

Does It Make Sense to Explain a Black Box With Another Black Box?

Anonymous ACL submission

Abstract

Although counterfactual explanations are a popular approach to explain ML black-box classifiers, they are less widespread in NLP. Most methods find those explanations by iteratively perturbing the target document until it is classified differently by the black box. We identify two main families of counterfactual explanation methods in the literature, namely, (a) *transparent* methods that perturb the target by adding, removing, or replacing words, and (b) *opaque* approaches that project the target document into a latent, non-interpretable space where the perturbation is carried out subsequently. This article offers a comparative study of the performance of these two families of methods on three classical NLP tasks. Our empirical evidence shows that opaque approaches can be an overkill for downstream applications such as fake news detection or sentiment analysis since they add an additional level of complexity with no significant performance gain. These observations motivate our discussion, which raises the question of whether it makes sense to explain a black box using another black box.

1 Introduction

The latest advances in machine learning (ML) have led to significant advances in various natural language processing (NLP) tasks (Devlin et al., 2019; Liu et al., 2019; Sanh et al., 2019), such as text generation, fake news detection, sentiment analysis, and spam detection. These notable improvements can be partly attributed to the adoption of methods that encode and manipulate text data using latent representations. Those methods embed text into high-dimensional vector spaces that capture the underlying semantics and structure of language, and that are suitable for complex ML models.

Despite the impressive gains in accuracy achieved by modern ML algorithms (Devlin et al., 2019; Brown et al., 2020), their utility can be diminished by their lack of interpretability (Shen et al.,

2020). This has, in turn, raised an increasing interest in ML explainability, the task of providing appropriate explanations for the answers of black-box ML algorithms (Jacovi, 2023). Indeed, a model could make correct predictions for the wrong reasons (Gururangan et al., 2018; McCoy et al., 2019). Unless the ML model is a white box, explaining the results of such an agent requires an explanation layer that elucidates the internal workings of the black box in a post-hoc manner.

While there are several ways to explain the outcomes of an ML model a posteriori, there has been a growing emphasis on counterfactual explanations, a domain that has experienced notable popularity over the last five years (Guidotti, 2022; Miller, 2019). A counterfactual explanation is a counter-example that is similar to the original text, but that elicits a different outcome in the black box (Wachter et al., 2018). Consider the classifier depicted in Figure 1, for sentiment analysis applied to the review “This is a good article” – classified as positive. In this toy example, a counterfactual could be the phrase “This is a **poor** article”. This explanation tells us that the adjective “good” was a possible reason for this sentence to be classified as positive, and changing the polarity of that adjective may change the classifier’s response.

Counterfactual explanation methods operate by increasingly perturbing the target text until the answer from the model – often a classifier – changes. Those perturbations can be conducted *transparently* by adding, removing, or changing words and syntactic groups (Martens and Provost, 2014; Yang et al., 2020; Ross et al., 2021) in the original target text as depicted in Figure 1. Since removing or adding words from a text can lead to unrealistic texts, more recent methods (Hase and Bansal, 2020; Robeer et al., 2021; S. Punla and C. Farro, 2022) embed the target text in a latent space that captures the underlying distribution of the model’s training corpus. Perturbations are then carried out

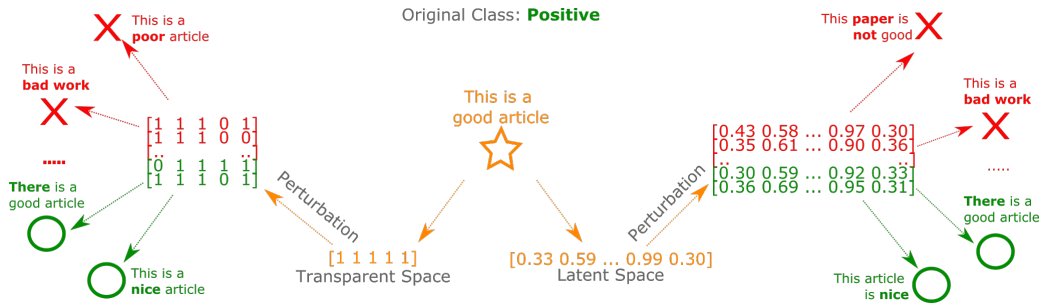


Figure 1: The mechanism employed to perturb the target documents by the transparent and opaque methods. Transparent techniques, on the left, convert the input text to a vector representation, where ‘1’ indicates the presence of the input word and ‘0’ denotes a replacement. Opaque methods, as on the right, embed words from the target text into a latent space and perturb the text in this high-dimensional space.

in this space and then brought back to the space of words to guarantee realistic counterfactual explanations. These explanation methods rely on *opaque* sophisticated techniques to compute those explanations (Li et al., 2021), which is tantamount to explaining a black box with another black box.

Based on this somehow paradoxical observation, we conduct a comparative study of various transparent and opaque post-hoc counterfactual explanation approaches. Rather than two distinct categories, the studied methods define a continuum, as some methods may combine transparent and non-interpretable techniques. Our study aims to understand whether it is worth resorting to latent approaches to explain complex ML models. The experimental results suggest that for some downstream NLP tasks, learning a latent representation for explanation purposes can be an overkill. To strengthen our point, we present and evaluate two novel transparent approaches for counterfactual explanations.

The paper is structured as follows. Section 2 surveys the existing counterfactual explanation methods. Section 3 introduces two novel transparent methods, which we then analyze in the light of the spectrum of existing transparent and opaque techniques (Section 4). We then elaborate on the experimental protocol of our comparative study in Section 5. The results of our experimentations are presented in Section 6. Section 7 discusses our findings and concludes the paper.

2 Related Works

Counterfactual explanation methods compute contrastive explanations for ML black-box algorithms by providing examples that resemble a target instance but that lead to a different answer in the black box (Wachter et al., 2018). These counterfactual explanations convey the minimum changes in

the input that would modify a classifier’s outcome. Social sciences (Miller, 2019) have shown that human explanations are contrastive and Wachter et al. (2018) have illustrated the utility of counterfactual instances in computational law. When it comes to NLP tasks, a good counterfactual explanation should be fluent (Morris et al., 2020), *i.e.*, read like something someone would say, and be sparse (Pearl, 2009), *i.e.*, look like the target instance.

Counterfactual approaches have gained popularity in the last few years. As illustrated by the surveys, first by Bodria et al. (2021) and later by Guidotti (2022), around 50 additional counterfactual methods appeared in a one-year time span. Despite this surge of interest in counterfactual explanations, their study for NLP applications remains underdeveloped (Ross et al., 2021). In the following, we elaborate on the existing counterfactual explanation methods for textual data along a spectrum that spans from transparent to opaque approaches.

Transparent Approaches. Given an ML classifier and a target text (also called a document), transparent techniques compute counterfactual explanations in a binary space. Each dimension represents the presence (1) or absence (0) of a word from a given vocabulary. Hence, to perturb a text, these methods toggle on and off 0s and 1s, where 0s are tantamount to adding, removing, or replacing words until the classifier yields a different answer. This was first proposed by Martens and Provost (2014) who introduced Search for Explanations for Document Classification (SEDC), a method that removes the words for which the classifier exhibits the highest *sensitivity*. More recently, Ross et al. (2021) developed Minimal Contrastive Editing (MICE), a method that employs a Text-To-Text Transfer Transformer to fill masked sentences.

Yang et al. (2020) presented Plausible Counterfactual Instances Generation (PCIG), which generates grammatically plausible counterfactuals through edits of single words with lexicons manually selected from the economics domain.

Opaque Methods. We define opaque approaches as those perturbing the input text in a latent space in \mathbb{R}^n . Methods such as Decision Boundary (Hase and Bansal, 2020), xSPELLS (S. Punla and C. Farro, 2022) or cfGAN (Robeer et al., 2021) operate in three phases. First, they embed the target text onto a latent space. This is accomplished by employing specific techniques such as Variational AutoEncoder (VAE) in the case of xSPELLS, or a pre-trained language model (LM) for cfGAN. Second, while the classifier’s decision boundary is not traversed, these methods perturb the latent representation of the target phrase. This is done by adding Gaussian noise in the case of xSPELLS, whereas cfGAN resorts to a Conditional Generative Adversarial Network. Finally, a decoding stage produces sentences from the latent representation of the perturbed documents.

There also exist methods such as Polyjuice (Wu et al., 2021), Generate Your Counterfactuals (GYC) (Madaan et al., 2021) and Tailor (Ross et al., 2022) that perturb text documents in a latent space, but can be instructed to change particular linguistic aspects of the target text, such as locality or grammar tense. Such methods are not particularly designed to compute counterfactual explanations but are rather conceived for other applications such as data augmentation.

Unlike pure word-based perturbation methods, latent representations are good at preserving *semantic closeness* for small perturbations. That said, these methods are not free of pitfalls. First, methods such as xSPELLS and cfGAN are deemed opaque since a latent space is not human-understandable (Shen et al., 2020). Moreover, existing latent-based approaches do not seem optimized for sparse counterfactual explanations – one of the defining features of a counterfactual. We show this through our experimental results that suggest that a minor alteration in the latent space can cause a significant alteration in the original space.

3 Two Novel Transparent Methods

Before elaborating on our study, we introduce two novel counterfactual explanation techniques, aimed to enrich the middle ground between fully opaque

and fully transparent approaches. The methods are called Growing Language and Growing Net, and both depend on an iterative process that replaces words within a target text $x = (x_1, \dots, x_d) \in X$ ($x_i \in \Sigma$ are words from a vocabulary Σ) until the predicted class of a given classifier $f : X \rightarrow Y$ changes. The goal of such a procedure is to compute sparse counterfactual explanations with the fewest modified words.

Algorithm 1 Explore

Require: target text $x = (x_1, \dots, x_d) \in X$, classifier f ;
SIMWORDS(\cdot) \rightarrow retrieves similar words
Hyper-parameters: $n = 2000$
Ensure: one or multiple counterfactual instances
1: **for** $i \leftarrow 1$ **to** d **do**
2: $W_i \leftarrow \text{SIMWORDS}(x_i, \text{POS}(x_i))$
3: **end for**
4: Initialize $Z = (z^1, \dots, z^n)$ as n copies of x
5: Initialize $C \leftarrow \emptyset$; $n_m \leftarrow 0$
6: **while** $n_m < d \wedge C = \emptyset$ **do**
7: $n_m \leftarrow n_m + 1$
8: **for** $j \leftarrow 1$ **to** n **do** \triangleright For each copy of x
9: **for** $l \leftarrow 1$ **to** n_m **do**
10: $k \leftarrow \text{random}(0, d)$ $\triangleright k : z_k^j = x_k$
11: $z_k^j \leftarrow \text{random word from } W_k$
12: **end for**
13: **if** $f(x) \neq f(z_j)$ **then**
14: $C \leftarrow C \cup \{z_j\}$
15: **end if**
16: **end for**
17: **end while**
18: **return** C

Algorithm 2 Growing Net

Require: a target text $x = (x_1, \dots, x_d) \in X$, classifier f ;
1: $C \leftarrow \text{explore}(x, f, \text{WN_SIMWORDS}_{d=1}(\cdot))$
2: **return** $\text{argmax}_{c \in C} \text{Wu-P}(c, x)$

Algorithm 3 Growing Language

Require: target text $x = (x_1, \dots, x_d) \in X$, classifier f ;
Hyper-parameters: $\tau = 0.02$; $\theta = 0.9$; $\theta_{\min} = 0.4$;
1: $C \leftarrow \emptyset$
2: **while** $\theta > \theta_{\min} \wedge C = \emptyset$ **do**
3: $C \leftarrow C \cup \text{explore}(x, f, \text{LM_SIMWORDS}_{\theta}(\cdot))$
4: $\theta \leftarrow \theta - \tau$
5: **end while**
6: **return** $\text{argmin}_{c \in C} \|x - c\|_0$

Algorithm 1 outlines the iterative exploration process employed by Growing Language and Growing Net. In the first step (lines 1 to 3), both approaches generate d sets of potential word replacements W_1, \dots, W_d for each word x_i in the target document x . Those replacements must have the same part-of-speech (POS) tag as x_i . The external module to obtain those word replacements depends on the method. These modules are detailed later. Subsequently, our methods create arti-

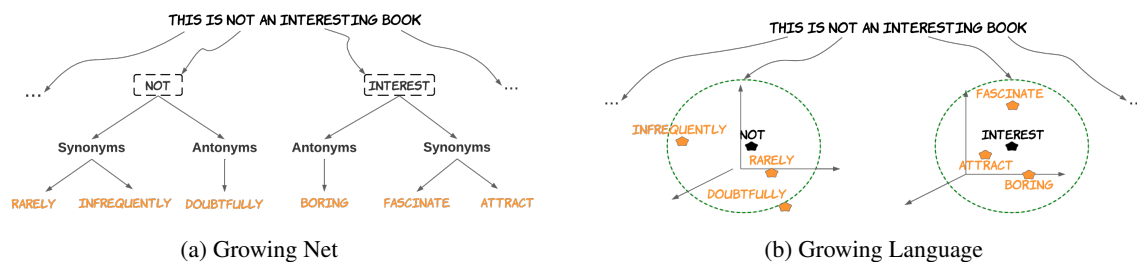


Figure 2: The mechanism to compute potential word replacements in Growing Net navigates the tree structure of WordNet. Conversely, Growing Language embeds words into a latent space on which it looks for nearby words.

227 ficial documents iteratively (lines 6 and 17) while
 228 some words in the original document remain non-
 229 replaced ($n_m < d$), or while we have not found any
 230 counterfactuals. At each iteration, the exploration
 231 keeps n copies of the original text (x) on which we
 232 replace n_m individual words (x_k) with randomly
 233 selected words from their respective sets of poten-
 234 tial replacements (W_k). Lines 13-15 check if the
 235 resulting phrases are counterfactual instances.

236 For example, consider the target review, “*This is*
 237 *not an interesting book*”, classified as negative by
 238 a sentiment analysis model. In the first round, our
 239 routine produces artificial reviews with only one
 240 modified word. Subsequent rounds will replace
 241 two words and so on (lines 9 to 12).

242 **Growing Net.** This method capitalizes on the
 243 rich structure of WordNet (Fellbaum, 1998) to
 244 identify potential word replacements. WordNet
 245 is a lexical database and thesaurus that organizes
 246 words and their meanings into a semantic tree of
 247 interrelated concepts. The method is described
 248 in Algorithm 2, and uses the module `WN_SIM-`
 249 `WORDSd`. In the exploration phase, Growing Net
 250 uses `WN_SIMWORDSd` to find words at a distance
 251 of at most d in the WordNet hierarchy among syn-
 252 onyms, antonyms, hyponyms, and hypernyms for a
 253 given word x_i to replace. This process is illustrat-
 254 ed in Figure 2a. In our experiments we set $d = 1$ as
 255 this value already yields good results – higher val-
 256 ues would incur longer runtimes. The exploration
 257 returns a set of counterfactuals, from which Grow-
 258 ing Net selects the one with the highest Wu-Palmer
 259 Similarity (Wu-P) (Wei and Ngo, 2007) as final
 260 explanation. This similarity score for text relies on
 261 Wordnet, and takes into account the relatedness of
 262 the concepts in the phrase, e.g., via the path length
 263 to their most common ancestor in the hierarchy.

264 **Growing Language.** This approach leverages
 265 the power of language models (LM) to restrict the

266 space of possible word replacements via the mod-
 267 ule `LM_SIMWORDSθ` (see Algorithm 3). Given
 268 a word x_i to replace, `LM_SIMWORDSθ` embeds
 269 the word onto the latent space of an LM, as il-
 270 lustrated in Figure 2b. Then `LM_SIMWORDSθ` re-
 271 trieves words whose latent representation is at a
 272 distance of $θ$ at most. In our experiments, we ini-
 273 tially set this threshold to 0.8 on a scale from 0
 274 to 1. If for a given $θ$, Growing Language cannot
 275 find counterfactual instances, the distance thresh-
 276 old is relaxed, i.e., reduced by $τ$ (set to 0.02 in
 277 our experiments), so that the exploration routine
 278 considers more words. Should multiple counter-
 279 factuals be found, Growing Language selects the
 280 one with the fewest modifications compared to the
 281 original document (minimal L0 distance). For our
 282 experiments, we employed Spacy (Honnibal and
 283 Montani, 2017), but any language model capable
 284 of embedding words and offering word distances
 285 could be applied in this context.

4 Interpretability Spectrum 286

287 We have presented counterfactual explanation tech-
 288 niques as either opaque or transparent. However,
 289 the landscape is more nuanced, for these techniques
 290 actually define a spectrum, which we depict in Fig-
 291 ure 3. The spectrum spans from the most transpar-
 292 ent methods on the left to the most opaque ones on
 293 the right. We elaborate on the various regions of
 294 this spectrum in the following.

295 **Fully Transparent.** At the leftmost end of the
 296 spectrum, we find the method SEDC (Martens
 297 and Provost, 2014), which perturbs text instances
 298 by hiding only highly sensitive words within
 299 the text. We place Growing Net on the right
 300 of SEDC, because it goes beyond simple word
 301 masking. Instead, it substitutes words judiciously
 302 via an external interpretable asset, namely Wordnet.
 303

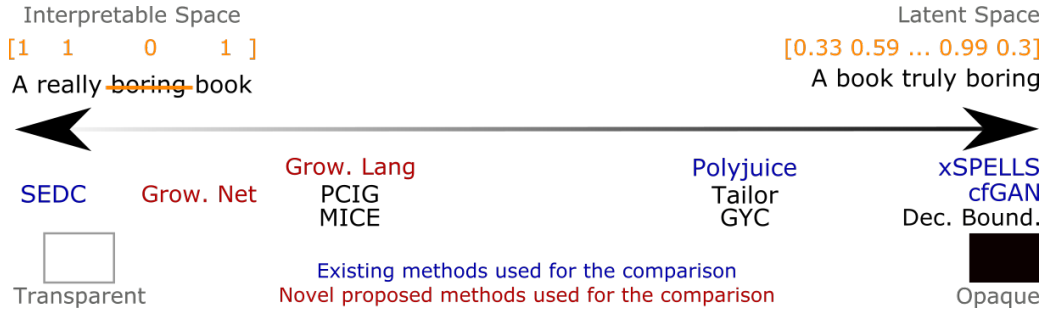


Figure 3: Spectrum for counterfactual explanation techniques that goes from the most transparent methods on the left to the most opaque on the right. Transparent methods perturb documents in a binary space; opaque methods do it in a latent space.

Transparent. Methods like PCIG (Yang et al., 2020), MICE (Ross et al., 2021), and Growing Language are considered more opaque than Growing Net, because they employ a latent space to identify semantically close word substitutions. Despite this reliance on black-box techniques, we consider them transparent because the search for counterfactuals is still carried out in the space of words.

Partially Opaque. Polyjuice, Tailor, and GYC fall in the category of partially opaque methods, as they leverage control codes to perturb the target document. Control codes are specific instructions that adapt the perturbation of the target text so that it complies with a specific task, such as translating, summarizing, or changing the tense of a text. While these modifications occur in a latent space, the inclusion of control codes provides some level of clarity regarding why a modification influences the model’s prediction.

Fully Opaque. On the far right of the interpretability spectrum, we encounter fully opaque approaches such as Decision Boundary, xSPELLS and cfGAN. These methods perturb instances in a latent space, making it challenging for users to discern the underlying process of counterfactual generation.

This interpretability spectrum provides valuable insights into the transparency and opacity of counterfactual explanation methods, allowing for a more nuanced understanding of their capabilities.

5 Experimental Protocol

Having introduced the spectrum of counterfactual explanation methods across the interpretability axis, we now describe the experimental setup designed to evaluate those methods. The code of the studied methods, the datasets, and the experimental results are available at <https://anonymous.4open.science/r/ebbwb-4B55/README.md>

5.1 Methods

We picked a set of representative domain-agnostic methods from all regions of the spectrum depicted in Figure 3. These include SEDG and Growing Net among the fully transparent methods, Growing Language among the transparent ones¹, Polyjuice among the partially opaque ones, and xSPELLS and cfGAN from the fully opaque group.

5.2 Tasks & Datasets

We conduct the evaluation on three popular downstream tasks: (a) spam detection in messages, (b) sentiment analysis, and (c) detection of fake news from newspaper headlines. The datasets associated to these tasks consist of two target classes, and contain between 4000 and 10660 textual documents. The average number of words in each document is between 11.8 and 20.8 as reported in Table 1. Except for the fake news dataset, we downloaded the data from Kaggle. The fake news dataset was constructed by us and its description is available in our repository <https://anonymous.4open.science/r/ebbwb-4B55/README.md>.

Dataset	No. of words		Instances	Accuracy (%)		
	Total	Average		MLP	RF	BERT
Fake	19419	11.8	4025	84	84	91
Polarity	11646	20.8	10660	72	67	82
Spam	15587	18.5	8559	100	100	100

Table 1: Information about the experimental datasets. The “average” column denotes the average number of words per instance (document).

5.3 Black-box Classifiers

Our evaluation uses two distinct black-box classifiers implemented using the scikit-learn library and

¹PCIG relies on domain specific rules from economics; MICE is computationally expensive according to the authors.

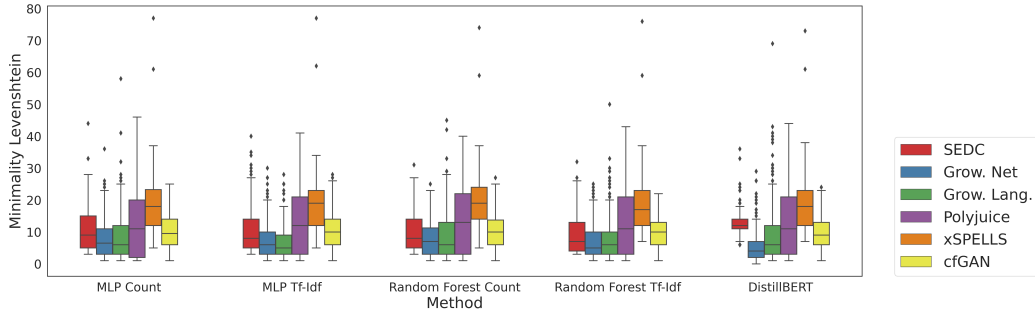


Figure 4: Minimality as the levenshtein edit distance between the closest counterfactual and the target text (\downarrow better).

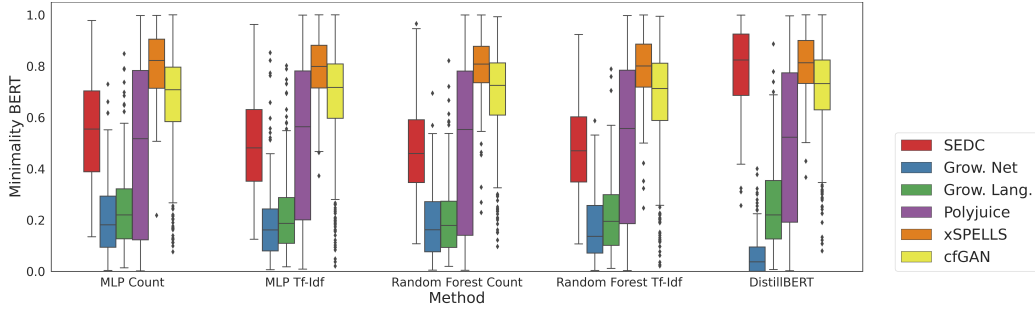


Figure 5: Minimality as the Sentence-BERT embedding distance between the closest counterfactual and the target text (\downarrow better).

367 already employed in (S. Punla and C. Farro, 2022).
 368 These black boxes are (i) a Random Forest (RF)
 369 consisting of 500 tree estimators, (ii) a multi-layer
 370 perceptron (MLP) with token counts as input, and
 371 (iii) a classifier based on DistillBERT². For the RF
 372 and the MLP, we employed both the *token count*
 373 and *tf-idf* vectorizers to convert text into proper
 374 inputs for the models.

375 We used 70% of the instances for training, and
 376 the remaining for testing. The classifiers’ test per-
 377 formances are shown in Table 1. The counterfactual
 378 explanations were computed for instances in
 379 those test sets.

380 6 Results

381 We now present the results of our evaluation, orga-
 382 nized in four rounds of experiments categorized ac-
 383 cording to two aspects. First, we assess the quality
 384 of the produced counterfactual explanations based
 385 on two essential criteria: (i) **minimality**, and (ii)
 386 **plausibility**. Second, we evaluate the methods
 387 themselves in terms of (iii) **flip change**, and (iv)
 388 **runtime**. For each evaluated method and black-
 389 box classifier, we computed counterfactual expla-
 390 nations for 100 target texts extracted from the test
 391 sets of our datasets.

²<https://is.gd/zljjJN>

392 6.1 Counterfactual Quality

393 A high-quality textual counterfactual explanation
 394 tells us what are the most sensitive parts or as-
 395 pects of the target phrase, that otherwise changed,
 396 would lead to a different classification outcome. It
 397 follows then that such an explanation must (i) in-
 398 cur minimal changes w.r.t the target phrase (sparse
 399 changes), and (ii) be linguistically plausible, i.e.,
 400 sound like something a person would naturally
 401 write or say (Guidotti, 2022).

402 **Minimality.** We quantify the minimality crite-
 403 rion by measuring the distance between the counterf-
 404 actual and the target sentence. Figure 4 and 5 display
 405 the results of our minimality assessments, consid-
 406 ering both the Levenshtein distance and the cosine
 407 similarity within the embedding space of the BERT-
 408 Sentence model (Reimers and Gurevych, 2019).
 409 This dual approach ensures a comprehensive eval-
 410 uation, accounting for both lexical similarity and
 411 latent features, including aspects of style.

412 Notably, our findings reveal that methods posi-
 413 tioned in the middle-ground, particularly Growing
 414 Net, performed favorably compared to opaque ap-
 415 proaches, both in terms of the number of words
 416 modified and semantic comparison. It is worth not-
 417 ing that xSPELLS introduced the most significant
 418 changes to the original text – contradicting one
 419 of the main functional requirements of a counter-
 420 factual explanation (Wachter et al., 2018). Simi-

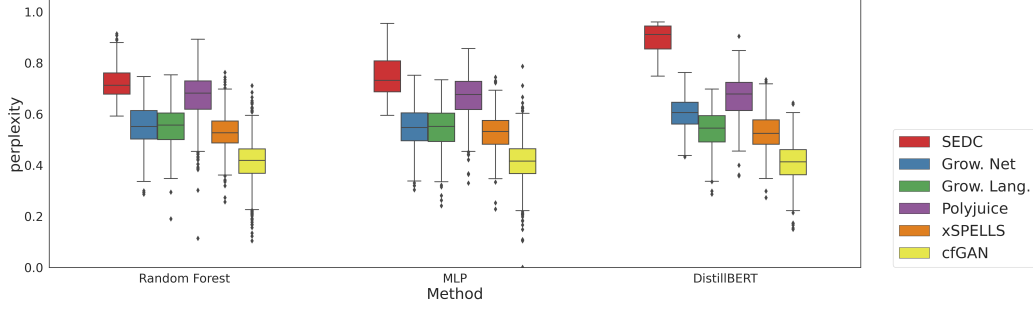


Figure 6: Perplexity as the MSE loss of a GPT model on the generated counterfactuals (\downarrow better).

Dataset	Fake			Spam			Polarity		
	MLP	RF	BERT	MLP	RF	BERT	MLP	RF	BERT
SEDC	0.95	0.82	1	0.47	0.42	0.56	0.92	0.93	0.98
Grow. Net	0.90	0.8	0.88	0.44	0.29	0.84	0.97	0.98	0.90
Grow. Lang.	0.84	0.84	0.77	0.58	0.61	0.17	0.92	0.92	0.92
Polyjuice	0.26	0.23	0.21	0.17	0.14	0.16	0.33	0.31	0.29
xSPELLS	0.68	0.78	0.77	0.98	0.95	0.91	0.91	0.76	0.91
cfGAN	0.18	0.12	0.09	0.14	0.05	0.03	0.50	0.50	0.48

Table 2: Average label flip per dataset and black box of the six counterfactual methods (\uparrow better).

larly, we observe a high variance in the minimality of the counterfactuals generated by Polyjuice, indicating that some counterfactuals were notably distant from their corresponding target instances. While these methods introduced minor perturbations to the original text, these modifications occurred within a latent space. Nothing guarantees, however, that these minor adjustments translate into visually subtle modifications of the target phrase when the resulting phrase is brought back to the original space. As an example, consider the target text “This is one of Polanski’s best films.” from the polarity dataset. For the DistillBERT classifier, cfGAN returns the counterfactual “this is one of **shot kingdom intelligence’ s all**”, which looks completely unrelated to the target text. Conversely, the transparent method SEDC produces the counterfactual “This is one of **MASK MASK MASK**”, whereas Growing Language outputs “This is one of Polanski’s **worst** films.”.

Additionally, we noted first that when the complexity of the classifier increases, the counterfactual explanations generated by SEDC lie farther from the original text. Secondly, we observe minor variations dependent on the vectorizer employed by the classifiers (*count* or *tf-idf*). Hence, for the subsequent phase of the evaluation, we present results exclusively for the *tf-idf* vectorizer.

Plausibility. While linguistic plausibility is typically evaluated through user studies (Madaan et al., 2021; Wu et al., 2021), we approximate it here following the techniques from Ross et al. (2021, 2022). Thus, we use perplexity scores based on a GPT language model (Brown et al., 2020), by calculating the average mean squared error (MSE) loss when predicting every token in the counterfactual from the previous ones. Figure 6 presents the plausibility of the counterfactuals. To enhance comparability, we normalized perplexity scores based on the maximum perplexity observed across the entire set of counterfactuals, where lower scores indicate higher plausibility. Notably, SEDC and Polyjuice generated texts with the lowest plausibility, which is expected since SEDC masks words, leading sometimes to nonsensical sentences. In contrast, cfGAN demonstrated the highest plausibility, while both Growing Net and Language achieved perplexity scores similar to those of xSPELLS.

6.2 Method Quality

We now compare the quality of the counterfactual explanation methods themselves based on (iii) label flip rate, which measures how frequently a method produces an instance classified differently by the model, and (iv) runtime, the time it takes for each method to generate a counterfactual explanation.

dataset	method	MLP	RF	BERT
fake	SEDC	31 (14)	13 (6)	15 (3)
	Grow. Net	2 (1)	1 (1)	7 (1)
	Grow. Lang.	55 (28)	55 (13)	34 (12)
	Polyjuice	38 (8)	70 (185)	29 (4)
	cfGAN	1 (0)	1 (0)	1 (0)
	xSPELLS	84 (6)	86 (7)	16 (1)
spam	SEDC	21 (13)	16 (9)	16 (6)
	Grow. Net	1 (1)	1 (1)	11 (4)
	Grow. Lang.	60 (16)	57 (14)	88 (43)
	Polyjuice	32 (7)	62 (184)	33 (15)
	cfGAN	1 (0)	1 (0)	1 (0)
	xSPELLS	219 (17)	198 (16)	22 (1)
polarity	SEDC	13 (10)	12 (9)	21 (6)
	Grow. Net	1 (1)	1 (1)	9 (2)
	Grow. Lang.	75 (33)	74 (32)	65 (29)
	Polyjuice	81 (30)	82 (48)	29 (4)
	cfGAN	1 (0)	1 (0)	1 (0)
	xSPELLS	136 (19)	115 (11)	24 (2)

Table 3: Average runtime in seconds of the studied counterfactual methods (and standard deviation).

Label flip rate. Table 2 provides an overview of the label flip results. It is noteworthy that except for the spam dataset, transparent methods achieve the highest label flip rate. This highlights the effectiveness of replacing words with antonyms as a means to discover counterfactuals. Additionally, xSPELLS exhibits strong performance for the spam dataset and similar label flip rates to transparent methods on polarity. We also emphasize that both Growing Net and Growing Language can be fine-tuned for a more exhaustive search by adjusting their parameters, for example by lowering the minimal similarity threshold (θ_{min} in Alg. 3) or by going further in WordNet’s tree structure (higher d in Alg. 2). While this can enhance the label flip rate, it may result in longer runtimes.

Runtime. Finally, Table 3 details the average and standard deviation of the runtime for each counterfactual explanation method across datasets and classifiers. Notably, cfGAN and Growing Net emerged as the fastest methods for generating counterfactuals. However, it is important to note that cfGAN requires the training of the Variational AutoEncoder (VAE) on each specific dataset, a process that incurs long training times. The time needed for fine-tuning varies, ranging from 4300 seconds for fake news title detection to 6755 seconds for spam detection. Furthermore, we observe that xSPELLS and Growing Language exhibit the slowest runtime performance. Growing Language, for instance, requires approximately 60 seconds to generate a single counterfactual, while xSPELLS exhibits run-

times, ranging from 16 seconds for fake news detection to 219 seconds for spam detection. These results reveal that, in contrast to opaque methods such as xSPELLS, transparent approaches like Growing Net are fast enough for real-time explainability.

7 Discussion & Conclusion

Our evaluation provides valuable insights into the landscape of counterfactual explanations for downstream NLP tasks. One of the most striking findings is that complexity, often associated with the use of neural networks and latent spaces, does not necessarily equate to superior performance in this context. Surprisingly, our results demonstrate that simpler approaches, characterized by a systematic and judicious strategy for word replacement, consistently yield satisfactory outcomes across all quality dimensions. The results of our study prompt a deeper reflection on the optimal strategies for generating counterfactual explanations in the field of NLP. It invites readers to embrace simplicity and transparency whenever the constraints of the application allow it.

Furthermore, our findings underscore the critical importance of transparency and interpretability in AI and ML, especially in high-stakes applications. The paradox of explaining a black box with another one calls into question the development of opaque approaches when transparent methods suffice, or when transparency is one of the goals in the first place. When focused on NLP applications, our results also call for reflection on the meaning and goal of explanations. If the task is to understand which aspects of a text should change to get a different outcome, a counterfactual explanation that drastically changes every word in the text may not be understandable. On the contrary, a counterfactual based on simple word-masking, albeit simple, may be perceived as implausible. This could hamper the goal of explanations as a means to elicit trust in users.

We therefore expect our findings to encourage the development of more transparent and interpretable AI systems that foster trust and accountability in every step of the AI-driven decision-making processes, either for prediction, recommendation, or explanation. Last but not least, we believe that the lessons drawn from this paper could be naturally ported to other explanation paradigms.

556 Limitations

557 We remind the reader that the evaluation was con-
558 ducted on three well-studied downstream applica-
559 tions, namely polarity analysis, fake news detec-
560 tion, and spam detection. Our results might there-
561 fore not generalize to other NLP tasks in special-
562 ized domains or different languages. While this
563 work puts transparent approaches in the spotlight,
564 our results suggest that plausible counterfactual ex-
565 amples need external domain-adapted knowledge
566 either in the form of language models or knowledge
567 graphs. These may not always be available though.
568 Finally, our evaluation was based on popular cri-
569 teria and metrics for counterfactual explanations.
570 Specialized applications may still take into account
571 additional criteria such as diversity or actionability.

572 References

573 Francesco Bodria, Fosca Giannotti, Riccardo Guidotti,
574 Francesca Naretto, Dino Pedreschi, and Salvatore
575 Rinzivillo. 2021. [Benchmarking and survey of ex-
576 planation methods for black box models](#). *CoRR*,
577 abs/2102.13076.

578 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
579 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
580 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
581 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
582 Gretchen Krueger, Tom Henighan, Rewon Child,
583 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
584 Clemens Winter, Christopher Hesse, Mark Chen, Eric
585 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
586 Jack Clark, Christopher Berner, Sam McCandlish,
587 Alec Radford, Ilya Sutskever, and Dario Amodei.
588 2020. [Language models are few-shot learners](#). In
589 *Proc. NeurIPS*.

590 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
591 Kristina Toutanova. 2019. [BERT: pre-training of
592 deep bidirectional transformers for language under-
593 standing](#). In *Proceedings of the 2019 Conference
594 of the North American Chapter of the Association
595 for Computational Linguistics: Human Language
596 Technologies, NAACL-HLT*, pages 4171–4186. Asso-
597 ciation for Computational Linguistics.

598 Christiane Fellbaum. 1998. *WordNet: An Electronic
599 Lexical Database*. Bradford Books.

600 Riccardo Guidotti. 2022. [Counterfactual explanations
601 and how to find them: literature review and bench-
602 marking](#). *Data Mining and Knowledge Discovery*.

603 Suchin Gururangan, Swabha Swayamdipta, Omer Levy,
604 Roy Schwartz, Samuel R. Bowman, and Noah A.
605 Smith. 2018. [Annotation artifacts in natural language
606 inference data](#). In *Proc. NAACL-HLT*. Association
607 for Computational Linguistics.

Peter Hase and Mohit Bansal. 2020. [Evaluating explain-
able AI: which algorithmic explanations help users
predict model behavior?](#) In *Proc. ACL*. Association
for Computational Linguistics. 608
609
610
611

Matthew Honnibal and Ines Montani. 2017. [spaCy 2:
Natural lanugage understanding with Bloom embed-
dings, convolutional neural networks and incremental
parsing](#). To appear. 612
613
614
615

Alon Jacovi. 2023. [Trends in explainable AI \(XAI\)
literature](#). *CoRR*, abs/2301.05433. 616
617

Ziqiang Li, Rentuo Tao, Jie Wang, Fu Li, Hongjing Niu,
Mingdao Yue, and Bin Li. 2021. [Interpreting the
latent space of gans via measuring decoupling](#). *IEEE
Trans. Artif. Intell.*, 2(1):58–70. 618
619
620
621

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta:
A robustly optimized BERT pretraining
approach](#). *CoRR*. 622
623
624
625
626

Nishtha Madaan, Inkit Padhi, Naveen Panwar, and Dip-
tikalyan Saha. 2021. [Generate your counterfactuals:
Towards controlled counterfactual generation for text](#).
In *Thirty-Fifth Conference on Artificial Intelligence,
AAAI, Conference on Innovative Applications of Ar-
tificial Intelligence, IAAI, The Symposium on Edu-
cational Advances in Artificial Intelligence, EAAI*.
AAAI Press. 627
628
629
630
631
632
633
634

David Martens and Foster J. Provost. 2014. [Explain-
ing data-driven document classifications](#). *MIS Q.*,
38(1):73–99. 635
636
637

Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right
for the wrong reasons: Diagnosing syntactic heuris-
tics in natural language inference](#). In *Proc. ACL*.
Association for Computational Linguistics. 638
639
640
641

Tim Miller. 2019. [Explanation in artificial intelligence:
Insights from the social sciences](#). *Artif. Intell.*, 267:1–
38. 642
643
644

John X. Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby,
Di Jin, and Yanjun Qi. 2020. [Textattack: A frame-
work for adversarial attacks, data augmentation, and
adversarial training in NLP](#). In *Proc. EMNLP*. Asso-
ciation for Computational Linguistics. 645
646
647
648
649

Judea Pearl. 2009. [Causal inference in statistics: An
overview](#). *Statistics Surveys*, 3(none):96 – 146. 650
651

Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert:
Sentence embeddings using siamese bert-networks](#).
In *Proc. EMNLP*. Association for Computational Lin-
guistics. 652
653
654
655

Marcel Robeer, Floris Bex, and Ad Feelders. 2021. [Gen-
erating realistic natural language counterfactuals](#). In
Findings EMNLP. Association for Computational
Linguistics. 656
657
658
659

- 660 Alexis Ross, Ana Marasovic, and Matthew E. Peters.
661 2021. [Explaining NLP models via minimal con-](#)
662 [trastive editing \(mice\)](#). In *Findings ACL/IJCNLP*.
663 Association for Computational Linguistics.
- 664 Alexis Ross, Tongshuang Wu, Hao Peng, Matthew E.
665 Peters, and Matt Gardner. 2022. [Tailor: Generating](#)
666 [and perturbing text with semantic controls](#). In *Proc.*
667 *ACL*. Association for Computational Linguistics.
- 668 Candida S. Punla and Rosemarie C. Farro. 2022. Are
669 we there yet?: An analysis of the competencies of
670 BEED graduates of BPSU-DC. *International Multi-*
671 *disciplinary Research Journal*, 4(3):50–59.
- 672 Victor Sanh, Lysandre Debut, Julien Chaumond, and
673 Thomas Wolf. 2019. Distilbert, a distilled version
674 of bert: smaller, faster, cheaper and lighter. *ArXiv*,
675 abs/1910.01108.
- 676 Yujun Shen, Jinjin Gu, Xiaou Tang, and Bolei Zhou.
677 2020. [Interpreting the latent space of gans for seman-](#)
678 [tic face editing](#). In *Proc. CVPR*. Computer Vision
679 Foundation / IEEE.
- 680 Sandra Wachter, Brent Mittelstadt, and Chris Russell.
681 2018. Counterfactual explanations without opening
682 the black box: Automated decisions and the GDPR.
683 *Harvard Journal of Law and Technology*, 31(2):841–
684 87.
- 685 Xiao-Yong Wei and Chong-Wah Ngo. 2007. [Ontology-](#)
686 [enriched semantic space for video search](#). In *Proc.*
687 *International Conference on Multimedia*. ACM.
- 688 Tongshuang Wu, Marco Túlio Ribeiro, Jeffrey Heer,
689 and Daniel S. Weld. 2021. [Polyjuice: Generating](#)
690 [counterfactuals for explaining, evaluating, and im-](#)
691 [proving models](#). In *Proc. ACL/IJCNLP*. Association
692 for Computational Linguistics.
- 693 Linyi Yang, Eoin M. Kenny, Tin Lok James Ng, Yi Yang,
694 Barry Smyth, and Ruihai Dong. 2020. [Generating](#)
695 [plausible counterfactual explanations for deep trans-](#)
696 [formers in financial text classification](#). In *Proc. COL-*
697 *ING*. International Committee on Computational Lin-
698 guistics.