
Data Augmentations for Improved (Large) Language Model Generalization

Amir Feder ^{*1,2}, Yoav Wald ^{*3}, Claudia Shi ¹, Suchi Saria ³ and David Blei ¹

¹ Columbia University, ² Google Research, ³ Johns Hopkins University

Abstract

The reliance of text classifiers on spurious correlations can lead to poor generalization at deployment, raising concerns about their use in safety-critical domains such as healthcare. In this work, we propose to use counterfactual data augmentation, guided by knowledge of the causal structure of the data, to simulate interventions on spurious features and to learn more robust text classifiers. We show that this strategy is appropriate in prediction problems where the label is spuriously correlated with an attribute. Under the assumptions of such problems, we discuss the favorable sample complexity of counterfactual data augmentation, compared to importance re-weighting. Pragmatically, we match examples using auxiliary data, based on diff-in-diff methodology, and use a large language model (LLM) to represent a conditional probability of text. Through extensive experimentation on learning caregiver-invariant predictors of clinical diagnoses from medical narratives and on semi-synthetic data, we demonstrate that our method for simulating interventions improves out-of-distribution (OOD) accuracy compared to baseline invariant learning algorithms.

1 Introduction

The reliance on spurious correlations is a significant challenge for Machine Learning (ML) safety as it can lead to performance degradation of deployed models. Spurious correlations are prevalent in various applications such as medical imaging [1, 2], text classification [3], and risk prediction systems [4]. Failures due to spurious correlations occur under distribution shift [5–7], which may result from differences in data recording protocols, shifts in the underlying population being monitored, or the way the ML tool is being used. In this paper, we focus on text classification and explore how using language models in a domain-informed way can help us avoid reliance on spurious correlations.

Consider a scenario where we want to make robust predictions about patients’ conditions, probability of readmission, etc., using clinical narratives written in hospitals [8–10]. In this setting, a common issue arises due to clinical practice, where patients with certain conditions are directed to specific caregivers in the hospital. When we train a predictor from a single dataset that exhibits some correlation between caregiver-specific style and clinical outcomes, the predictor may unintentionally rely on the style to make predictions. This leads to poor generalization on unseen hospitals, i.e. failure to generalize out of distribution(OOD), due to changes in clinical practice [7]. However, collecting a dataset that is large enough to avoid such spurious associations is infeasible due to various reasons such as rare conditions, privacy concerns, etc. To tackle this problem, we propose leveraging available auxiliary data (e.g., time, document type, demographics) and incorporating knowledge about the causal structure of the problem to build a more robust classifier. For example, in the note classification task, we can use our knowledge that some auxiliary data, such as the patient’s current state, can affect doctor assignment, to improve the classifier’s robustness.

*Equal Contribution. Correspondence to amir.feder@columbia.edu

Causal inference often makes use of such auxiliary data and has now been used in a variety of ways to improve OOD generalization [6, 11–14]. Data augmentation methods have demonstrated impressive performance in these tasks as well [15–17], and with recent improvements in generative models, forming additional principles to incorporate domain knowledge into data augmentations seems like a promising path forward.

In this work we pursue this and develop *causally-driven data augmentation methods*, that leverage auxiliary data and domain knowledge. Intuitively, generating versions of clinical narratives as if they had been written by different caregivers, de-correlates the writing style from the patient condition we wish to predict. However, such data generation can be difficult to achieve in practice and problem-specific traits must be taken into account [18]. Observing that data augmentation can be treated as counterfactual outcome estimation under a causal formalism, motivates the use of causal inference methods that are commonly used for such tasks across the sciences. While our approach can be applied to many modalities of data, in this work we focus on text classification and harness the recent advances in LLMs towards counterfactual estimation. Our contributions are:

1. Through extensive experiments, we show how the use of language models in a manner that is informed by causal knowledge improves model robustness in challenging safety-critical tasks in healthcare. Furthermore, our findings are reinforced by experiments that incorporate semi-synthetic scenarios, and simulations where there are ground-truth counterfactuals.
2. We formalize counterfactual data augmentation in a prediction setting as a method to deconfound the target and a spuriously correlated attribute. We show how deconfounding improves OOD generalization. In a setting where sample complexities for alternative methods (re-weighting and invariance penalties) can be derived, we show favorable generalization bounds for accurately performed data-augmentation.
3. Our data-augmentation methods rely on common assumptions in the causal inference literature such as no unmeasured confounding and parallel trends in diff-in-diff [19], applied with LLMs. We believe that leveraging auxiliary data and assumptions about causal structure, along with the use of LLMs and other generative models, can be a fruitful framework for addressing many out-of-distribution generalization problems.

Next, we provide a brief survey of relevant work (§2). We then present a formal setting motivating counterfactual augmentation for OOD generalization (§3), our methods for counterfactual estimation and reason formally about the preferable sample complexity of our approach (§4). Finally, we present our main experimental results (§5) and discuss limitations and future directions (§6).

2 Related Work

Invariant and Shift-stable Learning. This paper contributes to the growing literature on invariant and shift-stable learning, which tackles the problem of learning models that generalizes across different distributions or settings. Invariant learning through feature pruning was pioneered by Peters et al. [11], and has since been developed for variable selection [12, 20] and representation learning [13, 21–26]. These methods have been applied in a range of domains, including natural science [11, 12, 20], causal estimation [27, 28], computer vision [13, 23], and NLP [29–32]. However, recent studies have highlighted limitations in many invariant learning approaches, particularly in achieving conditional independence [33–36]. Others have investigated learning of stable models by leveraging causal methods through techniques like graph-surgery [6, 14], that come with generalization guarantees. Yet others have explored the advantages of data augmentation [37, 38]. In this work, we combine the latter two approaches to improve OOD generalization for text based classification.

Counterfactually Augmented Data. To learn invariant predictors, a popular and straightforward approach is data augmentation. When data augmentation involves actions that go beyond simple manipulations (e.g. image rotations, crops etc.), it is often referred to as *counterfactual data augmentation* [37]. Constructing counterfactual instances that involve perturbations to confounding factors [39], or to the label [37, 38, 40], and incorporating them into the training data, breaks up correlations that we do not wish our model to exploit towards prediction. Most work on counterfactual data augmentation in text involve manual editing by humans, heuristic keyword replacement, or automated text rewriting [37, 39, 41–50]. Manual editing is accurate and effective [38, 51] but expensive, hence our goal is to *make counterfactual data augmentation scalable*, demanding smaller

human effort. Keyword-based methods can be limited in coverage and difficult to generalize across languages [52]. Generative approaches offer a balance of fluency and coverage [53], but generating meaningful counterfactuals is challenging [54]. Our work departs from previous techniques by using *causal auxiliary data structure and LLMs* to alleviate this challenge and generate plausible counterfactual data augmentations.

Clinical Notes. Clinical notes are the backbone of electronic health records, often containing vital information not observed in other structured data Kreimeyer et al. [55]. Clinical NLP involves identifying this information, and standardized datasets and competitions exist for this purpose [56–60]. Best performing approaches have leveraged transformer architectures both for token-level classification tasks [61–64], and for using complete clinical records [65, 66]. Recently, large language models (LLMs), similar to those we use to generate counterfactual notes, were shown to have clear potential for improving clinical NLP systems [67, 68]. In our experiments, we follow recent papers in clinical NLP addressing challenges of degraded performance across different hospitals [69–71].

3 Problem Setting

To formally analyze how counterfactual data augmentation helps OOD generalization, we consider a setting where the label is spuriously correlated with a known attribute. This setting has been used previously to study learning with “shortcuts” [25] and spurious correlations [29]. We note that our approach is applicable and valid under additional settings and causal graphs (e.g. “purely spurious” problems defined in Wang and Veitch [72]) and we elaborate on this at ???. The data generating process used here motivates counterfactual data augmentation in a principled manner, as it describes the main problem we study and it is possible to analytically compare sample complexity with an alternative solution (see section 4.3).

Consider a classification problem with L classes, where the label Y is correlated with a certain attribute C in the training data and this correlation may change arbitrarily at test time (denoted by a red edge $C \leftrightarrow Y$ in fig. 1). In our medical notes example, C is the caregiver writing the note and Y is the underlying condition we wish to diagnose. We denote the number of caregivers in our training data by $[K]$. For a given loss function $\ell : \mathbb{R}^L \times [L] \rightarrow \mathbb{R}$ and distribution P , we denote the expected loss of a hypothesis $h : \mathcal{X} \rightarrow \mathbb{R}^L$ by $\mathcal{R}_P^\ell(h)$ and its expected accuracy by $\mathcal{R}_P^{\ell_{01}}(h)$. The data-generating process is depicted by the causal model in fig. 1, for our motivating example of clinical notes classification X is a vector representation of the clinical note and X^* is an unobserved sufficient statistic, representing all the relevant information about Y in the note that is unaffected by the writing style of the caregiver. Let us formally define this setting.

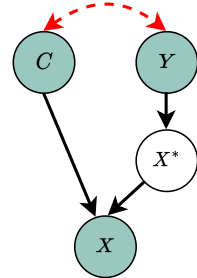


Figure 1: Prediction problem with a spuriously correlated attribute.

Definition 1. We denote the set of distributions induced by interventions on a causal model with the structure in fig. 1 by

$$\mathcal{P} = \{P(X | X^*, C)P(X^* | Y)P(Y)\tilde{P}(C | Y) : \tilde{P}(C | Y = y) \in \Delta^{K-1} \forall y \in [L]\},$$

where all distributions other than $\tilde{P}(C | Y)$ are fixed. In a prediction problem with a spuriously correlated attribute, the learner is provided with a set $\{(\mathbf{x}_i, y_i, c_i)\}_{i=1}^N$ sampled i.i.d from $P_{train} \in \mathcal{P}$. We assume that $X^* = e(X)$ almost surely for some $e : \mathbb{R}^d \rightarrow \mathbb{R}^{d^*}$.

In this problem, once X^* is recovered no additional information from X is needed to predict Y . We can also see from the graph that interventions on $\tilde{P}(C | Y)$ do not change the conditional distribution $P(Y | X^*)$. Therefore an optimal solution that does not rely on C is $h^*(\mathbf{x}) = \arg \max_{y \in [L]} P(Y = y | e(\mathbf{x}))$. In clinical note classification, X^* represents all the information in the note about the patient conditions, unswayed by the writing style of caretaker C . To obtain $h^*(\mathbf{x})$ we will rely on risk minimization w.r.t a distribution where Y and C are uncorrelated.

3.1 Learning Robust Classifiers when Counterfactuals are Available

Consider the unconfounded distribution $P_{\perp} \in \mathcal{P}$ that is given by intervening on C , setting it independent of Y and uniformly distributed, $\hat{P}(C | Y) = P_{\text{unif}}(C)$. An optimal classifier under P_{\perp} has the following min-max optimality guarantee.²

Lemma 1. *For the prediction problem in definition 1, the Bayes optimal classifier under the unconfounded distribution $P_{\perp} \in \mathcal{P}$ where C is uniformly distributed and independent of Y is $h^*(\mathbf{x}) = \arg \max_{y \in [K]} P_{\perp}(Y = y | X^* = e(\mathbf{x}))$. It is a minimizer of $\min_{h: \mathcal{X} \rightarrow [L]} \max_{P \in \mathcal{P}} \mathcal{R}_P^{\ell_{01}}(h)$ and $\mathcal{R}_P^{\ell_{01}}(h^*) = \mathcal{R}_{P_{\perp}}^{\ell_{01}}(h^*)$ for all $P \in \mathcal{P}$.*

Hence we would like to minimize risk w.r.t P_{\perp} and we cannot do that directly by via ERM since our training data is sampled from $P_{\text{train}} \neq P_{\perp}$. Instead we consider risk minimization over an augmented dataset that contains counterfactual instantiations of our training data under different values of C .

Minimizing $\mathcal{R}_{P_{\perp}}$ via Counterfactual Data Augmentation. Returning to our motivating example, assume that we could generate clinical notes for all alternative scenarios. That is, obtain the clinical notes that would have been written if each patient had been seen by all possible caregivers $c \in [K]$ and each caregiver had written their own version of the note $\mathbf{x}_i(c)$. Given these counterfactual clinical notes, we seek a hypothesis that minimizes the average loss over all such possible scenarios, denoted by $\widehat{\mathcal{R}}_{\text{aug}}^{\ell}(h)$.

Definition 2. *Consider a prediction problem with a spuriously-correlated attribute (see Definition 1). For a given example \mathbf{x}_i , we denote its counterfactual with attribute value $c \in [K]$ as derived from the corresponding causal model, by $\mathbf{x}_i(c)$. For estimates of the counterfactuals $\{\hat{\mathbf{x}}_i(c)\}_{i \in [N], c \in [K]}$ and a hypothesis $h \in \mathcal{H}$, the counterfactually augmented empirical risk is*

$$\widehat{\mathcal{R}}_{\text{aug}}^{\ell}(h) = \frac{1}{NK} \sum_{i \in [N], c \in [K]} \ell(h(\hat{\mathbf{x}}_i(c)), y_i). \tag{1}$$

We use approximate counterfactuals $\hat{\mathbf{x}}_i(c)$ in our definition to highlight that in practice we cannot obtain a precise estimate of $\mathbf{x}_i(c)$. In the ideal case where $\hat{\mathbf{x}}_i(c) = \mathbf{x}_i(c)$, the expected loss $\mathcal{R}_{\text{aug}}^{\ell}(h)$ where $N \rightarrow \infty$, satisfies $\mathcal{R}_{\text{aug}}^{\ell}(h) = \mathcal{R}_{P_{\perp}}^{\ell}(h)$. This follows by a simple derivation and it is part of a claim we give later in Lemma 2. Hence obtaining this dataset is useful for our goal of minimizing risk under P_{\perp} . Our main challenge is then to derive effective approximations for counterfactuals such as clinical notes under alternative writing styles.

4 Assumptions and Algorithms for Estimating Counterfactuals

Perfectly capturing writing style is a strong assumption. Even if we could perfectly model writing styles, we only observe a limited set of variables - the actual notes x , outcomes y , and assigned caregivers c . We do not observe all factors that could influence what each caregiver would write. To alleviate this problem, we make use of auxiliary data M that is available during training, but might not be available in deployment.

As an example, consider two caregivers c and \tilde{c} , where a note \mathbf{x}_i was written by $c_i = \tilde{c}$. We want to estimate what $\mathbf{x}_i(c)$, the note caregiver c would have written, might look like. To this end we will build a model $\tau_c(\cdot)$ that takes data and generates a note in caregiver c 's style. Now suppose caregiver c usually sees patients with high blood pressure and always includes blood pressure values in notes, while \tilde{c} rarely does. A naive model estimating $\hat{\mathbf{x}}_i(c) = \tau_c(\mathbf{x}_i)$ based only on c 's notes may fill in false blood pressure information, conflating that with c 's style. Including vitals data like blood pressure, typically recorded in a patient's health record, can provide additional context for our model. This extra information can assist the model in reasoning about external/background variables, leading to more accurate estimates.

²This claim is shown in Makar et al. [25], appendix A includes a proof for completeness. We set the distribution over C in P_{\perp} as uniform for simplicity, the derivation for non-uniform distributions is analogous.

4.1 Identification of the Counterfactual Distributions

To make effective use of this data, we suggest that the input to the model $\tau_c : \mathcal{X} \times \mathcal{M} \rightarrow \mathcal{X}$ will include a baseline text to be edited and auxiliary data \mathbf{m} . Intuitively, accounting for confounding between the identity of the caregiver C and the text X , with auxiliary data M should result in improved augmentation.

We formalize this intuition using an assumption from causal inference. To identify the counterfactual text distributions using the observed distribution, we assume strong ignorability [73–75]

Assumption 1 (Strong ignorability). *For all $P \in \mathcal{P}$ it holds that $X(c) \perp\!\!\!\perp C \mid M$, and for all values of $\mathbf{m} \in \mathcal{M}$, $P(\mathbf{m}) > 0$.*

Under this assumption, we can rewrite the counterfactual distribution with the observed distribution,

$$P(X(c)) = \int P(X(c) \mid M = \mathbf{m})P(M = \mathbf{m})d\mathbf{m} = \int P(X \mid C = c, M = \mathbf{m})P(M = \mathbf{m})d\mathbf{m}.$$

However, in practice, we do not observe many samples from $P(X \mid C = c, M = \mathbf{m})$, making it a poor approximation for the counterfactual distribution. We address this by using counterfactual data augmentation [37]. Formally, we assume that for all possible counterfactual distributions $c \in [K]$, there exist a function τ_c that maps from the observed distribution $P(X \mid M = \mathbf{m})$ to the target counterfactual distribution $P(X(c) \mid M = m)$.

We approximate the loss under the counterfactual distributions through the empirical loss produced by data augmentation. That is, for a hypothesis $h \in \mathcal{H}$

$$\mathbb{E}_{P(X(c))}[\ell(h(\mathbf{x}), y)] \approx \frac{1}{N} \sum_{i \in [N]} \tau_c(\mathbf{x}_i, \mathbf{m}_i).$$

Note that whenever the text in the training set is already written by caregiver c , i.e. $c_i = c$, we will simply keep the original text \mathbf{x}_i

Evaluation of Augmented Distribution. The right hand-side of the above equation is a Monte-Carlo estimator of the distribution of augmented notes, which averages the distributions $\tau_{*,c}(P_{\text{train}}(X, M))$ over all caregivers $c \in [K]$. The distribution $\tau_{*,c}(P_{\text{train}}(X, M))$ is aimed to follow the style of caregiver c . While the observed samples from one counterfactual distribution may not be sufficient to approximate the whole distribution, they can be used to assess the quality of the counterfactual augmentation algorithm τ_c .

High-quality counterfactual estimation, as measured by small distributional divergence between our estimator and the target distribution, will help in lowering the upper bound on the risk $\mathcal{R}_{P_1}^\ell(h)$ (see lemma 2 in section 4.3). Then to estimate divergences between these two distributions, we may use validation sets from our training data. A sample from $\tau_{*,c}(P_{\text{train}}(X, M))$ is obtained simply by running training data through τ_c , while a sample from $P(X(c))$ can be obtained either by adjusting for M , or we can obtain a sample from $P(X \mid C = c, M = \mathbf{m})$ for each value of \mathbf{m} and compare that to a sample obtained by augmenting validation data where $M = \mathbf{m}$. In both cases two-sample tests can be applied and obtain estimates of divergences between the two distributions. That is of course as long as positivity holds, i.e. the second part of the assumption, as otherwise we will not be able to obtain samples of $P(X \mid C = c, M = \mathbf{m})$ for certain values of \mathbf{m} and c .

We now describe the estimation methods that obtain τ_c . The methods are based on classical causal inference methods, applied to our high-dimensional setting, and relying on the auxiliary data M .

4.2 Methods for Estimation of Counterfactuals

Counterfactual estimation is an established problem in causal effect estimation [74, 76, 77]. Here we adapt identification strategies and estimation procedures in the causal literature to estimate $\mathbf{x}_i(c)$. Our framework for estimating counterfactuals *CATO* (Causal-structure Driven Augmentations for Text OOD Generalization) involves the use of an LLM to model the conditional probability distribution of text. Counterfactuals are formed by matching similar auxiliary data examples or manipulating texts’ vector representations, as described below.

Prompting with matched examples. Our first estimation method in Algorithm 1(B) draws insights from matching [76]. We construct a prompt for an LLM, that given an original text \mathbf{x} and a set of

Algorithm 1 *CATO*

Input: Training set $\{(\mathbf{x}_i, y_i, c_i, \mathbf{m}_i)\}_{i=1}^N$
Hypothesis class \mathcal{H}
Version $\in \{(A), (B)\}$
Optional pre-treatment data $\{(\mathbf{x}_{\text{pre},i})\}_{i=1}^N$
Output: A hypothesis $h_{\text{aug}}(\mathbf{x})$
1: **if** Version = (A) **then**
2: Get $\tau_c(\mathbf{m}, \mathbf{x})$ with preprocess (A)
3: Get $\hat{\mathbf{x}}_i(c) = \tau_c(\mathbf{x}_{i,\text{pre}}, \mathbf{m}_i) \forall i \in [N]$
4: **else**
5: Get $\tau_c(\mathbf{m}, \mathbf{x})$ with preprocess (B)
6: Get $\hat{\mathbf{x}}_i(c) = \tau_c(\mathbf{x}_i, \mathbf{m}_i) \forall i \in [N]$
7: **end if**
8: **return** $h_{\text{aug}} \in \mathcal{H}$ that minimizes $\widehat{\mathcal{R}}_{\text{aug}}^\ell$.

Pre-process *CATO* (A)

Assume: \mathbf{m} includes the label y and pre-treatment attribute c_{pre} , among other auxiliary data. We are given $\{\mathbf{x}_{j,\text{pre}}\}_{j=1}^N$.

- 1: Set $\rho(c_j, \mathbf{m}_j) = \mathbf{x}_j - \mathbf{x}_{j,\text{pre}}$ for $j \in [N]$.
 - 2: **return** $\tau_c(\mathbf{x}, \mathbf{m}) := \mathbf{x}_{\text{pre}} + \rho(c, \mathbf{m})$
-

Pre-process *CATO* (B)

Assume: \mathbf{m} includes the label y among other auxiliary data.

- 1: **return** prompt $\tau_c(\mathbf{x}, \mathbf{m})$ that rewrites \mathbf{x} in the style of matching examples with attribute c , i.e. $\{\mathbf{x}_j : (\mathbf{m}_j, c_j) = (\mathbf{m}, c)\}$.
-

context notes, asks the LLM to rewrite \mathbf{x} in their style. Now given text \mathbf{x} with auxiliary data \mathbf{m} that we wish to estimate with counterfactual value c (i.e. writing style), $\tau_c(\mathbf{x}, \mathbf{m})$ runs this prompt with context notes whose auxiliary data is similar to \mathbf{m} and their attribute value equals the desired c .

Diff-in-diff estimation. The procedure we use for medical note generation relies on additional structure involving panel data (i.e. data collected over time intervals across several individuals). In our case of clinical narratives, a narrative is usually consisted of several notes taken over the course of a patient’s visit and each may be written by a different caregiver. Prediction is made using the release note from the hospital whose embedding consists our features \mathbf{x} . For simplicity let us consider a single note \mathbf{x}_{pre} taken prior to \mathbf{x} . Difference-in-difference [19, 78, 79] estimation of causal effect is based on the parallel-trends, or constant effect assumption that two units i, j with similar pre-treatment conditions would have seen the same effect had they been given the same treatment. In our case, the treatment is an assignment to a certain caregiver. Hence we assume our auxiliary data \mathbf{m} includes c_{pre} , the caregiver assigned pre-treatment.

Assumption 2 (constant effect). *Let $\mathbf{x}_{i,\text{pre}}$ be the pre-treatment features for unit i , and assume \mathbf{m}_i includes the pre-treatment attribute $c_{i,\text{pre}}$. There exists a function $\rho : [K] \times \mathcal{M} \rightarrow \mathcal{X}$ such that $\mathbf{x}_i(c) = \mathbf{x}_{i,\text{pre}} + \rho(c, \mathbf{m}_i)$.*


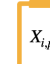



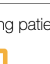


Under this assumption, to calculate $\mathbf{x}_i(c)$ we can use any unit j for which $\mathbf{m}_i = \mathbf{m}_j$ and has $c_j = c$ to estimate $\rho(c, \mathbf{m}_i) = \mathbf{x}_j - \mathbf{x}_{\text{pre},j}$. The resulting estimation procedure is given in algorithm 1(B) and illustrated in section 4.2.

Before empirically evaluating our methods, we discuss alternatives for learning robust classifiers in our setting, and how their properties fair compared to counterfactual augmentation.

4.3 Why Bother with Counterfactual Data Augmentation?

Reasoning about counterfactuals with problem-specific domain knowledge is a considerable challenge, and it is interesting to see whether this has any advantage in learning robust classifiers compared to methods that rely on less stringent assumptions. A simple alternative to approximating counterfactuals involves re-weighting the loss function (see e.g. Makar et al. [25], Shimodaira [80]).

Reweighting baseline. Intuitively, re-weighting samples from the uncorrelated distribution $P(Y, C) = P(Y)P(C)$ by setting for each example i a weight $w_i = P_{\text{train}}(Y = y_i)P_{\text{train}}(C = c_i)/P_{\text{train}}(Y = y_i, C = c_i)$ and

Time \ Patient	$T - 1$ (Progress)		T (Discharge)	
	Caretaker	Note	Caretaker	Note
i				
j				

Panel A: Matching patients using auxiliary data

$$\hat{\mathbf{x}}_i(c_j) = \mathbf{x}_{i,\text{pre}} + (\mathbf{x}_j - \mathbf{x}_{j,\text{pre}})$$

$$\hat{\mathbf{x}}_j(c_i) = \mathbf{x}_{j,\text{pre}} + (\mathbf{x}_i - \mathbf{x}_{i,\text{pre}})$$

Panel B: Generating counterfactual discharge summaries

Figure 2: Generating counterfactual clinical notes for patients using auxiliary data with Algorithm 1(A).

minimizing the weighted empirical risk:

$$\widehat{\mathcal{R}}_{\mathbf{w}}^{\ell}(h) = \frac{1}{m} \sum_{i \in [m]} w_i \ell(h(\mathbf{x}_i), y_i).$$

It can be proved that at the limit of infinite data the method learns a min-max optimal hypothesis, as it also effectively minimizes $\mathcal{R}_{P_{\perp}}^{\ell}$ (see [25]). While augmentations may not seem advantageous for identifying the correct hypothesis, reweighting can demand a larger sample to identify the correct hypothesis, particularly when Y and C are highly correlated.³

Comparing sample complexities. To make this statement precise, we can apply the bounds from Cortes et al. [81] and compare them with an upper bound that we will derive for our method in Lemma 2. To this end, let us consider the exponent of the Rényi divergence as a measure of dependence between Y and C in the training data. The divergence is given by $d_{\alpha, \text{train}}(Y, C) = [\sum_{y \in [L], c \in [K]} P_{\text{train}}^{\alpha}(Y = y, C = c) / P_{\text{train}}^{\alpha-1}(Y = y) P_{\text{train}}^{\alpha-1}(C = c)]^{\frac{1}{\alpha-1}}$, and we may derive the following bound for a hypothesis $h \in \mathcal{H}$ and any $\delta \in [0, 1]$:

$$\mathcal{R}_{P_{\perp}}^{\ell}(h) \leq \widehat{\mathcal{R}}_{\mathbf{w}}^{\ell}(h) + \sqrt{\frac{2d_{2, \text{train}}(Y, C) \cdot \log(1/\delta)}{N}} + \frac{d_{\infty, \text{train}}(Y, C) \cdot \log(1/\delta)}{N}. \quad (2)$$

A complementary lower bound on $\widehat{\mathcal{R}}_{\mathbf{w}}^{\ell}(h)$ can also be derived based on results in Cortes et al. [81]. To compare this with counterfactual augmentations, denote our augmentation model by $\tau : \mathcal{X} \times \mathcal{M} \rightarrow \mathcal{X}^K$, which is some measurable function whose output’s c -th coordinate is the counterfactual estimate w.r.t. caregiver c , i.e. $\hat{\mathbf{x}}(c) = \tau_c(\mathbf{x}, \mathbf{m})$. The following statement quantifies the relation between the accuracy of $\tau(\cdot)$ in approximating counterfactuals and the classification accuracy of a model learned from the augmented data, via minimization of $\widehat{\mathcal{R}}_{\text{aug}}^{\ell}(h)$ in eq. (1).

Lemma 2. *Consider a prediction problem with a spuriously-correlated attribute (definition 1), a measurable function $\tau : \mathcal{X} \times \mathcal{M} \rightarrow \mathcal{X}^K$, and let $d_1(P, Q)$ denote the total variation distance between two distributions P, Q . Further let h^*, h_{aug}^* denote the optimal hypotheses w.r.t $\mathcal{R}_{P_{\perp}}^{\ell_{01}}, \mathcal{R}_{\text{aug}}^{\ell_{01}}$ respectively and let $\lambda_{\text{aug}} = [R_{P_{\perp}}^{\ell_{01}}(h_{\text{aug}}^*) - R_{P_{\perp}}^{\ell_{01}}(h^*)]$. For any hypothesis $h \in \mathcal{H}$, and any $\delta \in (0, 1)$ it holds that with probability at least $1 - \delta$ over the draw of the training set,*

$$\mathcal{R}_{P_{\perp}}^{\ell_{01}}(h) \leq \widehat{\mathcal{R}}_{\text{aug}}^{\ell_{01}}(h) + \sqrt{\frac{\log(1/\delta)}{N}} + K^{-1} \cdot \sum_{c \in [K]} d_1(\tau_{c, *}(P_{\text{train}}(X, M)), P(X(c))) + \lambda_{\text{aug}}.$$

The divergence $d_1(\tau_{c, *}(P_{\text{train}}(X, M)), P(X(c)))$ is a distance between the true distribution over counterfactual instances $P(X(c))$ and our augmented data $\tau_{c, *}(P_{\text{train}}(X, M))$.⁴ Divergences other than total-variation can be used, resulting in tighter bounds, e.g. see Ben-David et al. [82]. As we generate better counterfactuals this divergence decreases, and it can also be shown that h^* and h_{aug}^* coincide. Hence λ_{aug} vanishes and the bound scales with $N^{-\frac{1}{2}}$, resulting in a gain of factor $d_{2, \text{train}}(Y, C)$ over the upper bound on $\widehat{\mathcal{R}}_{\mathbf{w}}^{\ell_{01}}(h)$ in Equation (2). We discuss the details in the appendix, and in Section 5 we show this empirically through simulations.

Takeaways and additional baselines. We emphasize that that the counterfactual datapoints should not be interpreted as “more data” in the sense of i.i.d training examples, they rather embody knowledge about how the causal mechanism that generates features X acts under interventions on the attribute C (as formalized in e.g. [74, 83]). This translates into an improved sample complexity towards risk minimization on P_{\perp} . Counterfactuals are not the only type of causal knowledge that may be leveraged for learning more stable models. Many data dependent penalty terms have been proposed to impose conditional independence constraints drawn from the causal structure of the problem. Theory on these methods usually shows improved OOD performance under infinite data [13, 22, 24, 29]. Our baselines include a method based on the Maximum-Mean Discrepancy (MMD) from Makar et al. [25] who show improved sample complexity under a linear hypothesis class.

³We remark that other works discuss the potential benefits of data augmentation for identification in other problem settings, e.g. [72, Thm. 9] and [17].

⁴The notation $\tau_{c, *}(.)$ denotes the pushforward measure. We note that in our implementation τ_c is data dependent and we ignore this dependence to enable a simple analysis.

5 Experiments

We empirically study the following questions: (1) Can *CATO* enhance OOD performance of downstream classifiers? (2) Does it surpass the combination of reweighting and invariance penalties? (3) Is it more effective than alternative augmentation techniques, thus demonstrating the usefulness of the causal graph? (4) How sensitive is *CATO* to quality of counterfactuals?

These questions seek to establish causally-motivated augmentations as a practical approach for improving OOD performance. We address Q#1,#2 and #3 through our theoretical foundation and across all empirical studies, while Q#4 is explored in the synthetic experiments. Further details about the experimental setup, including data statistics, model hyperparameters, and data splits, can be found in Appendix B. Table 1 provides an overview of the tasks we experiment with.

Input (x)	Label (y)	ID Data	OOD Data	Spurious Feature (c)	auxiliary data (m)
Clinical Narratives	Condition Prediction Note Segmentation Demographic Traits	MIMIC-III	i2b2-2010 partner data i2b2-2006	Caregiver ID	Medications, Lab Results, Vitals
Restaurant Reviews	Restaurant Rating	CEBaB	CeBAB- Spurious	Food-mention	Service, Noise, Ambiance, Food
Synthetic Data	$\{0, 1\}$	Gaussians		$\{0, \dots, 7\}$	-

Table 1: Description of all our tasks and their corresponding experimental setup.

Baselines. We compare *CATO* to several baselines:

- Observational - Baseline model trained on the original data. *PubMed BERT* [84] for *clinical narratives*, logistic regression for the *restaurant reviews* and *synthetic* experiments.⁵
- Reweighting - Baseline model with sample reweighting as in Makar et al. [25].
- MMD - Baseline model with an MMD penalty as in Makar et al. [25], Veitch et al. [29].
- IRM - Baseline model with the IRMv1 penalty as in Arjovsky et al. [13].
- GroupDRO - Baseline model trained with the GroupDRO objective as in Sagawa et al. [85].
- Naive Augmentations - Baseline model on a dataset that also includes augmentations, generated by prompting an LLM to create more examples (without matching or diff-in-diff).
- Conditional Augmentations - Augmentations are generated by matching on auxiliary data and prompting an LLM to create one example in the the style of the other.

The reweighting and MMD approaches are discussed and contrasted to counterfactual augmentation in Section 4. IRM and GroupDRO are the most well-known principled methods for OOD generalization that are used in the literature. The augmentation approaches are compared here to demonstrate the importance of using the causal structure of the data.

5.1 Clinical Narratives

Data. We consider three representative clinical NLP tasks, *clinical condition* prediction, *note segmentation* and *demographic traits* identification⁶, for which we have both ID and OOD data. We utilize several electronic health records (EHR) datasets. We train on MIMIC-III [86], a widely-used medical dataset containing over 2 million notes from 38,597 adult patients, 49,785 hospital admissions, and 3,500 healthcare professionals between 2001 and 2012. MIMIC-III is commonly used in NLP research for clinically-related tasks and for pre-training language models for the medical domain [87]. When available, we use i2b2 2006 and 2010 competitions as our held-out hospital dataset. In the note segmentation task, we use private held-out data.

Generating notes from counterfactual caregivers. To generate augmentations, we select caregivers with multiple patients and notes for more than one patient. For each caregiver-patient pair where both their last progress note and discharge summary were written by that caregiver⁷, we match them to similar patients having the same initial caregiver but a different one for their discharge summary. In matching, we select patients with similar medications and lab results (denoted as patient’s

⁵Appendix B includes results where the Baseline model is also BioBERT, SentenceBERT or GPT3.

⁶See Appendix B for results on the *demographic traits* identification task.

⁷During a patient’s stay, progress notes capture its current state. When leaving the hospital, a discharge summary is written.

auxiliary data m in Table 1). We then generate counterfactual discharge summaries for matched patients using Algorithm 1(A) and train the model using original data and generated counterfactuals.

Figure 3 presents results for *CATO* (A) using language model representations generated using these matched examples. See Appendix B for training details and results for *CATO* (A) with LLM prompts, and Appendix C for synthetic note examples and the prompts used.

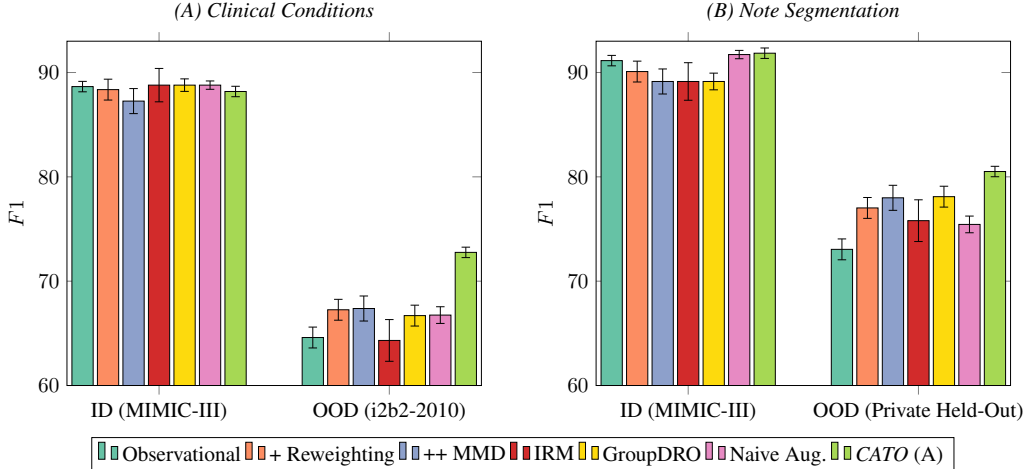


Figure 3: Results ($F1$ averaged across 5 runs) for predicting *clinical conditions* (A) and for *clinical note segmentation* (B) from the text narratives. *CATO* (A) outperforms all baselines on OOD data.

Clinical Condition Prediction. *Clinical condition* prediction is a concept extraction task focused on medical concepts in patient reports [88]. Here we trained *PubMed BERT* models on a subset of MIMIC-III, labelled using the same annotation guidelines as in i2b2-2010, the OOD dataset the models are tested on. As can be seen in the Figure 3(A), in the ID setting only the naive augmentations improve performance slightly. In the OOD setting, all OOD methods help (*reweighting*, *MMD*, *IRM*, *GroupDRO*, *CATO* (A)), but our causally-motivated augmentation approach is substantially better than the alternatives. On average (across 5 runs), *CATO* (A) improves precision above the baseline by more than 7% (absolute), and recall by more than 8%. The naive augmentation approach improves over the vanilla *PubMed BERT* model, but is outperformed by all OOD methods.

Note Segmentation. In this task, models need to recognize sections in free-form clinical notes [89]. Given that section headers vary between hospitals, the models must discern sections based solely on the note content, excluding headers. As can be seen in Figure 3(B), similarly to *clinical condition* prediction, the diff-in-diff approach to augmentations (*CATO* (A)) substantially improved OOD performance, and as expected does not help ID. The naive augmentations are the best performing method ID, but is again outperformed by all other methods OOD.

5.2 Restaurant Reviews

Data. We use the *CeBaB* dataset [49], which consists of short restaurant reviews and ratings from *OpenTable*, including evaluations for food, service, noise, ambiance, and an overall rating. We used the train-exclusive split of the dataset, which contains 1,755 examples. We construct two experimental settings: the original *CeBaB* dataset, and a modified version, denoted as *CeBaB-Spurious*, where there’s a spurious correlation between training and deployment.

To construct *CeBaB-Spurious*, we leverage the availability of both the original and perceived ratings for each review in *CeBaB*. The original rating represents the reviewer’s initial thoughts when writing the review, while the perceived rating indicates whether the review contains information

Method	<i>CeBaB</i>	<i>CeBaB-Spur.</i>
Observational	0.85	0.64
Reweighting	0.84	0.68
Naive Aug.	0.80	0.62
Conditional Aug.	0.84	0.70
<i>CATO</i> (B)	0.84	0.75

Table 2: Accuracy on *CeBaB* and *CeBaB-Spurious*. *CATO* (B) outperforms all baselines when we introduce a spurious correlation.

about various restaurant attributes (e.g., food, service, noise, ambiance) and their associated sentiment. We utilize this unique data structure to capture reviewers’ writing styles. Some reviewers are concise and provide limited descriptions, while others are more descriptive and include more information. To incorporate this variability, we introduce a new attribute called *food-mention* to signify the presence of food-related information in a review. If the perceived food rating is either negative or positive, we assign a value of 1 to the *food-mention* attribute; otherwise, it is set to 0. We subsample the data such that there is a correlation of 0.72 between *food-mention* and the outcome.

Generating reviews with counterfactual food mentions. Following Algorithm 1, we generate counterfactual restaurant reviews conditional on food and overall ratings. We find matched examples for each review, select those with different food-mentions, and prompt an LLM to rewrite them, reflecting how the reviews would appear if the reviewer was more/less concise.

Results. As shown in Table 2, adding counterfactual augmentations leads to better OOD generalization, while naive data augmentation hurts model performance. In line with the sample complexity argument in Section 4, conditional augmentation effectively doesn’t add new data and therefore doesn’t improve model performance.

5.3 Synthetic Data

To test sensitivity of *CATO* to quality of counterfactuals (Q#4), we generate synthetic data for a binary classification problem where $K = 8$ (cardinality of C). We sample $\hat{P}(C | Y)$ to simulate varying degrees of spurious correlations. Then we draw $\mathbf{x} = [\mathbf{x}^*, \mathbf{x}_{\text{spu}}]$ from a Gaussian distribution,

$$\mathbf{x}_i = \begin{bmatrix} \mathbf{x}_i^* \\ \mathbf{x}_{\text{spu},i} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_{y_i} \\ \boldsymbol{\mu}_{c_i} \end{bmatrix}, \begin{bmatrix} \sigma^2 \mathbf{I}_{d^*} & 0 \\ 0 & \sigma_{\text{spu}}^2 \mathbf{I}_{d_c} \end{bmatrix} \right).$$

In this case $\hat{\mathbf{x}}_i(c)$ is obtained by adding $\mu_c - \mu_{c_i}$ to $\mathbf{x}_{\text{spu},i}$. To corrupt our augmentation, we instead add $\xi_i(\mu_c - \mu_{c_i})$ where ξ_i is drawn from a truncated Gaussian centered at $\lambda \in (0, 1)$. We train models with a fixed sample size (in the appendix we also examine varying sample sizes and additional types of corruption) and evaluate the trained models’ accuracy on P_1 to examine the interplay between spurious correlation strength (measured by mutual information $I(Y; C)$), and counterfactual augmentation quality. As can be seen in Figure 4, corruptions degrade performance under stronger spurious correlations, though a strong corruption is required for reweighting to become preferable.

6 Discussion

In this work, we have presented a data augmentation approach based on the causal structure of auxiliary data for improving OOD generalization, specifically focusing on text classification tasks. However, our approach is not without limitations. The validity of our assumptions, the specification of the causal graph and the quality of the counterfactual approximation all present challenges to address in future work. Further, our results suggest that performing data augmentation in an unprincipled manner can also hurt model performance. Utilizing additional techniques for OOD generalization, learning the causal structure directly from the data, and improving quality and reliability of the counterfactual approximation process can help mitigate these concerns. Overall, we believe that causally-motivated data augmentation methods like ours can help address challenges in developing robust and reliable machine learning systems, particularly in safety-critical applications.

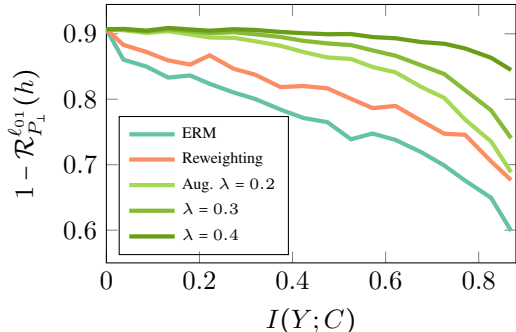


Figure 4: OOD accuracy ($1 - \mathcal{R}_{P_1}^{l_{01}}(h)$) and Y, C correlation strength ($I(Y; C)$). Lower values of λ correspond to stronger corruptions of the augmentations. Even with substantial corruption ($\lambda = 0.2$) and strong correlation, augmentations outperform baselines.

References

- [1] John R Zech, Marcus A Badgeley, Manway Liu, Anthony B Costa, Joseph J Titano, and Eric Karl Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS medicine*, 15(11):e1002683, 2018.
- [2] Alex J DeGrave, Joseph D Janizek, and Su-In Lee. Ai for radiographic covid-19 detection selects shortcuts over signal. *Nature Machine Intelligence*, 3(7):610–619, 2021.
- [3] R Thomas McCoy, Ellie Pavlick, and Tal Linzen. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. *arXiv preprint arXiv:1902.01007*, 2019.
- [4] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1721–1730, 2015.
- [5] Joaquin Quinonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. Mit Press, 2008.
- [6] Adarsh Subbaswamy, Peter Schulam, and Suchi Saria. Preventing failures due to dataset shift: Learning predictive models that transport. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3118–3127. PMLR, 2019.
- [7] Samuel G Finlayson, Adarsh Subbaswamy, Karandeep Singh, John Bowers, Annabel Kupke, Jonathan Zittrain, Isaac S Kohane, and Suchi Saria. The clinician and dataset shift in artificial intelligence. *New England Journal of Medicine*, 385(3):283–286, 2021.
- [8] Peter Spyns. Natural language processing in medicine: an overview. *Methods of information in medicine*, 35(04/05):285–301, 1996.
- [9] Li Zhou and George Hripcsak. Temporal reasoning with medical data—a review with emphasis on medical natural language processing. *Journal of biomedical informatics*, 40(2):183–202, 2007.
- [10] Stephen Wu, Kirk Roberts, Surabhi Datta, Jingcheng Du, Zongcheng Ji, Yuqi Si, Sarvesh Soni, Qiong Wang, Qiang Wei, Yang Xiang, et al. Deep learning in clinical natural language processing: a methodical review. *Journal of the American Medical Informatics Association*, 27(3):457–470, 2020.
- [11] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2016.
- [12] Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M Mooij. Domain adaptation by using causal inference to predict invariant conditional distributions. In *Neural Information Processing Systems (NeurIPS)*, pages 10869–10879, 2018.
- [13] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- [14] Adarsh Subbaswamy, Bryant Chen, and Suchi Saria. A unifying causal framework for analyzing dataset shift-stable learning algorithms. *Journal of Causal Inference*, 10(1):64–89, 2022.
- [15] Alexander Robey, George J. Pappas, and Hamed Hassani. Model-based domain generalization. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Neural Information Processing Systems (NeurIPS)*, 2021. URL <https://openreview.net/forum?id=J0xB9h40A-1>.
- [16] Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. Improving out-of-distribution robustness via selective augmentation. In *International Conference on Machine Learning*, pages 25407–25437. PMLR, 2022.

- [17] Irena Gao, Shiori Sagawa, Pang Wei Koh, Tatsunori Hashimoto, and Percy Liang. Out-of-domain robustness via targeted augmentations. *arXiv preprint arXiv:2302.11861*, 2023.
- [18] Murat Kocaoglu, Christopher Snyder, Alexandros G. Dimakis, and Sriram Vishwanath. CausalGAN: Learning causal implicit generative models with adversarial training. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=BJE-4xWOW>.
- [19] Alberto Abadie. Semiparametric difference-in-differences estimators. *The review of economic studies*, 72(1):1–19, 2005.
- [20] Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018.
- [21] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 624–639, 2018.
- [22] Yoav Wald, Amir Feder, Daniel Greenfeld, and Uri Shalit. On calibration and out-of-domain generalization. *Neural Information Processing Systems (NeurIPS)*, 34:2215–2227, 2021.
- [23] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, pages 5815–5826. PMLR, 2021.
- [24] Aahlad Manas Puli, Lily H Zhang, Eric Karl Oermann, and Rajesh Ranganath. Out-of-distribution generalization in the presence of nuisance-induced spurious correlations. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=12RoR2o32T>.
- [25] Maggie Makar, Ben Packer, Dan Moldovan, Davis Blalock, Yoni Halpern, and Alexander D’Amour. Causally motivated shortcut removal using auxiliary labels. In *International Conference on Artificial Intelligence and Statistics*, pages 739–766. PMLR, 2022.
- [26] Yibo Jiang and Victor Veitch. Invariant and transportable representations for anti-causal domain shifts. *arXiv preprint arXiv:2207.01603*, 2022.
- [27] Claudia Shi, Victor Veitch, and David M Blei. Invariant representation learning for treatment effect estimation. In *Uncertainty in Artificial Intelligence*, pages 1546–1555. PMLR, 2021.
- [28] Mingzhang Yin, Yixin Wang, and David M Blei. Optimization-based causal estimation from heterogeneous environments. *arXiv preprint arXiv:2109.11990*, 2021.
- [29] Victor Veitch, Alexander D’Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual invariance to spurious correlations in text classification. *Neural Information Processing Systems (NeurIPS)*, 34:16196–16208, 2021.
- [30] Yana Dranker, He He, and Yonatan Belinkov. Irm—when it works and when it doesn’t: A test case of natural language inference. *Advances in Neural Information Processing Systems*, 34:18212–18224, 2021.
- [31] Amir Feder, Katherine A Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E Roberts, et al. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond. *Transactions of the Association for Computational Linguistics*, 10:1138–1158, 2022.
- [32] Amir Feder, Guy Horowitz, Yoav Wald, Roi Reichart, and Nir Rosenfeld. In the eye of the beholder: Robust prediction with causal user modeling. In *Neural Information Processing Systems (NeurIPS)*, 2022.
- [33] Pritish Kamath, Akilesh Tangella, Danica Sutherland, and Nathan Srebro. Does invariant risk minimization capture invariance? In *International Conference on Artificial Intelligence and Statistics*, pages 4069–4077. PMLR, 2021.

- [34] Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. *arXiv preprint arXiv:2010.05761*, 2020.
- [35] Ruocheng Guo, Pengchuan Zhang, Hao Liu, and Emre Kiciman. Out-of-distribution prediction with invariant risk minimization: The limitation and an effective fix. *arXiv preprint arXiv:2101.07732*, 2021.
- [36] Yoav Wald, Gal Yona, Uri Shalit, and Yair Carmon. Malign overfitting: Interpolation can provably preclude invariance. *arXiv preprint arXiv:2211.15724*, 2022.
- [37] Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. Learning the difference that makes a difference with counterfactually-augmented data. *arXiv preprint arXiv:1909.12434*, 2019.
- [38] Divyansh Kaushik, Amrith Setlur, Eduard Hovy, and Zachary C Lipton. Explaining the efficacy of counterfactually-augmented data. *arXiv preprint arXiv:2010.02114*, 2020.
- [39] Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H Chi, and Alex Beutel. Counterfactual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 219–226, 2019.
- [40] Rohan Jha, Charles Lovering, and Ellie Pavlick. Does data augmentation improve generalization in nlp? *arXiv preprint arXiv:2004.15012*, 2020.
- [41] Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. Evaluating models’ local decision boundaries via contrast sets. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.117. URL <https://aclanthology.org/2020.findings-emnlp.117>.
- [42] Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. FOIL it! find one mismatch between image and language caption. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 255–265, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1024. URL <https://aclanthology.org/P17-1024>.
- [43] Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. Causalm: Causal model explanation through counterfactual language models. *Computational Linguistics*, 47(2):333–386, 2021.
- [44] Ran Zmigrod, Sabrina J. Mielke, Hanna Wallach, and Ryan Cotterell. Counterfactual data augmentation for mitigating gender stereotypes in languages with rich morphology. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1651–1661, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1161. URL <https://aclanthology.org/P19-1161>.
- [45] Parker Riley, Noah Constant, Mandy Guo, Girish Kumar, David Uthus, and Zarana Parekh. Textsettr: Label-free text style extraction and tunable targeted restyling. *arXiv preprint arXiv:2010.03802*, 2020.
- [46] Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel S Weld. Polyjuice: Automated, general-purpose counterfactual generation. *arXiv preprint arXiv:2101.00288*, 2021.
- [47] Chengzhi Mao, Augustine Cha, Amogh Gupta, Hao Wang, Junfeng Yang, and Carl Vondrick. Generative interventions for causal learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3947–3956, 2021.
- [48] Daniel Rosenberg, Itai Gat, Amir Feder, and Roi Reichart. Are vqa systems rad? measuring robustness to augmented data with focused interventions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 61–70, 2021.

- [49] Eldar D Abraham, Karel D’Oosterlinck, Amir Feder, Yair Gat, Atticus Geiger, Christopher Potts, Roi Reichart, and Zhengxuan Wu. CEBaB: Estimating the causal effects of real-world concepts on NLP model behavior. *Neural Information Processing Systems (NeurIPS)*, 35: 17582–17596, 2022.
- [50] Zhengxuan Wu, Karel D’Oosterlinck, Atticus Geiger, Amir Zur, and Christopher Potts. Causal proxy models for concept-based model explanations. In *International Conference on Machine Learning*, pages 37313–37334. PMLR, 2023.
- [51] Nitish Joshi and He He. An investigation of the (in) effectiveness of counterfactually augmented data. In *60th Annual Meeting of the Association for Computational Linguistics, ACL 2022*, pages 3668–3681. Association for Computational Linguistics (ACL), 2022.
- [52] Maria Antoniak and David Mimno. Bad seeds: Evaluating lexical methods for bias measurement. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1889–1904, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.148. URL <https://aclanthology.org/2021.acl-long.148>.
- [53] Xiaoling Zhou and Ou Wu. Implicit counterfactual data augmentation for deep neural networks. *arXiv preprint arXiv:2304.13431*, 2023.
- [54] Nitay Calderon, Eyal Ben-David, Amir Feder, and Roi Reichart. DoCoGen: Domain Counterfactual Generation for Low Resource Domain Adaptation. In *Proceedings of the 60th Annual Meeting of the Association of Computational Linguistics (ACL)*, 2022.
- [55] Kory Kreimeyer, Matthew Foster, Abhishek Pandey, Nina Arya, Gwendolyn Halford, Sandra F Jones, Richard Forshee, Mark Walderhaug, and Taxiarchis Botsis. Natural language processing systems for capturing and standardizing unstructured clinical information: a systematic review. *Journal of biomedical informatics*, 73:14–29, 2017.
- [56] Özlem Uzuner. Recognizing obesity and comorbidities in sparse data. *Journal of the American Medical Informatics Association*, 16(4):561–570, 2009.
- [57] Guergana K Savova, James J Masanz, Philip V Ogren, Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-Schuler, and Christopher G Chute. Mayo clinical text analysis and knowledge extraction system (ctakes): architecture, component evaluation and applications. *Journal of the American Medical Informatics Association*, 17(5):507–513, 2010.
- [58] Peter B Jensen, Lars J Jensen, and Søren Brunak. Mining electronic health records: towards better research applications and clinical care. *Nature Reviews Genetics*, 13(6):395–405, 2012.
- [59] Elizabeth Ford, John A Carroll, Helen E Smith, Donia Scott, and Jackie A Cassell. Extracting information from the text of electronic medical records to improve case detection: a systematic review. *Journal of the American Medical Informatics Association*, 23(5):1007–1015, 2016.
- [60] Henghui Zhu, Ioannis Ch Paschalidis, and Amir Tahmasebi. Clinical concept extraction with contextual word embedding. *arXiv preprint arXiv:1810.10566*, 2018.
- [61] Yifan Peng, Shankai Yan, and Zhiyong Lu. Transfer learning in biomedical natural language processing: an evaluation of bert and elmo on ten benchmarking datasets. *arXiv preprint arXiv:1906.05474*, 2019.
- [62] Vikas Yadav and Steven Bethard. A survey on recent advances in named entity recognition from deep learning models. *arXiv preprint arXiv:1910.11470*, 2019.
- [63] Yuqi Si, Jingqi Wang, Hua Xu, and Kirk Roberts. Enhancing clinical concept extraction with contextual embeddings. *Journal of the American Medical Informatics Association*, 26(11): 1297–1304, 2019.
- [64] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.

- [65] Dmitri Roussinov, Andrew Conkie, Andrew Patterson, and Christopher Sainsbury. Predicting clinical events based on raw text: from bag-of-words to attention-based transformers. *Frontiers in Digital Health*, 3:214, 2022.
- [66] Tom M Seinen, Egill A Fridgeirsson, Solomon Ioannou, Daniel Jeannetot, Luis H John, Jan A Kors, Aniek F Markus, Victor Pera, Alexandros Rekkas, Ross D Williams, et al. Use of unstructured text in prognostic clinical prediction models: a systematic review. *Journal of the American Medical Informatics Association*, 29(7):1292–1302, 2022.
- [67] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022.
- [68] John W Ayers, Adam Poliak, Mark Dredze, Eric C Leas, Zechariah Zhu, Jessica B Kelley, Dennis J Faix, Aaron M Goodman, Christopher A Longhurst, Michael Hogarth, et al. Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum. *JAMA Internal Medicine*, 2023.
- [69] Amir Feder, Itay Laish, Shashank Agarwal, Uri Lerner, Avel Atias, Cathy Cheung, Peter Clardy, Alon Peled-Cohen, Rachana Fellingner, Hengrui Liu, et al. Building a clinically-focused problem list from medical notes. In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 60–68, 2022.
- [70] Fan Zhang, Itay Laish, Ayelet Benjamini, and Amir Feder. Section classification in clinical notes with multi-task transformers. In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 54–59, 2022.
- [71] Amir Feder, Danny Vainstein, Roni Rosenfeld, Tzvika Hartman, Avinatan Hassidim, and Yossi Matias. Active deep learning to detect demographic traits in free-form clinical notes. *Journal of Biomedical Informatics*, 107:103436, 2020.
- [72] Zihao Wang and Victor Veitch. A unified causal view of domain invariant representation learning. *arXiv preprint arXiv:2208.06987*, 2022.
- [73] Guido W Imbens and Jeffrey M Wooldridge. Recent developments in the econometrics of program evaluation. *Journal of economic literature*, 47(1):5–86, 2009.
- [74] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [75] Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *International conference on machine learning*, pages 3076–3085. PMLR, 2017.
- [76] Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- [77] Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- [78] David Card and Alan B Krueger. Minimum wages and employment: A case study of the fast food industry in new jersey and pennsylvania, 1993.
- [79] Joshua D Angrist and Jörn-Steffen Pischke. *Mostly harmless econometrics: An empiricist’s companion*. Princeton university press, 2009.
- [80] Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- [81] Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance weighting. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Neural Information Processing Systems (NeurIPS)*, volume 23. Curran Associates, Inc., 2010.
- [82] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.

- [83] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.
- [84] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021.
- [85] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2019.
- [86] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [87] Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*, 2019.
- [88] Özlem Uzuner, Brett R South, Shuying Shen, and Scott L DuVall. 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text. *Journal of the American Medical Informatics Association*, 18(5):552–556, 2011.
- [89] Alexandra Pomares-Quimbaya, Markus Kreuzthaler, and Stefan Schulz. Current approaches to identify sections within clinical narratives from electronic health records: a systematic review. *BMC medical research methodology*, 19:1–20, 2019.
- [90] Terence Tao. *An introduction to measure theory*, volume 126. American Mathematical Soc., 2011.
- [91] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [92] Koby Crammer, Michael Kearns, and Jennifer Wortman. Learning from multiple sources. *Journal of Machine Learning Research*, 9(8), 2008.
- [93] Aahlad Puli, Nitish Joshi, He He, and Rajesh Ranganath. Nuisances via negativa: Adjusting for spurious correlations via data augmentation. *arXiv preprint arXiv:2210.01302*, 2022.
- [94] Aahlad Puli, Lily Zhang, Yoav Wald, and Rajesh Ranganath. Don’t blame dataset shift! shortcut learning due to gradients and cross entropy. *arXiv preprint arXiv:2308.12553*, 2023.
- [95] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [96] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [97] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- [98] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.
- [99] OpenAI. Gpt-4 technical report, 2023.

Appendix

A Proofs of Formal Claims

Notation. We will use random variables C, Y, M, X with images $[K], \mathcal{Y} = [L], \mathcal{M}, \mathcal{X}$ respectively in our probabilistic causal models. For a function $\tau_c : \mathcal{X} \times \mathcal{M} \rightarrow \mathcal{X}$, and measure P over sets in $\mathcal{X} \times \mathcal{M}$, we denote by $\tau_{c,*}P(X, M)$ the pushforward measure [90, §1.4]. $\tau_c(\cdot)$ will be used to refer to the c -th coordinate of the output of a function $\tau : \mathcal{X} \times \mathcal{M} \rightarrow \mathcal{X}^K$. The notation \mathcal{H} will be used for hypothesis classes where $h : \mathcal{X} \rightarrow \mathcal{Y}$ for any $h \in \mathcal{H}$. The 0-1 loss $\ell_{01} : \mathcal{Y} \times \mathcal{Y} \rightarrow \{0, 1\}$ is given by $\ell_{01}(\hat{y}, y) = 1_{\hat{y} \neq y}$. For a node V in a causal graph we will use $pa(V)$ for its causal parents.

For completeness we rewrite the definition of our data generating process from the main paper, this time adding the auxiliary data M into our model.

Definition 1. Consider a probabilistic causal model with endogenous random variables X, X^*, Y, C, M taking on values in $\mathcal{X}, \mathcal{X}^*, [L], [K], \mathcal{M}$ and exogenous independent random variables [83] $N_X, N_{X^*}, N_Y, N_C, N_M$, where the induced graph is a DAG that satisfies the following,

- Y is d -separated from X by X^*, C, M and also by X^*, C .
- Y, X^* are not descendants of C .

An anti-causal prediction problem with a spuriously-correlated attribute is a set of distributions \mathcal{P} obtained by all interventions on C that replaces the distribution of exogenous noise N_C , mechanism $f_C(pa(C), N_C)$ with another mechanism (i.e. a measurable function $\tilde{f}(pa(C), N_C)$), or sets a fixed value (i.e. $do(C = c)$). Under the settings of this problem, a learner is provided with a set $\{(\mathbf{x}_i, y_i, c_i)\}_{i=1}^N$ sampled i.i.d from $P_{\text{train}} \in \mathcal{P}$.

We denote by $P_{\perp} \in \mathcal{P}$ the distribution obtained by intervening on C and setting it to a uniform distribution, i.e. $P_{\perp}(X, X^*, Y, C, M) = K^{-1} \sum_{c \in [K]} P(Y, X, X^*, M | do(C = c))$. Note that the problem described by fig. 1 and definition 1 of the main paper is a special case of this setting where M is discarded, and P_{\perp} coincides with setting $\tilde{P}(C | Y)$ to a uniform distribution.

Recall our assumption about perfect recovery of X^* .

Assumption 3. For an anti-causal prediction problem with a spuriously correlated attribute, we assume that $X^* = e(X)$ a.e. for some $e : \mathcal{X} \rightarrow \mathcal{X}^*$.

Under these conditions $h(\mathbf{x}) = \arg \max_{y \in [L]} P_{\perp}(Y = y | X = \mathbf{x})$ is an optimal risk-invariant predictor as described below.

Lemma 1. For the prediction problem in definition 1, the Bayes optimal classifier under the unconfounded distribution $P_{\perp} \in \mathcal{P}$ where C is uniformly distributed and independent of Y is $h^*(\mathbf{x}) = \arg \max_{y \in [K]} P_{\perp}(Y = y | X^* = e(\mathbf{x}))$. It is a minimizer of $\min_{h: \mathcal{X} \rightarrow [L]} \max_{P \in \mathcal{P}} \mathcal{R}_P^{\ell_{01}}(h)$ and $\mathcal{R}_P^{\ell_{01}}(h^*) = \mathcal{R}_{P_{\perp}}^{\ell_{01}}(h^*)$ for all $P \in \mathcal{P}$.

Proof. Assume $P_{\text{train}} \in \mathcal{P}$ is the distribution from which our training data is obtained. We will show that any hypothesis satisfying $h(X) = g \circ e(X)$ for some $g : \mathcal{X}^* \rightarrow \mathcal{Y}$ (i.e. that only depends on X^*) achieves the same risk over all $P \in \mathcal{P}$. To this end note that for such a hypothesis we have,

$$\begin{aligned}
 R_{P_{\text{train}}}^{\ell_{01}}(h) &= \int \ell_{01}(h(X), Y) P_{\text{train}}(X | Y, C, X^*, M) P_{\text{train}}(Y, C, X^*, M) dX^* dX dY dC dM \\
 &= \int \ell_{01}(g \circ e(X), Y) P_{\text{train}}(X | C, X^*, M) P_{\text{train}}(Y, C, X^*, M) dX^* dX dY dC dM \\
 &= \int \ell_{01}(g(X^*), Y) P_{\text{train}}(X | C, X^*, M) P_{\text{train}}(Y, C, X^*, M) dX^* dX dY dC dM \\
 &= \int \ell_{01}(g(X^*), Y) P_{\text{train}}(X^*, Y) dX^* dY \\
 &= \int \ell_{01}(g(X^*), Y) P(X^*, Y) dX^* dY.
 \end{aligned}$$

The first line writes down the expected risk explicitly, the second removes conditioning on Y in the distribution on X since we assumed Y is d -separated from X by C, X^*, M . In the third line we

make it explicit that h depends on X^* alone, then we integrate out X, C, M . On the last line we remove the subscript train to denote that this distribution is fixed across $P \in \mathcal{P}$ as we assumed that X^*, Y are non-descendants of C (and members of \mathcal{P} are obtained by interventions on C). Now for any $\tilde{P} \in \mathcal{P}$ we may repeat this derivation for $R_{\tilde{P}}^{\ell_{01}}(h)$ and we will obtain the same term (since $P(X^*, Y)$ are fixed regardless of the intervention applied in P , as we just argued), and we may conclude $R_{P_{\text{train}}}^{\ell_{01}}(h) = R_{\tilde{P}}^{\ell_{01}}(h)$.

Next to show that the Bayes optimal classifier over P_{\perp} is the min-max optimal classifier w.r.t \mathcal{P} , consider the interventional distribution where C is set to some fixed value $c \in [K]$, i.e. $P(X, X^*, Y \mid do(C = c))$. Under the graph we obtain from this intervention, Y is d -separated from X given X^* . Hence,

$$\begin{aligned} P(Y \mid X = \mathbf{x}, do(C = c)) &= \int_{X^*} P(Y \mid X^*, X = \mathbf{x}, do(C = c))P(X^* \mid X = \mathbf{x}, do(C = c))dX^* \\ &= P(Y \mid X^* = e(\mathbf{x}), X = \mathbf{x}, do(C = c)) \\ &= P(Y \mid X^* = e(\mathbf{x}), do(C = c)), \end{aligned}$$

where the first equality holds since $X^* = e(X)$ and the second from d -separation. Hence the Bayes optimal classifier under $P(Y, X \mid do(C = c))$ is $h^*(\mathbf{x}) = g \circ e(\mathbf{x}) = \arg \max_{y \in [L]} P(Y = y \mid e(\mathbf{x}), do(C = c))$. As argued earlier, since Y, X^* are non-descendants of C , it holds that $P(Y \mid e(X), do(C = c))$ is fixed across all $c \in [K]$. Hence $h^*(\mathbf{x})$ is the Bayes optimal classifier for all such interventional distributions and also for $P_{\perp}(X, Y) = \frac{1}{K} \sum_{c \in [K]} P(X, Y \mid do(C = c))$, and from our earlier discussion it is risk-invariant, i.e. $R_{P_{\perp}}^{\ell_{01}}(h^*) = R_P^{\ell_{01}}(h^*)$ for all $P \in \mathcal{P}$, which also means $\max_{P \in \mathcal{P}} R_P^{\ell_{01}}(h^*) = R_{P_{\perp}}^{\ell_{01}}(h^*)$. It is the min-max optimal classifier w.r.t \mathcal{P} since any $h \neq h^*$ will have $\max_{P \in \mathcal{P}} R_P^{\ell_{01}}(h) \geq R_{P_{\perp}}^{\ell_{01}}(h) \geq R_{P_{\perp}}^{\ell_{01}}(h^*)$. \square

Next we turn to prove a bound on sample complexity of counterfactual data augmentations.

Lemma 2. *Consider an anti-causal prediction problem with a spuriously-correlated attribute (definition 1), a measurable function $\tau : \mathcal{X} \times \mathcal{M} \rightarrow \mathcal{X}^K$, and let $d_1(P, Q)$ denote the total variation distance between two distributions P, Q . Further let h^* denote the optimal hypothesis w.r.t $\mathcal{R}_{P_{\perp}}^{\ell_{01}}$ and let $\lambda_{\text{aug}} = [R_{\text{aug}}^{\ell_{01}}(h^*) + R_{P_{\perp}}^{\ell_{01}}(h^*)]$. For any hypothesis $h \in \mathcal{H}$, and any $\delta \in (0.5, 1)$ it holds that with probability at least $1 - \delta$ over the draw of the training set,*

$$\mathcal{R}_{P_{\perp}}^{\ell_{01}}(h) \leq \widehat{\mathcal{R}}_{\text{aug}}^{\ell_{01}}(h) + \sqrt{\frac{\log(1/\delta)}{N}} + K^{-1} \cdot \sum_{c \in [K]} d_1(\tau_{c,*}(P_{\text{train}}(X, M)), P(X(c))) + \lambda_{\text{aug}}.$$

Proof. Our first step is to show that for any hypothesis $h \in \mathcal{H}$, if our augmentation process is exact in the sense that $\tau_c(X, M) = X(c)$ a.e., then the expected risk (i.e. risk taken over an infinitely large sample) on the augmented data coincides with that over the unconfounded distribution $P_{\perp}(X, Y) = P_{\text{unif}}(C)P(X, Y \mid do(C))$.

$$\begin{aligned} \mathcal{R}_{\text{aug}}^{\ell_{01}}(h) &= \mathbb{E}_{P_{\text{train}}(C, Y, M, X)} \left[K^{-1} \sum_{c \in [K]} \ell_{01}(h(\tau_c(X, M)), Y) \right] \\ &= K^{-1} \sum_{c \in [K]} \mathbb{E}_{P_{\text{train}}(C, Y, M, X)} [\ell_{01}(h(X(c)), Y)] \\ &= K^{-1} \sum_{c \in [K]} \mathbb{E}_{P_{\text{train}}(C, Y, X)} [\ell_{01}(h(X(c)), Y(c))] \\ &= K^{-1} \sum_{c \in [K]} \mathbb{E}_{P(Y, X \mid do(C=c))} [\ell_{01}(h(X), Y)] \\ &= \mathcal{R}_{P_{\perp}}^{\ell_{01}}(h). \end{aligned} \tag{3}$$

To bound $\mathcal{R}_{\text{aug}}^{\ell_{01}}(h) - \widehat{\mathcal{R}}_{\text{aug}}^{\ell_{01}}(h)$ we note that $\{\mathbf{x}_i, y_i, \mathbf{m}_i\}_{i=1}^N$ are *i.i.d* samples from a joint distribution, where we may consider the loss on each example as $K^{-1} \sum_{c \in [K]} \ell_{01}(h(\tau_c(\mathbf{x}_i, \mathbf{m}_i), y_i))$, then by

standard results using the Hoeffding inequality, e.g. Mohri et al. [91, Corollary 2.11], we get that for $\delta \in (0.5, 1)$,

$$\mathcal{R}_{\text{aug}}^{\ell_{01}}(h) \leq \widehat{\mathcal{R}}_{\text{aug}}^{\ell_{01}}(h) + \sqrt{\frac{\log(1/\delta)}{N}}. \quad (4)$$

Finally, to obtain our result consider any $c \in [C]$. Denote

$$\begin{aligned} \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h) &:= \mathbb{E}_{P_{\text{train}}(Y,M,X)}[\ell_{01}(h(\tau_c(X,M))Y)], \\ \mathcal{R}_{P_1,c}^{\ell_{01}}(h) &:= \mathbb{E}_{P(Y,X|d_{\mathcal{O}}(C=c))}[\ell_{01}(h(X),Y)], \end{aligned}$$

and for h^* denote $\mathcal{R}_{\text{aug},c}^{\ell_{01}}(h, h^*) := \mathbb{E}_{P_{\text{train}}(Y,M,X)}[\ell_{01}(h(\tau_c(X,M)), h^*(\tau_c(X,M)))]$ and respectively for $\mathcal{R}_{P_1,c}^{\ell_{01}}(h, h^*)$. The rest of our derivation is along the lines of Ben-David et al. [82, Theorem 2]. We use the distance

$$d_{\mathcal{H}\Delta\mathcal{H}}(\tau_{c,*}P_{\text{train}}(X, M), P(X(c))) = 2 \sup_{g \in \mathcal{H}\Delta\mathcal{H}} |P_{\text{train}}(g(\tau_c(X, M)) = 1) - P(g(X(c)) = 1)|,$$

where $\mathcal{H}\Delta\mathcal{H} = \{g(\mathbf{x}) = 1_{h(\mathbf{x}) \neq h'(\mathbf{x})} \mid h, h' \in \mathcal{H}\}$ is a set of binary hypotheses, i.e. functions that mark disagreements between hypotheses in \mathcal{H} . It is easy to see that $d_{\mathcal{H}\Delta\mathcal{H}}$ lower bounds d_1 which takes the supremum w.r.t all measurable subsets for the two measures, since the sets of inputs where $h(\mathbf{x}) = 1$ are contained in those subsets. Also from [82, Lemma 3] we have that for any hypotheses $h, h' \in \mathcal{H}$ it holds that

$$|\mathcal{R}_{\text{aug},c}^{\ell_{01}}(h, h') - \mathcal{R}_{P_1,c}^{\ell_{01}}(h, h')| \leq \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\tau_{c,*}P_{\text{train}}(X, M), P(X(c)))$$

Then following the proof in Ben-David et al. [82, Theorem 2], where the first and third inequalities will rely on the triangle inequality for classification errors [92], we may get:

$$\begin{aligned} \mathcal{R}_{P_1,c}^{\ell_{01}}(h) &\leq \mathcal{R}_{P_1,c}^{\ell_{01}}(h^*) + \mathcal{R}_{P_1,c}^{\ell_{01}}(h, h^*) \\ &\leq \mathcal{R}_{P_1,c}^{\ell_{01}}(h^*) + \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h, h^*) + [\mathcal{R}_{P_1,c}^{\ell_{01}}(h, h^*) - \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h, h^*)] \\ &\leq \mathcal{R}_{P_1,c}^{\ell_{01}}(h^*) + \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h, h^*) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\tau_{c,*}P_{\text{train}}(X, M), P(X(c))) \\ &\leq \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h) + \mathcal{R}_{P_1,c}^{\ell_{01}}(h^*) + \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h^*) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\tau_{c,*}P_{\text{train}}(X, M), P(X(c))) \\ &= \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h) + \mathcal{R}_{P_1,c}^{\ell_{01}}(h^*) + \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h^*) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\tau_{c,*}P_{\text{train}}(X, M), P(X(c))) \end{aligned}$$

Finally, we note that $\mathcal{R}_{P_1}^{\ell_{01}}(h) = K^{-1} \sum_{c \in [K]} \mathcal{R}_{P_1,c}^{\ell_{01}}(h)$ and similarly we have that $\mathcal{R}_{\text{aug}}^{\ell_{01}}(h) = K^{-1} \sum_{c \in [K]} \mathcal{R}_{\text{aug},c}^{\ell_{01}}(h)$, hence applying the above inequality for all $c \in [K]$ and averaging we get:

$$\begin{aligned} \mathcal{R}_{P_1}^{\ell_{01}}(h) &\leq \mathcal{R}_{\text{aug}}^{\ell_{01}}(h) + \frac{1}{2} K^{-1} \sum_{c \in [K]} d_{\mathcal{H}\Delta\mathcal{H}}(\tau_{c,*}P_{\text{train}}(X, M), P(X(c))) + \lambda_{\text{aug}} \\ &\leq \mathcal{R}_{\text{aug}}^{\ell_{01}}(h) + K^{-1} \sum_{c \in [K]} d_1(\tau_{c,*}P_{\text{train}}(X, M), P(X(c))) + \lambda_{\text{aug}}. \end{aligned}$$

Combining with eq. (4) we get the desired result. \square

A.1 Additional Causal Structures Where our Approach may be Used

The problem setting we analyze in this work (see definition 1) captures a few interesting problems, mainly described as shortcut learning in the literature [25, 93, 94]. However counterfactual data augmentation, and subsequently our approach of using auxiliary data to perform it, are applicable to additional problem settings. Wang and Veitch [72] formalize domain-invariant learning under many data generating processes they refer to as Causally Invariant with Spurious Associations (CISA), where Z (in our setting the caregiver C) is called the spurious factor of variation. These settings include a variety of causal and anti-causal prediction problems, and they assume that there exists some part of the input X , referred to as $X_{\perp Z}^{\frac{1}{2}}$, that holds all the information in X that is not caused by Z . Whenever it holds that $Y \perp\!\!\!\perp X \mid X_{\perp Z}^{\frac{1}{2}}, Z$ the association between Z and Y is called ‘‘purely

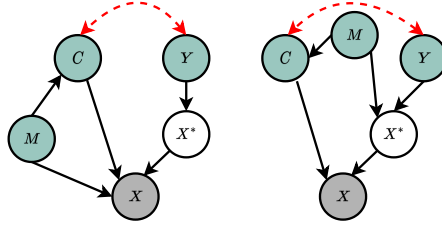


Figure 5: Possible causal structures that involve the auxiliary data M , where unobserved M corresponds to unobserved confounding between X and C .

spurious" and Thm. 9 in Wang and Veitch [72] states that for all such problems counterfactual data augmentation learns the optimal invariant predictor over the training distribution. Hence in all such settings, improving counterfactual data augmentation with *CATO* can be beneficial towards OOD generalization. We refer the interested reader to [72] for further details on CISA problems and their properties.

We further note that in our work we excluded the auxiliary data M from the causal model as we are agnostic to its specific causal relation with other factors in the data, so long as it satisfies ?? 1 of strong ignorability. fig. 5 depicts two potential structures that may adhere to this assumption.

B Experimental Details

We provide here further details about the experimental setup, the datasets we use, hyperparameters chosen for training the models, and data splits. We also include additional experiments that were omitted from the main paper for brevity, including experiments on identifying *demographic traits* in clinical narratives.

B.1 Clinical Narratives

B.1.1 Data

We describe here the *MIMIC-III i2b2-2006* and *i2b2-2010* datasets.

MIMIC-III. The *MIMIC-III* (Medical Information Mart for Intensive Care III) dataset is a large, publicly available database containing detailed and anonymized health-related data associated with over 40,000 patients who stayed in critical care units at the Beth Israel Deaconess Medical Center in Boston, Massachusetts between 2001 and 2012. *MIMIC-III* is a rich resource for researchers in various fields, such as medicine, data science, artificial intelligence, and healthcare analytics. The dataset contains a diverse range of data types, including demographics, vital signs, laboratory test results, medications, and clinical notes. The dataset contains over 2 million clinical notes contributed by over 3,500 distinct healthcare professionals, including doctors, nurses, and other clinicians, with an average of 571 notes per author.

The notes in the *MIMIC-III* dataset come in various types, reflecting the diverse aspects of patient care and documentation in the intensive care setting. Some of the most common note types include:

- Nursing/Progress notes: These are daily notes written by nurses or other care providers, documenting the patient’s progress, condition, and care provided.
- Radiology reports: Reports written by radiologists after interpreting medical imaging studies (e.g., X-rays, MRIs, CT scans).
- ECG reports: Reports documenting the interpretation of electrocardiogram results.
- Discharge summaries: Comprehensive summaries written by physicians when a patient is discharged from the hospital, outlining the patient’s hospital course, treatments, and follow-up instructions.
- Physician consult notes: Notes written by specialists when consulted by the primary care team to provide their expert opinion on specific medical issues.

- Pharmacy notes: Notes documenting medication-related information, including dosing, administration, and potential drug interactions.
- Social work notes: Notes related to the patient’s psychosocial status, including social and family support, living arrangements, and other relevant factors.

i2b2-2006. The i2b2 (Informatics for Integrating Biology and the Bedside) initiative is a collaborative effort that aims to develop new methods and tools for biomedical research. It focuses on the development of a scalable computational infrastructure that can be used to accelerate the translation of basic research findings into clinical applications. As part of this effort, i2b2 has hosted several shared tasks and challenges related to natural language processing and machine learning in healthcare.

In 2006, the first i2b2 challenge, known as the *i2b2-2006* challenge, was conducted, focusing on the identification of obesity and its comorbidities in discharge summaries. The dataset provided for the challenge contained 694 de-identified discharge summaries, which were randomly selected from the Research Patient Data Registry (RPDR) at Partners HealthCare. The dataset was divided into a training set of 514 discharge summaries and a test set of 180 discharge summaries. It is important to mention that the *i2b2-2006* dataset is relatively small compared to the *MIMIC-III* dataset and does not provide detailed information about the number of distinct authors or the average number of notes per author.

However, the discharge summaries typically include various sections such as patient demographics, admission and discharge dates, admission diagnoses, hospital course, procedures, medications, and follow-up plans. These summaries are generally written by physicians at the time of patient discharge, providing an overview of the patient’s medical condition, treatment received, and overall hospital stay.

i2b2-2010. The *i2b2-2010* challenge, also known as the i2b2/VA challenge, was a shared task organized by the i2b2 (Informatics for Integrating Biology and the Bedside) initiative in collaboration with the US Department of Veterans Affairs (VA). The challenge aimed to encourage the development of natural language processing (NLP) and machine learning techniques for extracting medical concepts from clinical narratives. Specifically, the *i2b2-2010* challenge focused on the identification of medical problems, tests, and treatments from free-text clinical records.

The dataset provided for the *i2b2-2010* challenge contained 826 de-identified clinical records, which were sourced from three different institutions: Partners HealthCare, the University of Pittsburgh Medical Center (UPMC), and the VA. The dataset was divided into a training set of 349 records and a test set of 477 records.

Similar to the *i2b2-2006* challenge, the *i2b2-2010* dataset is relatively small compared to the *MIMIC-III* dataset and does not provide detailed information about the number of distinct authors or the average number of notes per author. The clinical records in the dataset are composed of diverse note types, such as discharge summaries, progress notes, radiology reports, and pathology reports, contributed by physicians, nurses, and other healthcare professionals.

While the dataset does not provide specific information about the number of distinct authors, the fact that the notes were contributed by different types of healthcare professionals across multiple institutions increases the dataset’s diversity, making it more representative of real-world clinical settings.

B.1.2 PubMed BERT

In our clinical narratives experiments, we use *PubMed BERT* [84], a variant of the original BERT model [95], as our vanilla model. That is, all of the baselines and *CATO* all use it either for embedding clinical text or for predicting *conditions*, *demographic traits* and *note segments*.

PubMed BERT is a BERT-based (Bidirectional Encoder Representations from Transformers) model that has been pre-trained specifically on biomedical and scientific text data [84]. The model leverages the BERT architecture, which is a transformer-based deep learning model that has gained significant attention in natural language processing (NLP) for its state-of-the-art performance across a wide range of tasks.

PubMed BERT is pre-trained on a large corpus of approximately 14 million biomedical abstracts from the PubMed database, which is a comprehensive repository of biomedical literature. By pre-training the model on domain-specific data, *PubMed BERT* is expected to have a better understanding of biomedical concepts, terminology, and language patterns compared to general domain models like BERT-base and BERT-large [95].

The main advantage of using *PubMed BERT* for biomedical text mining tasks is its domain-specific knowledge, which can lead to improved performance and more accurate results when fine-tuned on various downstream tasks, such as named entity recognition, relation extraction, document classification, and question answering. Since *PubMed BERT* is pre-trained on a large corpus of biomedical text, it is better suited to capturing the unique language patterns, complex terminology, and the relationships between entities in the biomedical domain.

Hyperparameters for Fine-Tuning PubMed BERT on MIMIC-III. In our study, we leveraged a pre-trained *PubMed BERT* model and fine-tuned it on the *MIMIC-III* dataset. During pre-training, the model employed masked language modeling and next sentence prediction objectives. The architecture consisted of 12 layers, 768 hidden units, and 12 attention heads. For task-specific optimization, we used the following hyperparameters: a $3e-5$ learning rate with a linear warmup during the initial 10% of training steps, a batch size of 32, a maximum sequence length of 512 tokens, and a dropout rate of 0.1. The AdamW optimizer was applied with a 0.01 weight decay and a 1.0 gradient clipping threshold. To prevent overfitting, early stopping was based on validation loss and used a 3-epoch patience. The fine-tuning process ran for up to 20 epochs, unless early stopping criteria were met sooner.

The fine-tuning process was executed on a high-performance computing cluster with multiple NVIDIA Tesla V100 GPUs, each equipped with 32 GB of memory, using the *PyTorch* deep learning framework [96]. The dataset was preprocessed and tokenized using the *HuggingFace Transformers* library [97].

B.1.3 Demographic Traits Detection

Demographic Traits detection is the task of identifying residual private information in the clinical note, after removing the known identifier types (names, ages, dates, addresses, ID’s, etc.) [71]. We train all models on a subset of *MIMIC-III* and test on *i2b2-2006*. Table 3 presents our results. While performance gains from the Causal Augmentation approach are not as large as in the other clinical NLP tasks, it is still the best method in terms of *F1* score on out-of-distribution examples.

	ID (<i>MIMIC-III</i>)			OOD (<i>i2b2-2006</i>)		
	P	R	F1	P	R	F1
<i>PubMed BERT</i>	80.61	78.12	79.34	53.32	90.1	66.92
+ <i>Re-Weighting</i>	81.31	78.57	79.92	56.75	91.38	70.02
++ <i>MMD</i>	80.68	78.84	79.75	56.19	91.49	69.62
<i>Bio BERT</i>	79.5	77.63	78.55	53.32	89.84	66.71
<i>Sentence BERT</i>	79.29	76.18	76.53	52.22	89.82	65.04
<i>GPT3</i>	78.31	76.01	77.18	52.73	88.52	63.98
<i>Naive Aug.</i>	81.45	79.35	80.39	52.9	89.58	66.52
<i>Causal Aug.</i>	80.65	78.84	79.73	59.76	90.16	71.88

Table 3: Results (averaged across 5 runs) for predicting demographic traits from the text narratives on in-distribution and out-of-distribution data.

B.2 Restaurant Reviews

Data. We use the *CEBaB* dataset [49], which consists of short restaurant reviews and ratings from *OpenTable*, including evaluations for food, service, noise, ambiance, and an overall rating. For our experiments, we used the train-exclusive split of the dataset, which contains 1,755 examples.

To analyze the data, we transformed the overall rating into a binary outcome. The original rating scale ranges from 1 to 5, and we classified a rating of 3 or higher as 1, and anything below as 0. We

utilized a bag-of-words model with *CountVectorizer* and fitted logistic regression models from the *sklearn* library [98].

To investigate these questions, we construct two experimental settings: the original *CeBAB* dataset, and a modified version, denoted as *CeBAB-Spurious*, where there’s a spurious correlation between training and deployment.

The data is randomly split into a training set with 1,000 examples and a test set with 755 examples. We explore two data augmentation schemes:

1. Naive data augmentation: This approach involves randomly selecting two reviews from the dataset and prompting *GPT-4* [99] to rewrite one restaurant review in the style of the other. By applying the naive augmentation, we obtain an additional 1,000 training examples.
2. Conditional data augmentation : We match the ratings and sub-ratings in the reviews to create pairs. We then prompt *GPT-4* to rewrite one review to match the style of the other. Because not all pairs have matches in this case, the conditional data augmentation generates 926 augmentations. See Appendix B for details of the prompt.

Generating reviews with counterfactual food mentions. Following the counterfactual generation procedure in Algorithm 1, we generate counterfactual restaurant reviews conditional on food rating and overall rating. For each review, we first find a set of matched examples. We then select the subset that has different food-mention attribute and prompt *GPT-4* to rewrite. This results in 2,537 augmentations. The counterfactual augmentation should capture what the reviews should look like had a reviewer been more/less concise. Following Algorithm 1, we generate counterfactual restaurant reviews conditional on food and overall ratings. We find matched examples for each review, select those with different food-mentions, and prompt a *GPT-4* to rewrite them, reflecting how the reviews would appear if the reviewer was more/less concise.

Prompt Example.

```
helper_prompt = """
you are a very helpful, diligent, and intelligent language model assistant,
your task to generate counterfactual restaurant reviews,
that is what the restaurant review would be if it is given a different rating.
You will be given an original restaurant review and a comparator review
Your task is to rewrite the original review, such that it will have the same
review score as the comparator review.
The rating is with respect to ambiance, food, noise, and service.
---- EXAMPLE INPUT - START ----

original_review: [],
original_ratings: [
rating_ambiance: score,
rating_food: score,
rating_noise: score,
rating_service: score
]

compare_reviews: []
compare_ratings: [
rating_ambiance: score,
rating_food: score,
rating_noise: score,
rating_service: score
]

---- EXAMPLE INPUT - END ----
ANSWER FORMAT:
{
```

```

original_review: [],
original_score: [],
rewrite_review: [],
}

```

""

B.3 Synthetic Data

As described in the main paper we study a binary classification problem where $K = 8$ (cardinality of C), and sample $\tilde{P}(C | Y)$ to simulate varying degrees of the spurious correlation (specifically, we draw $\mathbf{x} = [\mathbf{x}^*, \mathbf{x}_{\text{spu}}]$ from a Gaussian distribution,

$$\mathbf{x}_i = \begin{bmatrix} \mathbf{x}_i^* \\ \mathbf{x}_{\text{spu},i} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_{y_i} \\ \boldsymbol{\mu}_{c_i} \end{bmatrix}, \begin{bmatrix} \sigma^2 \mathbf{I}_{d^*} & 0 \\ 0 & \sigma_{\text{spu}}^2 \mathbf{I}_{d_c} \end{bmatrix} \right).$$

In our simulations, we set $d^* = 10, d_{\text{spu}} = 300$ and $\sigma_{\text{spu}}^2 = 0.05, \sigma = 0.01d^*$ to make the max-margin classifiers depend on the spurious features. The parameters μ_{y_i}, μ_{c_i} are drawn uniformly from a sphere of norm 1/3 and 60, respectively. For the corruptions of augmentations where we add $\xi_i(\mu_c - \mu_{c_i})$, the ξ_i variables are drawn from a truncated Gaussian centered at λ with standard deviation 0.1.

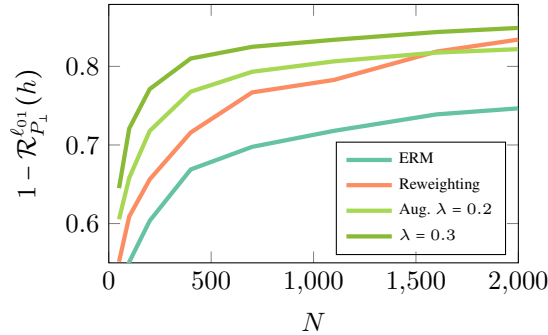


Figure 6: OOD accuracy ($1 - \mathcal{R}_{P_{\perp}}^{l_{01}}(h)$) for growing size of i.i.d training set N . We run 15 repetitions where $\tilde{P}(C | Y)$ are drawn randomly with correlation strength $I(Y; C) = 0.743 \pm 0.019$. With large amounts of data, the reweighting method approaches optimal performance and may outperform solutions based on corrupted data augmentation (e.g. it surpasses the more heavily corrupted data augmentation with $\lambda = 0.2$).

For the results in fig. 4 of the main paper we set the number of training examples N at 600 and the distributions $\tilde{P}(C | Y)$ are sampled such that for each interval of size 0.05 between 0 and 0.9 for the values of $I(Y; C)$, we draw 30 instances within that interval. In fig. 6 we give results for another experiment where we plot curves for reweighting, ERM and corrupted augmentation under several values of N under a strong spurious correlation. We draw values for $\tilde{P}(C | Y)$ such that that $I(Y; C)$ is in $[0.7, 0.8]$ (mean 0.743 and standard deviation 0.019 with 15 repetitions). Considering the bounds in eq. (2) and the one in lemma 2, we expect that as N grows the reweighting method will approach optimal accuracy, while the augmentation methods will incur some bias due to the imperfect approximation of counterfactuals. In fig. 6 we see that the augmentation based methods begin to plateau, while the reweighting method gains larger improvement from the growing N , and outperforms the more heavily corrupted augmentation method ($\lambda = 0.2$) with a large enough training set. The best possible mean accuracy in this simulation is 0.934 (as we gather from training a classifier from \mathbf{x}^* alone), the best accuracy obtained by the methods in the figure is 0.849 (Aug. $\lambda = 0.3$ with $N = 1999$).

C Synthetic Note Examples

While it is not possible to share real notes from the *MIMIC-III* or any of the *i2b2* datasets, we can share how a typical note might look like by generating synthetic notes with LLMs. Here, we generate a synthetic note using *GPT-4*, and then generate a counterfactual note using our prompt:

Discharge Summary Note

Dr. Jane Doe

March 29, 2023

Patient Name: John Doe

Admission Date: March 25, 2023

Discharge Date: March 29, 2023

Diagnosis: Pneumonia

History of Present Illness: Mr. Doe is a 58-year-old male with a past medical history of hypertension who presented to the emergency department with a 3-day history of fever, cough, and shortness of breath. Chest X-ray showed bilateral infiltrates, and he was admitted for treatment of community-acquired pneumonia.

Hospital Course: The patient was started on intravenous antibiotics and supplemental oxygen. Over the course of his hospitalization, his oxygen requirements decreased, and his fever and cough resolved. He was able to tolerate a regular diet and was mobilized with physical therapy.

Consultations: Infectious disease was consulted, and the patient's antibiotics were adjusted based on culture results.

Discharge Medications:

1. Levofloxacin 750mg daily for 10 days
2. Acetaminophen 650mg every 6 hours as needed for fever or pain
3. Albuterol inhaler 2 puffs every 4 hours as needed for shortness of breath

Follow-up: The patient was instructed to follow up with his primary care physician in 1 week for a repeat chest X-ray and evaluation of his symptoms.

Discharge Instructions:

1. Finish the entire course of antibiotics as prescribed.
2. Use the albuterol inhaler as needed for shortness of breath.
3. Rest and avoid strenuous activity until cleared by your doctor.
4. Follow up with your primary care physician in 1 week.

Discharge Condition: The patient was discharged in stable condition with improved symptoms.

Here's a counterfactual discharge summary note written by a doctor with a different writing style:

<p>Patient Name: John Doe Admission Date: March 25, 2023 Discharge Date: March 29, 2023 Diagnosis: Pneumonia HPI: John is a middle-aged man who came in with a fever and cough that he's had for three days. He was diagnosed with pneumonia and admitted. Hospital Course: John was given antibiotics and oxygen. He slowly got better and was able to eat normally and move around more. He saw an infectious disease specialist who adjusted his treatment. Consultations: Infectious disease saw John and changed his medicine. Discharge Medications:</p> <ol style="list-style-type: none">1. Levofloxacin 750mg once a day for 10 days2. Acetaminophen 650mg every 6 hours as needed for fever or pain3. Albuterol inhaler 2 puffs every 4 hours as needed for shortness of breath <p>Follow-up: Follow up with PCP in 1 week. Discharge Instructions:</p> <ol style="list-style-type: none">1. Finish your antibiotics.2. Use the inhaler if you need it.3. Rest and avoid heavy activity until you feel better.4. Follow up with your doctor next week. <p>Discharge Condition: Stable, going home.</p>
--

As can be seen from these examples, the counterfactual note is much more concise and to-the-point than the original example. The language used is more direct and less descriptive, and there is less detail provided about the patient's course of treatment.

D Possible Limitations of LLMs in Generating Augmented Datasets

As mentioned in our discussion, there are several possible limitations that should be carefully considered before applying our approach in practice, especially in high-stakes applications such as medical notes classification. We list some of the main possible limitations and points to consider, along with a short discussion on each.

- *LLM generation quality:* LLMs vary in their ability to generate realistic text. It is possible that LLMs introduce biases into our problem, inherited from their own training data. This requires further study, however from our manual examination we found their quality satisfactory (see appendix C for generation examples) and that OOD generalization also improved for models trained on the augmented data they generate. We also include experiments with several types of LLMs in appendix B to verify that our findings are consistent across the types of LLMs we considered.
- *Counterfactual approximation:* Other than generation quality, the additional challenge in using LLMs for counterfactual data augmentation is our ability to elicit a good approximation to the counterfactual text. Our methods rely on principles from causal inference to advance disciplined approaches for this task. While further studies are required (e.g. systematically comparing small sets of manual re-writes of texts to the elicited LLM output), we view our work as a promising first step in this direction, which we expect to be significantly extended and improved in future work.
- *Effect of biases on OOD generalization:* Since we focus on OOD generalization, the limitations and possible biases mentioned above must be weighed within this context. Namely, we should bear in mind that even though generation may be biased, this bias is only harmful when it affects the generalization of a downstream classifier, and this is what we evaluate. Further, in OOD generalization we consider cases where the training data is biased in the first place, and training a standard predictive model also results in a biased solution. Hence we must weigh risks and limitations of alternative solutions vs. those of LLMs.