data
595 science: Introducing caafe for context-aware automated feature engineering. *Advances in Neural*
596 *Information Processing Systems*, 36, 2024.
- 597 David R Hunter. Mm algorithms for generalized bradley-terry models. *The annals of statistics*, 32
598 (1):384–406, 2004.
- 600 Mahammed Kamruzzaman and Gene Louis Kim. Prompting techniques for reducing social bias in
601 llms through system 1 and system 2 cognitive processes. *arXiv preprint arXiv:2404.17218*, 2024.
- 602 Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, and Yuichi Ike. Counterfactual explanation
603 trees: Transparent and consistent actionable recourse with decision trees. In *International Con-*
604 *ference on Artificial Intelligence and Statistics*, pp. 1846–1870. PMLR, 2022.
- 606 Amir-Hossein Karimi, Gilles Barthe, Borja Balle, and Isabel Valera. Model-agnostic counterfactual
607 explanations for consequential decisions. In *International conference on artificial intelligence*
608 *and statistics*, pp. 895–905. PMLR, 2020.
- 609 Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A survey of algorithmic
610 recourse: contrastive explanations and consequential recommendations. *ACM Computing*
611 *Surveys*, 55(5):1–29, 2022.
- 613 Mark T Keane and Barry Smyth. Good counterfactuals and where to find them: A case-based tech-
614 nique for generating counterfactuals for explainable ai (xai). In *Case-Based Reasoning Research*
615 *and Development: 28th International Conference, ICCBR 2020, Salamanca, Spain, June 8–12,*
616 *2020, Proceedings 28*, pp. 163–178. Springer, 2020.
- 617 Mark T Keane, Eoin M Kenny, Eoin Delaney, and Barry Smyth. If only we had better counterfactual
618 explanations: Five key deficits to rectify in the evaluation of counterfactual xai techniques. *arXiv*
619 *preprint arXiv:2103.01035*, 2021.
- 620 Eoin Kenny and Weipeng Huang. The utility of “even if” semifactual explanation to optimise posi-
621 tive outcomes. *Advances in Neural Information Processing Systems*, 36, 2024.
- 623 Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop
624 Kang. Benchmarking cognitive biases in large language models as evaluators. *arXiv preprint*
625 *arXiv:2309.17012*, 2023.
- 626 Avni Kothari, Bogdan Kulynych, Tsui-Wei Weng, and Berk Ustun. Prediction without preclusion:
627 Recourse verification with reachable sets. In *The Twelfth International Conference on Learning*
628 *Representations*, 2024. URL <https://openreview.net/forum?id=SCQfYpdoGE>.
- 630 Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with language
631 models. *arXiv preprint arXiv:2303.00001*, 2023.
- 632 Cassidy Laidlaw and Stuart Russell. Uncertain decisions facilitate better preference learning. *Ad-*
633 *vances in Neural Information Processing Systems*, 34:15070–15083, 2021.
- 635 Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg
636 evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.
- 637 Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey
638 on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35, 2021.
- 640 Mstz. Heloc dataset. <https://huggingface.co/datasets/mstz/heloc>, 2024. Ac-
641 cessed: 2024-09-13.
- 642 Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yasmin Schmitt, Jörg
643 Schlötterer, Maurice Van Keulen, and Christin Seifert. From anecdotal evidence to quantitative
644 evaluation methods: A systematic review on evaluating explainable ai. *ACM Computing Surveys*,
645 55(13s):1–42, 2023.
- 646 Duy Nguyen, Bao Nguyen, and Viet Anh Nguyen. Cost-adaptive recourse recommendation by
647 adaptive preference elicitation. *arXiv preprint arXiv:2402.15073*, 2024.

- 648 OpenAI. Hello gpt-4 turbo. <https://openai.com/index/hello-gpt-4o/>, 2024. Accessed: 2024-09-13.
649
650
- 651 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
652 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
653 instructions with human feedback. *Advances in neural information processing systems*, 35:
654 27730–27744, 2022.
- 655 Goutham Ramakrishnan, Yun Chan Lee, and Aws Albarghouthi. Synthesizing action sequences
656 for modifying model decisions. In *Proceedings of the AAAI conference on artificial intelligence*,
657 volume 34, pp. 5462–5469, 2020.
- 658 Kaivalya Rawal and Himabindu Lakkaraju. Beyond individualized recourse: Interpretable and inter-
659 active summaries of actionable recourses. *Advances in Neural Information Processing Systems*,
660 33:12187–12198, 2020.
- 661 Kaivalya Rawal and Himabindu Lakkaraju. Learning recourse costs from pairwise feature compar-
662 isons. *arXiv preprint arXiv:2409.13940*, 2024.
- 663 Tomu Tominaga, Naomi Yamashita, and Takeshi Kurashima. Reassessing evaluation functions in
664 algorithmic recourse: An empirical study from a human-centered perspective. *arXiv preprint*
665 *arXiv:2405.14264*, 2024.
- 666 Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In
667 *Proceedings of the conference on fairness, accountability, and transparency*, pp. 10–19, 2019.
- 668 Julius Von Kügelgen, Amir-Hossein Karimi, Umang Bhatt, Isabel Valera, Adrian Weller, and Bern-
669 hard Schölkopf. On the fairness of causal algorithmic recourse. In *Proceedings of the AAAI*
670 *conference on artificial intelligence*, volume 36, pp. 9584–9594, 2022.
- 671 Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening
672 the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.
- 673 Jayanth Yetukuri, Ian Hardy, and Yang Liu. Actionable recourse guided by user preference.
674
675

676 A ACTIONABILITY CONSTRAINTS AND FEATURES USED

677
678
679 The following are the datasets and features used, alongside any actionability constraints employed
680 throughout the paper:

681
682 **HELOC Dataset:** Here, the actionability constraints were to clamp feature mutations at the highest
683 and lowest values observed in the dataset.

- 684 • *MSinceMostRecentInqexcl7days*: Number of months passed since the last credit inquiry on
685 the individual.
- 686 • *NumRevolvingTradesWBalance*: The number of the individual’s current credit accounts
687 (e.g. credit cards) that have balances on them.
- 688 • *NumTradesOpeninLast12M*: The number of new credit accounts opened in the last 12
689 months.
- 690 • *NumInqLast6M*: The number of credit inquiries carried out on the individual in the last 6
691 months.

692
693 **Adult Census Dataset:** Here, the actionability constraints were to clamp feature mutations at the
694 highest and lowest values observed in the dataset. Also, age and education number were only al-
695 lowed to move upwards.

- 696 • *isMale*: If the person is male, or female, represented as 1 or 0, respectively.
- 697 • *age*: The person’s age, represented as a floating point number.
- 698 • *native-country-United-States*: If the person’s birthplace is the United States, or not, repre-
699 sented as 1 or 0, respectively.

- *marital-status-Married*: If the person is married, or not, represented as 1 or 0, respectively.
- *education-num*: The person’s level of education, represented by a positive integer, where higher numbers are higher levels of education.
- *hours-per-week*: The number of hours the person works per week, represented by a positive integer.
- *workclass-Private*: If the person works for a private company, or is self-employed, represented as 1 or 0, respectively.
- *isCaucasian*: Is the person white or not, represented as 1 or 0, respectively.

German Credit Dataset: Here, the actionability constraints were to clamp numeric feature mutations at the highest and lowest values observed in the dataset.

- *status*: Status of existing checking account.
- *duration*: The proposed duration of the loan in months, expressed as an integer.
- *credit history*: The person’s credit history with the options.
- *purpose*: The purpose of the loan.
- *amount*: The size of the loan asked for.

B PERTURBATION FUNCTION

In the context of feature vector perturbation, we employ a probabilistic approach to introduce controlled mutations to the feature set. Specifically, we perturb a feature vector by altering a random subset of its components. The number of features to be perturbed, denoted as k , is selected from the discrete set $\{1, 2, 3, 4\}$ with a predefined probability distribution. The probability mass function (PMF) for k is given by:

$$P(K = k) = \begin{cases} 0.8 & \text{if } k = 1 \\ 0.1 & \text{if } k = 2 \\ 0.05 & \text{if } k = 3 \\ 0.05 & \text{if } k = 4 \end{cases}$$

This distribution ensures that perturbing a single feature is the most probable event, while perturbing four features is the least probable. The purpose was to focus on sparsity for the cost function training, but also have some robustness. When perturbing numeric features, they have four possible values in our tests. All numeric features can be perturbed upwards one standard deviation, or half a standard deviation. If actionability constraints allow, they can also be perturbed down the same two values, they were then rounded to the nearest integer.

C OUT OF DISTRIBUTION EXPERIMENT

In addition, we were interested in how our cost functions, which were trained to specialise in sparse single feature modifications performed out of distribution when scoring multiple feature mutations in recourse. Hence, we trained a custom tree and MLP model on Adult Census with only 2 feature perturbations allowed, which achieved 88.2% and 87.5% test accuracy on the LLM labels, respectively, and dropped by 82.1% and 81%, respectively. This drop in performance constituted an average of 6.15%, and shows that performance is largely maintained out of distribution, but for maximum effect the training data should represent what is desired in deployment.

D PROMPTS

Here are the prompts for HELOC, all other datasets followed the exact same pattern, and all can be seen in the code base if desired. Note in the actual prompts we instructed the LLM to use 1, 2, 0,

756 to select Recourse 1, Recourse 2, and neither, respectively, although in the main paper we used 1, 0,
757 and 0.5, as this more accurately reflected the Bradley-Terry model.

758
759 **Standard prompt (\mathcal{B}):**

760 You are a helpful assistant to a data scientist to help them
761 label data. You will be shown a data point representing a person
762 Alex, and a mutation of it, You will also be shown a data point
763 representing a person Jaden, and a mutation of it, your task is
764 to label which of the two mutations would take more effort to
765 achieve.

766 The data will be the HELOC Dataset which uses these features:

767 MSinceMostRecentInqexcl7days: Number of months passed
768 since the last credit inquiry on the individual.
769 NumRevolvingTradesWBalance: The number of the individual's
770 current credit accounts (e.g. credit cards) that have balances on
771 them. NumTradesOpeninLast12M: The number of new credit accounts
772 opened in the last 12 months. NumInqLast6M: The number of credit
773 inquiries carried out on the individual in the last 6 months.

774 The data is represented in array form like ['MSinceMostRecentInqexcl7days',
775 'NumRevolvingTradesWBalance', 'NumTradesOpeninLast12M',
776 'NumInqLast6M']

777
778 Now consider the following individual Alex: ""+str(x1)"" Now
779 consider this mutation of Alex: ""+str(x1p)""

780 Now consider the following individual Jaden: ""+str(x2)"" Now
781 consider this mutation of Jaden: ""+str(x2p)""

782 Which of these two mutations would take more effort? You must
783 provide an answer.

784
785 Remember the following 4 rules and use them in your decision:

786 1. Some features are naturally harder to change than others, use
787 this logic.

788 2. For numerical features, the difficulty of changing them can
789 often depend on their starting values.

790 3. Apart from the mutated features, consider the other features
791 which are different between Alex and Jaden, and how this may
792 affect difficulty.

793 4. Do not ever use demographic features (e.g., age, gender, race)
794 when considering the difficulty of mutating other features.

795
796 Outline your reasoning process step by step, before giving your
797 answer as 1, 2, or 0 in the tags <answer>...</answer>, where 1
798 means you think the first mutation requires more effort, 2 means
799 you think the second mutation requires more effort, and 0 means
800 you think there is no difference.

801
802 **Custom Prompt ($\mathcal{B} + \mathcal{B}'$):**

803 You are a helpful assistant to a data scientist to help them
804 label data. You will be shown a data point representing a person
805 Alex, and a mutation of it, You will also be shown a data point
806 representing a person Jaden, and a mutation of it, your task is
807 to label which of the two mutations would take more effort to
808 achieve.

809 The data will be the HELOC Dataset which uses these features:

810 MSinceMostRecentInqexcl7days: Number of months passed
811 since the last credit inquiry on the individual.
812 NumRevolvingTradesWBalance: The number of the individual's
813 current credit accounts (e.g. credit cards) that have balances on
814 them. NumTradesOpeninLast12M: The number of new credit accounts
815 opened in the last 12 months. NumInqLast6M: The number of credit
816 inquiries carried out on the individual in the last 6 months.
817 The data is represented in array form like ['MSinceMostRecentInqexcl7days',
818 'NumRevolvingTradesWBalance', 'NumTradesOpeninLast12M',
819 'NumInqLast6M']
820 Now consider the following individual Alex: ""+str(x1)"" Now
821 consider this mutation of Alex: ""+str(x1p)""
822
823 Now consider the following individual Jaden: ""+str(x2)"" Now
824 consider this mutation of Jaden: ""+str(x2p)""
825 Which of these two mutations would take more effort?
826
827 Remember the following:
828 1. The hardest features to change, in order from the
829 hardest to easiest are [MSinceMostRecentInqexcl7days,
830 NumRevolvingTradesWBalance, NumTradesOpeninLast12M, NumInqLast6M]
831 2. For the numerical features, they are all harder to increase
832 the higher they get.
833 3. If NumInqLast6M is greater than zero, then increasing
834 'NumTradesOpeninLast12M' becomes more difficult.
835
836 Outline your reasoning process step by step, before giving your
837 answer as 1, 2, or 0 in the tags <answer>...</answer>, where 1
838 means you think the first mutation requires more effort, 2 means
839 you think the second mutation requires more effort, and 0 means
840 you think there is no difference.
841 **Prompt to elicit numerical response from LLM:**
842 You are a helpful assistant to a data scientist that helps them
843 label data. You will be shown a data point representing a person
844 Alex, and a mutation of it. your task is to label how much effort
845 this mutation was to achieve using a number between 0 and 1, where
846 0 is no effort, and 1 is the most possible effort.
847
848 The data will be the HELOC Dataset which uses these features:
849 MSinceMostRecentInqexcl7days: Number of months passed
850 since the last credit inquiry on the individual.
851 NumRevolvingTradesWBalance: The number of the individual's
852 current credit accounts (e.g. credit cards) that have balances on
853 them. NumTradesOpeninLast12M: The number of new credit accounts
854 opened in the last 12 months. NumInqLast6M: The number of credit
855 inquiries carried out on the individual in the last 6 months.
856 The data is represented in array form like ['MSinceMostRecentInqexcl7days',
857 'NumRevolvingTradesWBalance', 'NumTradesOpeninLast12M',
858 'NumInqLast6M']
859 Now consider the following individual Alex: ""+str(x1)"" Now
860 consider this mutation of Alex: ""+str(x1p)""
861
862 Using a floating point number between 0 and 1, how much effort was
863 this to achieve? You must provide an answer.

864 Outline your reasoning process step by step before giving your
865 answer in the tags <answer>...</answer>

866 **Human Study Prompt:**

867
868 You are a helpful assistant to a data scientist to help them
869 label data. You will be shown a data point representing a person
870 Alex, and a mutation of it, You will also be shown a data point
871 representing a person Jaden, and a mutation of it, your task is
872 to label which of the two mutations would take more effort to
873 achieve.

874 The data will be the HELOC Dataset which uses these features:

875 MSinceMostRecentInqexcl7days: Number of months passed
876 since the last credit inquiry on the individual.
877 NumRevolvingTradesWBalance: The number of the individual's
878 current credit accounts (e.g. credit cards) that have balances on
879 them. NumTradesOpeninLast12M: The number of new credit accounts
880 opened in the last 12 months. NumInqLast6M: The number of credit
881 inquiries carried out on the individual in the last 6 months.

882 The data is represented in array form like ['MSinceMostRecentInqexcl7days',
883 'NumRevolvingTradesWBalance', 'NumTradesOpeninLast12M',
884 'NumInqLast6M']

885
886 Now consider the following individual Alex: ""+str(x1)+" Now
887 consider this mutation of Alex: ""+str(x1p)+"

888 Now consider the following individual Jaden: ""+str(x2)+" Now
889 consider this mutation of Jaden: ""+str(x2p)+"

890 Which of these two mutations would take more effort? You must
891 provide an answer.

892 Remember the following 4 rules and use them in your decision:

893
894 1. Some features are naturally harder to change than others, use
895 this logic.

896
897 2. For numerical features, the difficulty of changing them can
898 often depend on their starting values.

899
900 3. Apart from the mutated features, consider the other features
901 which are different between Alex and Jaden, and how this may
affect difficulty.

902
903 4. Do not ever use demographic features (e.g., age, gender, race)
when considering the difficulty of mutating other features.

904 Outline your reasoning process step by step, before giving your
905 answer as 1, 2, or 0 in the tags <answer>...</answer>, where 1
906 means you think the first mutation requires more effort, 2 means
907 you think the second mutation requires more effort, and 0 means
908 you think there is no difference.

909 **Prompt to acquire ground truth dependencies:**

910 ... (insert the previous standard prompt here) ...

911
912 In the above problem, what are the primary feature dependencies
913 that may effect effort?

914
915
916
917

E HUMAN STUDY QUESTION EXAMPLE

Here we supply an example question from the human study for the HELOC dataset. The full survey can be seen in the supplement.

Perceived Effort Required in Dataset Feature Mutations

HELOC Dataset

Data Description:
 The FICO HELOC dataset contains anonymized information about home equity line of credit (HELOC) applications made by real homeowners. The customers in this dataset have requested a credit line in the range of USD 5,000 - 150,000.

Selected Features:

1. **Months Since Recent Inquiries:** Number of months passed since the last credit inquiry on the individual.
2. **Number of Credit Accounts with Balances:** The number of the individual's current credit accounts (e.g. credit cards) that have balances on them
3. **Number of New Credit Accounts:** The number of new credit accounts opened in the last 12 months
4. **Number of Inquiries:** The number of credit inquiries carried out on the individual in the last 6 months

* 2.

Alex	Months since Recent Inquiry	Number of Credit Accounts with Balances	Number of New Credit Accounts	Number of Inquiries
Features	6	3	2	4
Change(s) to Make			4	

Jaden	Months since Recent Inquiry	Number of Credit Accounts with Balances	Number of New Credit Accounts	Number of Inquiries
Features	6	3	2	4
Change(s) to Make		5		

Which individual's proposed change would require more effort?

Alex
 Jaden

* 3. In your opinion, how much difference in effort do you perceive between the two changes (for Alex & Jaden) in the scenario above?

1 (Almost No Difference) 7 (Really Large Difference)

Figure 6: Human study question example.

F UNCERTAINTY IN HUMAN STUDY

Here we consider the percentage of replies from the LLM in the Section 4 human study for each question which were uncertain (i.e., it chose the third option rather than Recourse 1 or 2), alongside the average response humans gave for the distance between the two recourse in Figure 6. For both lists, we normalized each to be between 0-1, and plotted them in a scatter plot to see the correlation. Both lists represent each group’s uncertainty in choosing a recourse, and shows how the LLMs and humans correlate in this aspect to a high degree (Person’s $r=0.48$; $p < 0.04$). What this tells us is that how uncertain humans are on these recourse questions in the human study strongly correlate. Note that the correlation is identical with un-normalized scores also, we just do so here for clarity and visual purposes.

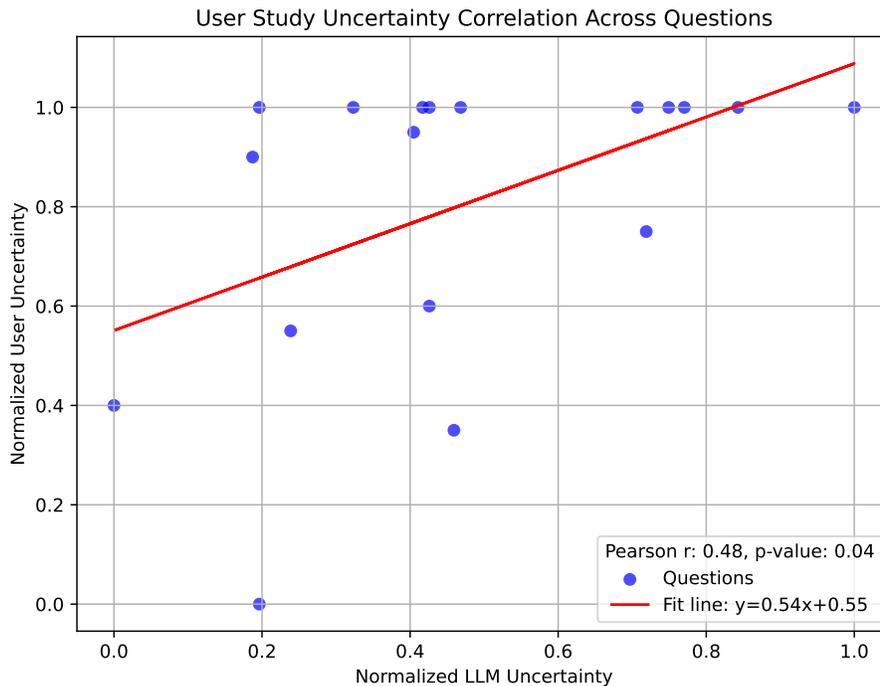


Figure 7: The correlation between LLM uncertainty and human uncertainty in the human study shows both groups were similarly uncertain on each question.

G FEATURE DEPENDENCIES

In Section 4.3 we evaluated how well various prompt types picked up on feature dependencies. To acquire these dependencies in an objective way, we queried Claude Sonnet 3.5 to list all feature dependencies in all 3 datasets using the standard prompt in Section D and adding:

...

If I change the same feature in Alex and Jaden the same amount, but another feature is different which effects the effort involved, what would be the most likely dependencies like this to happen?

Question	P-value	Same Most Common Response	LLM Uncertainty %
1	0.27	True	0.00
2	0.00	True	0.00
3	0.26	True	7.14
4	0.18	False	42.86
5	1.00	True	28.57
6	1.00	True	46.43
7	1.00	True	0.00
8	0.53	True	0.00
9	0.35	True	0.00
10	0.53	True	3.57
11	0.47	True	0.00
12	0.59	False	71.43
13	0.00	False	0.00
14	0.03	True	0.00
15	0.06	True	0.00
16	0.56	True	32.14
17	0.12	True	0.00
18	0.24	True	17.86

Table 2: Results of 18 Question in Human Study: We are looking to see which have statistically similar distributions or the same most common response as a sign of LLM alignment with humans in judgement of cost. Overall, 17/18 show one metric or the other with positive results.

We repeated this 10 times and took the three dependencies which occurred most often, these were:

HELOC:

1. If NumTradesOpeninLast12M is low, it makes it more challenging to increase NumRevolvingTradesWBalance.
2. If NumInqLast6M is high, it suggests it should be more difficult to increase NumTradesOpeninLast12M.
3. If NumInqLast6M is high, it should be more difficult to increase NumRevolvingTradesWBalance.

Adult:

1. Increasing working hours should be more difficult if you are married.
2. Changing marital status should be more difficult the older you are.
3. Increasing working hours should be more difficult if working for a private company.

German Credit:

1. A bad credit history should make it harder to increase your loan amount.
2. A bad credit history should make it harder to increase your duration.
3. The effort to decrease the duration of a loan should be harder for larger loan amount.

H GROUND TRUTH FOR SECTION 4.4

We had to define a ground truth for our fine-tuning experiments to see how we could manipulate the four desiderata outlined previously. Note Desideratum 4 (i.e., fair cost) was only evaluate on Adult due to its numerous demographic features. The ground truth defined in B' for each dataset was:

HELOC:

... 1. The hardest features to change, in order from the hardest to easiest are [MSinceMostRecentInqexcl7days, NumRevolvingTradesWBalance, NumTradesOpeninLast12M, NumInqLast6M]

1080 2. For the numerical features, they are all harder to increase
1081 the higher they get.

1082 3. If NumInqLast6M is greater than zero, then increasing
1083 'NumTradesOpeninLast12M' becomes more difficult.
1084

1085 **Adult Census:**

1086 ... 1. The hardest features to change, in order from the hardest
1087 to easiest are [native-country-United-States, isWhite, isMale,
1088 age, marital-status-Married, education-num, workclass-Private,
1089 hours-per-week]

1090 2. For age, education-num, and hours-per-week, they are all
1091 harder to increase the higher they get.

1092 3. Increasing hours-per-week is more effort if the person works
1093 for a private company.

1094 4. Never use demographic information (i.e., isMale, age, isWhite)
1095 when calculating the effort of other feature changes.
1096

1097 **German Credit:**

1098 ... 1. The hardest features to change, in order from the hardest
1099 to easiest are [credit history, status, purpose, duration, amount]
1100

1101 2. For the numerical features, they are all harder to increase
1102 the higher they get.

1103 3. Having bad credit history or bad status makes it harder to
1104 increase amount.
1105

1106 These dependencies were chosen to be fine-tuned because they performed badly in Section 4.3.
1107

1108 I BASELINE IMPLEMENTATION DETAILS

1109 This section serves to give full details about the implementation of Keane & Smyth (2020) and
1110 Wachter et al. (2017) in Section 4.6. The data used was Adult Census with 30,000 for training, and
1111 6,000 for testing the recourse generation.
1112

1113 I.1 KEANE AND SMYTH

1114 This method is data driven and works by defining a case-base of recourse options for training
1115 data (Keane & Smyth, 2020). In practice, each training data has its nearest unlike neighbor found in
1116 the case-base and the difference between the two is taken as one recourse option. Recourses of 2 or
1117 less feature changes are preferred by the authors, we focus on single feature changes. At test time, a
1118 query has its nearest neighbour found in the case base and its recourse is applied to the query, this is
1119 repeated for all nearest neighbours to find the best recourse option adhering to some constraints. For
1120 us, these constraints are a single feature mutation, and that the result must be a valid counterfactual.
1121 Finally, we also considered the 100 nearest neighbours as possible recourses.
1122

1123 I.2 WACHTER ET AL.

1124 A heavily implemented framework in research (Wachter et al., 2017), the method works by gener-
1125 ating a set of random recourses which optimize to be closer to the query, while optimizing to also
1126 be the counterfactual class. The second constraint is gradually up-weighted with a lambda term to
1127 be more important throughout several optimization steps. We implement the method as normal with
1128 300 possible counterfactuals during optimization, categorical features are snapped to the closest real
1129 value, the results are filtered to those which are valid counterfactuals, and the closest chosen as the
1130 answer. Because we are interested in sparse explanations, we also clamp each possible counterfac-
1131 tual to have one possible feature mutation, which in practice is done allowing the largest currently
1132 mutated feature to be the recommended recourse action.
1133

J CASE-STUDY EXTRA RESULTS

To complete our case study in Section 4.6, we add the two other datasets in the paper. We focused on Adult Census in the main paper because it is less debatable what the most actionable features are.

	MSinceMostRecentInqexcl7days	NumRevolvingTradesWBalance	NumTradesOpeninLast12M	NumInqLast6M
Keane and Smyth (2020) - Data Driven				
L_1	336	9	1	0
Ours	270	43	8	25
Wachter et al. (2017) - SGD Driven				
MAD	167	1	0	0
Ours	5	10	5	148

Table 3: Heloc Results: On average the baselines favored Months Since Most Recent Inquiry Excluding 17 days, in contrast to our cost function which favored Number of inquiries in the last 6 months and number of revolving trades with balance as a trade-off. Considering the first feature has a time constraint, it is immediately more actionable to modify the feature our cost function chose. Our cost function also generally offers a more diverse set of explanations.

	Repayment Term	Loan Amount	Status	Credit History	Purpose
Keane and Smyth (2020) - Data Driven					
L_1	0	37	0	0	0
Ours	0	37	0	0	0
Wachter et al. (2017) - SGD Driven					
MAD	1	1	22	34	51
Ours	19	1	1	51	38

Table 4: German Credit Results: Keane and Smyth performed poorly on this dataset because (1) the dataset itself is smaller than the others (666 training), and is heavily imbalanced (95/5%), hence because it is a data driven method which directly uses the data for computation, there was sparse examples of how to generate counterfactuals. In Wachter et al. (2017) our method favored *Repayment Term* and *Credit History* in comparison to MAD which focused on *Status* and *Purpose*. Arguably, *Repayment Term* is the easiest feature to modify, as *Status* involves changing your savings amount which is quite costly when increasing, *Credit History* by comparison is easier to change but takes time, and *Purpose* which involves major changes to ones future plans.

K DEPENDENCY EXPERIMENT WITH SYNTHETIC DATA

Due to our datasets in the paper being popular recourse datasets, it is reasonable to assume that there are counterfactual pairs the LLM has seen during pre-training. Hence, to verify that it can detect causal dependencies without having seen direct examples from datasets, we create a synthetic dataset which does not exist anywhere, and thus cannot be part of the LLM’s pre-training data.

We generated a medical dataset of personal information which is unlikely to have similar publicly available datasets used in recourse papers online. The features we used were:

- Saturated Fat Intake: the amount of saturated fat they eat, from 0 to 100 grams..
- Salt Intake: the amount of salt fat they eat, from 0 to 10 grams.
- Processed Food Intake: the amount of saturated fat they eat, from 0 to 500 grams.
- Cholesterol Level: Cholesterol level in mg/dl.
- Blood Pressure: their systolic blood pressure.
- Weight: weight in KG.

1188 We chose these features because there are three scientifically known dependencies we can use as a
 1189 ground truth.

- 1190 1. It is harder to lower cholesterol if your saturated fat is too high.
- 1191 2. It is harder to lower blood pressure if your salt intake is too high.
- 1192 3. It is harder to lose weight if your intake of processed food is too high.

1193 We generated this dataset according to reasonable values seen in the below script

```

1194
1195 import numpy as np
1196 import pandas as pd
1197
1198 np.random.seed(42)
1199 n_samples = 1000
1200
1201 saturated_fat_intake = np.random.randint(0, 101, n_samples) # 0 to 100 grams
1202 salt_intake = np.random.randint(0, 11, n_samples) # 0 to 10 grams
1203 processed_food_intake = np.random.randint(0, 501, n_samples) # 0 to 500 grams
1204
1205 cholesterol_level = np.round(200 + 0.5 * saturated_fat_intake +
1206 np.random.normal(0, 10, n_samples)).astype(int)
1207 blood_pressure = np.round(120 + 2 *
1208 salt_intake + np.random.normal(0, 5, n_samples)).astype(int)
1209 weight = np.round(70 + 0.1 * processed_food_intake +
1210 np.random.normal(0, 5, n_samples)).astype(int)
1211
1212 mortality_risk = ((cholesterol_level > 240) |
1213 (blood_pressure > 140) | (weight > 100)).astype(int)
1214
1215 data = pd.DataFrame({
1216     'Saturated Fat Intake': saturated_fat_intake,
1217     'Salt Intake': salt_intake,
1218     'Processed Food Intake': processed_food_intake,
1219     'Cholesterol Level': cholesterol_level,
1220     'Blood Pressure': blood_pressure,
1221     'Weight': weight,
1222     'Mortality Risk': mortality_risk
1223 })
1224
1225

```

1226 With this dataset in hand we iterated each datum 3 times and made the following 3 mutations each
 1227 time to test each dependency. In each case the datum was duplicated to control for other features
 1228 and focus only on the dependencies across a diverse range of data.

- 1229 1. We set saturated fat to 10g v 100g, and took the action of lowering cholesterol by 10.
- 1230 2. We set salt intake to 1g or 15g, and took the action of decreasing blood pressure by 10.
- 1231 3. We set processed food intake to 10g v. 1500g, and took the action of losing 5KG of weight.

1232 In all cases these mutations had additional random noise added to them for robustness. The LLM
 1233 was allowed to answer that the first, second, or neither recourse was higher effort. In all cases,
 1234 recourse 2 was the ground truth, hence, if most of the LLM responses are Recourse 2, then it is
 1235 picking up on the causal dependencies. Lastly, we compared a standard prompting scheme with no
 1236 information, and the same prompt with a high-level overview of the desiderata in Section 2. Results
 1237 are shown in Figure 8. Overall, when adding the high-level desiderata to the prompt, the LLM can
 1238 always detect these known causal dependencies with very high accuracy. This shows the LLM is
 1239 capable of reasoning about causal dependencies without being exposed to similar training data in
 1240 the past. Moreover, what is particularly interesting is that by explicitly telling the LLM to consider
 1241

other dependencies (by adding the desiderata to the prompt), it is able to do that. However, without being told to consider dependencies in the prompt, it is not able to reason correctly.

In short, this experiment tells us two important things. First, LLMs can reason about causal dependencies it has not been exposed to before in terms of counterfactual data available on the internet. Secondly, in order to do this, the desiderata from Section 2 must be added to the prompt. Note, this is a general desiderata, not dataset specific, it does not make the method less general.

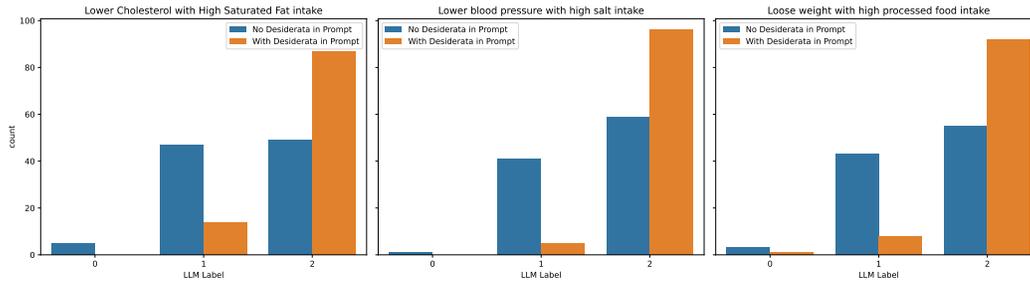


Figure 8: The labeling of the LLM on our synthetic causal relationships. The ground truth always corresponds to Recourse option 2. In general, the LLM was capable of modeling the causal dependencies with 90% accuracy. When ablating the desiderata from the prompt, this reduced to near random guessing between the two recourse options. Note, 0 corresponds to the LLM assigning equal cost to both recourses.