
The CF Score: Evaluating Counterfactuals

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 ▶Write a proper abstract◀

2 **1 Introduction**

3 With the rising popularity of machine learning applications within the image domain, an urgent need
4 for understanding machine learning models arise. As a tool for understanding machine learning
5 models, many techniques for explaining such models have been proposed. In some cases, when fully
6 automated algorithms affects peoples lives, an explanation can even required by law [11]. To name
7 but a few techniques, some methods are based on generating heatmaps that identify salient input
8 features which drive model predictions [30, 5, 24, 6], others rely on simpler surrogate models that
9 are easier to explain [14], and some produce contrastive examples [38, 8]. In this work, we focus on
10 counterfactual explanations for the image domain.

11 Counterfactual explanations aims at identifying specific changes to specific inputs, such that the
12 predicted outcome of a machine learning model changes to a desired outcome. Such explanations
13 are of particular interest in situations where users of systems based on machine learning needs to
14 analyze the inner workings of the system. For example, a surveillance system picking out candidates
15 for screening could be assessed for biases.

16 Work on counterfactual explanations has become increasingly popular [26, 29, 39, 38, 13, 12, 22,
17 14, 37, 7, 28, 9, 8, 21, 18, 10, 17, 2, 40]. However, no one common interpretation of properties
18 that counterfactual explanations should possess, neither one standard metric exists for evaluating
19 counterfactual explanations. In this work, we consolidate the desiderata of image counterfactuals into
20 two categories. In terms of the two categories, we discuss existing metrics for image counterfactuals
21 and demonstrate experimentally how such metrics sometimes fail in capturing the intentions of the
22 metrics. We further propose two additional metrics, which contributes to a more complete profile of
23 counterfactual methods. While we demonstrate that existing metrics do have flaws, we argue that
24 reporting all metrics to give a more complete profile of counterfactuals is essential for the field of
25 counterfactual explanations to move forward properly.

26 ▶ Add that when data complexity increases, then the performance of the metrics decrease. ◀

27 **2 Counterfactual Examples**

28 In their original setting, counterfactual examples are closely connected to causality and structural
29 causal models. The counterfactual example relates to “what if?” questions.

30 **Definition 1 Causal counterfactual:** A counterfactual question is one where we seek information
31 about how the world would have been different given some specific change. For example; in a
32 situation where X did happened, what if X had not happened? ▶cite Causation book or somewhere
33 else for a definition.◀

34 With structural causal models, one can compute probabilities of events happening in a parallel
35 universe with an identical causal model and identical random variable outcomes, had variables (X)
36 been forced to take other values. However, in deep learning, we are typically lacking such causal
37 model, and a theoretical analysis becomes impossible. As an alternative, we come up with examples,
38 that resembles the original input to the model, but with specific changes that makes the model under
39 consideration predict a desired outcome. For example, if we have a model that predicts digit labels
40 from images, we might turn a four into a nine by connecting the two top-lines of the four.

41 Such an example has a close resemblance to Definition 1, but does not answer the question as such.
42 Instead, it answers a different type of counterfactual question.

43 **Definition 2 None-causal counterfactual:** A none-causal counterfactual question is one where we
44 seek what should have been changed to obtain a desired outcome. For example; input X yielded
45 outcome Y , how should we change X to obtain outcome Z .

46 An answer to this question typically comes in the form “Had values v_1, \dots, v_m been $\hat{v}_1, \dots, \hat{v}_m$ and
47 all other values remained the same, then outcome Y would have been Z .” Naturally, such answers
48 come in various forms and might not convey the information that we would expect. As such, much
49 related work has defined additional criteria that counterfactuals need to possess in order to be useful
50 for humans. We devote the next section to summarize such criteria.

51 3 The Good Counterfactual

52 For image counterfactual explanations to be useful for humans, both changing too little and too
53 much of the input may be undesirable. For example, much research has gone into adversarial attacks,
54 where images are changed slightly to change the predicted outcome, while not being perceptually
55 different from the input [34]. In relation to explainability, such attacks are undesirable as they
56 convey no information that humans will understand. On the contrary, changing too much could be
57 counterfactuals that maintains no information from the original input. Such example does not reveal
58 any actionable information either. In turn, related work describes various criteria to limit the output
59 space of valuable counterfactuals. In the following, we have categorized the most pronounced criteria.

60 **Minimal changes.** The first criterion is that changes to the input should be minimal [?]. If a
61 counterfactual example does not resemble the original input, it is of little or no use. In [16], the
62 minimality criterion is slightly adapted to state that only class-dependent features should change, *i.e.*,
63 features like connecting the two top-lines when changing a four to a nine can change, but features
64 like tilt and font weight should not.

65 Minimal changes are also important in the user context. A classic example is the bank loan application.
66 Imagine a customer who is not able to get a loan. If the counterfactual example says is that the
67 customer needs to get a pay raise, have larger savings, and ten other changes, then it may not be
68 actionable for the customer.

69 Especially for simpler datasets like tabular data, the amount of change can simply be counted [36, 13,
70 15]. However, as we will demonstrate in our experiments, simple counting and similarly manhattan
71 or euclidean distances may be useful but have drawbacks on high-dimensional data like images.
72 Dhurandhar et al. [8] also computes proxies for the optimal counterfactual features to be changed on
73 tabular data and quantifies how many of such features are identified by their method. Computing such
74 proxies are computationally heavy to identify and arguably of little relevance as pixels are inherently
75 dependent on each other.

76 **Realistic changes.** Another criterion is that the counterfactual examples should be likely to stem
77 from the same data distribution as the training data [27]. That is, the counterfactual examples
78 should be realistic or feasible as phrased by Poyiadzi et al. [27]. [25] describes how counterfactual
79 explanations should be faithful by both being in the proximity of the input and being connected to
80 samples of the counterfactual class. Similarly, [22] presents the causal edge score and [37] the *IM1*
81 score to compare densities. One can see such attempts of forcing changes to be realistic as an attempt
82 to guards against adversarial attacks and other artificial changes to the input.

83 For some domains, realistic changes also means that the counterfactual examples should obey, *e.g.*,
84 physical laws. For example, loan applicants cannot increase their educational level without increasing
85 their age [3]. This particular aspect is sometimes denoted both actionable [27] and feasible [?].

86 As for minimal changes, realistic changes are hard to quantify in the image domain. Is an adversarial
87 attack for example realistic? Although being important and useful for consolidating the field of
88 counterfactual explanations, we find in experiments that the metrics we studies all have pitfalls and
89 should be used with care.

90 **Computation time.** If counterfactual examples are to be used in a real world setting, it is also
91 important that they can be computed efficiently [36, 37, 16]. Assume that it takes 20 seconds every
92 time an explanation has to be generated. In such case, interaction with a system to find answers will
93 be very poor. According to [?], every process that takes more than 0.1 second is too slow to be used
94 in an interactive setting.

95 Other quantities that affect computation time that are more method specific, such as the number of
96 gradient steps needed [37] or number of patches swapped between the input and a sample from the
97 target class [12]. As such quantities cannot be used in a broad evaluation, we do not include them in
98 this work.

99 **Additional criteria.** Some methods are based on simpler surrogate models, which mimic the more
100 complex model under consideration. The surrogate models are used to explain the more complex
101 models. Although we do not consider such approaches in this work, we briefly include some additional
102 criteria here for completeness. It is obvious that such surrogate models should behave similarly to
103 how the complex models behaves. To quantify this behavior, scores like percent wise agreement,
104 accuracy, and F1-scores between surrogate and complex models are often used [14, 36, 26].

105 It is of course also important that humans can actually use the generated counterfactuals. In turn,
106 multiple works have been testing their methods for, *e.g.*, human teaching [12] or expert verification [8].
107 The tests were both domain and method specific and thus prohotit a generalized test. Therefor, we do
108 not include them for further evaluation here.

109 4 Evaluating Counterfactuals

110 As a basis for comparing methods in terms of minimal and realistic changes, multiple metrics have
111 been proposed. In this section, we present those metrics which have been applied to images and
112 discuss their applicability in terms of how they measure changes. We argue that the metrics are useful,
113 but results should be evaluated as a whole and not stand-alone.

114 4.1 Existing Metrics

115 **Simple distance metrics.** A natural first approach to measuring the amount of change between
116 inputs and counterfactuals is metrics like L_1 and L_2 -norms. However, as is broadly agreed in literature,
117 such norms work poorly on high-dimensional data like images [18]. We include it here because
118 both metrics (or variants there of) are used in objective functions for gradient based counterfactual
119 methods [38, 7, 37] and are thus also a natural first choice for metrics.

120 In the experiment section, we include a hybrid metric denoted the elastic net distance (EN), which is
121 defined as $EN(x, c) = \|x - c\|_1 + \beta\|x - c\|_2$, where x is the input and c is the counterfactual. We
122 use this metric because it has both the properties of the L_1 and the L_2 -norm.

123 **►I am not sure how much to write here. It is very simple norms that everyone knows, so I
124 don't want to state the obvious too much. The main thing thay may be missing is what the
125 pros and cons are in terms of minimal and realistic changes.◀**

126 **Target-class validity.** A central property that is sometimes reported is the target-class validity [22],
127 which tells the percentage of the generated counterfactuals that are predicted to be of the target class
128 by the classifier under consideration. The score tells how effective a given method is in generating
129 counterfactual examples.

130 **FakeMNIST.** Despite not being a metric as such, Hvilsted et al. [16] introduced an image dataset
 131 which allows evaluating whether only class dependent features are changed, when producing counterfactual
 132 examples. MNIST images [20] were shuffled and assigned new random labels in order
 133 to generate the dataset. The top-left 10×1 pixels were then colored according to the new labels
 134 (see first row of Figure 1a in the experiments for examples). In turn, the digits present in the images
 135 are independent of the labels while only the top-left pixels are label-dependent. If a counterfactual
 136 method operates properly, only the ten top left pixels should change, as the rest of the image contains
 137 no class relevant information.

138 In a similar spirit as how methods for producing saliency maps should be able to pass the sanity
 139 checks described in [1], we find it reasonable to argue the counterfactual methods should be able to
 140 produce counterfactual examples that only change the ten top left pixels on the FakeMNIST dataset.

141 ***IM1 score.*** One way to measure realistic changes is proposed by [37], which employs auto-
 142 encoders to quantify how well counterfactual examples follow the distribution of the training data.

143 The *IM1score* is based on variational auto-encoders trained on the same training data as the classifier
 144 to be explained. The score shows the ratio between how well the counterfactual example c of target
 145 class q can be reconstructed by an auto-encoder trained on data from the target class AE_q and an
 146 auto-encoder AE_p trained on the data of the input class p .

$$IM1(x, c) := \frac{\|c - AE_q(c)\|_2^2}{\|c - AE_p(c)\|_2^2 + \epsilon}. \quad (1)$$

147 A value below one means that c has less error when reconstructed with AE_q than with AE_p . According-
 148 ing to Van Looveren and Klaise [37], a lower value means that $\hat{x}^{(q)}$ follows the distribution of the
 149 target class better than the distribution of the input class.

150 As argued in the previous section, it is important to measure how well counterfactual examples
 151 follows the distribution of the training data. The *IM1score* is a valuable tool for assessing such
 152 properties. As the score is quantitative, it also allows comparing different methods across publications.
 153 Finally, the score is somewhat established as a metric, as multiple published (and unpublished) papers
 154 report the score [37, 22, 32].

155 There are, however, some concerns about the score as well. First, we note that the L_2 -norm can
 156 be deceiving when used on images. For example, we find in our experiments, that methods which
 157 make small changes to the input can get good scores, presumably because tiny changes that are not
 158 recovered yields only tiny additional errors in the score. Second, we have a concern that some classes
 159 may be easier to reconstruct than others, resulting in skewed scores for different inputs. For example,
 160 images of the digits one is arguably easier to reconstruct than, e.g., images of the eight digit. As
 161 such, converting a one to an eight may yield a higher score than the opposite way. Finally, we have
 162 not identified any publicly available pre-trained auto-encoders for computing the scores. When new
 163 auto-encoders need to be trained for every publication, results may vary and may not be comparable
 164 across publications. In our experiments, we support this final concern by demonstrating how, e.g.,
 165 data-normalization will yield incomparable scores.

166 ***IM2 [37].*** Van Looveren and Klaise [37] also propose the *IM2score* which utilize the discrepancy
 167 between reconstructions made by a class specific auto-encoder AE_q and an auto-encoder trained on
 168 the entire training set AE . The metric is defined as

$$IM2(c) := \frac{\|AE_q(c) - AE(c)\|_2^2}{\|c\|_1 + \epsilon}. \quad (2)$$

169 In Equation (2), the squared discrepancy norm is normalized by the L_1 -norm of the counterfactual.
 170 According to the authors, a low value of *IM2* indicates an interpretable counterfactual because the
 171 counterfactual follows the distribution of the target class as well as the distribution of the whole data
 172 set.

173 The quality of the score is however debatable. Mahajan et al. [22] argues that both *IM1* and *IM2* is
 174 better reported by displaying both the denominator and numerator of each score. Recently, Schut
 175 et al. [32] also demonstrated that the *IM2* score fails to identify out-of-sample images. In a similar

176 vein, we find in our experiments that when comparing sufficiently complex datasets, the *IM2* score
177 becomes statistically insignificant amongst three different methods.

178 ►There may be more papers but I didn't find they yet.◀

179 **Fréchet Inception Distance.** The Fréchet Inception Distance (FID) is a metric widely used for
180 evaluating generative models [23]. The metric compares how much two datasets align by comparing
181 statistics of embeddings of the Inception V3 network [35]. In terms of image counterfactuals, the
182 score has been used to evaluate how well counterfactuals align with the original dataset [31, 33]. The
183 FID is defined from mean Inception embeddings μ_1 and μ_2 , and covariance matrices Σ_1 and Σ_2 of
184 the embeddings:

$$\text{FID} = \|\mu_1 - \mu_2\|_2^2 + \text{tr}[\Sigma_1 + \Sigma_2] - 2\text{tr}\left[\sqrt{\Sigma_1 \Sigma_2}\right] \quad (3)$$

185 In this work, we consider images that are smaller (64×64) compared to inputs of the Inception V3
186 model (299×299). As a consequence, we compute the same score, but for a different network. In
187 particular, we use the embeddings from the *AE* used in the *IM2* score and denote it the Fréchet
188 auto-encoder distance (FAED). Although the score is highly dependent on the embedding network,
189 we believe that our experimental results will also extend to the inception network.

190 The score is well suited for evaluating whether generated counterfactuals follow the distribution of
191 the training data. However, we find in experiments, that methods generating tiny changes to the input
192 may be deemed of high quality because tiny changes are either filtered out by the auto-encoder or
193 does not affect summarizing statistics.

194 4.2 New Proposed Metrics.

195 **Oracle score.** Through our experiments, we find that tiny changes similar to adversarial attacks
196 often yield undesirable but good scores when evaluated with the presented metrics. To mitigate this
197 issue, we introduce the oracle score, which is similar to the target-class validity but evaluated on an
198 external classifier. The objective is to make a score which is less vulnerable to adversarial attacks.

199 The score is based on the intuition that adversarial attacks are highly specific to a particular classifier
200 ►make sure this is the case◀. On the other hand, proper counterfactuals that change only class
201 relevant features should generalize to other classifiers. The oracle score is based on training an
202 additional classifier – the oracle – which is used to classify the counterfactuals. The score is then the
203 percentage of counterfactuals that are predicted to be of the target class. ►Should we discuss pros
204 and cons further here?◀

205 **Independent feature comparison.** For multi-label datasets, we can compute counterfactual exam-
206 ples for one label and use classifiers trained on the other labels to evaluate how much other classes
207 change for the counterfactual examples. Our proposed method is similar to how Singla et al. [33]
208 evaluate how a multi-class classifier changes the predicted outcome on the counterfactuals. We
209 propose to use the Jensen-Shannon (JS) divergence to monitor predicted outcome of different classes.
210 Over an entire test set, we compute counterfactuals for a particular label, and then we report average
211 JS divergences between predictions on inputs and counterfactuals.

212 Intuitively, non-related labels should remain unchanged, while labels that correlate may co-vary. For
213 example, if we change an image of a face without makeup to a face with makeup, we would expect
214 the face to be smiling as much as the input, but the outcome of predicting whether the face is wearing
215 lipstick should roughly follow the makeup label as lipstick is arguably a subset of wearing makeup.
216 Experimentally, we find that more realistically looking counterfactuals behave more realistically than,
217 e.g., examples with tiny adversarial-like changes.

218 One particular example of a dataset, where our metric is applicable is the FakeMNIST dataset. On
219 that dataset, we can train a classifier to predict the new “fake” labels and another to predict the original
220 labels. This way, we can evaluate whether counterfactual examples affect only the “fake” labels. In
221 our experiments, our proposed metric yields scores that align well with human interpretation on two
222 different datasets.

223 **5 Experiments**

224 In this section, we demonstrate how existing metrics need to accompany each other, when counterfac-
225 tual methods are proposed. In particular, we demonstrate through multiple experiments that present
226 metrics have vulnerabilities and should not be used in isolation.

227 **Mehtds compared.** Throughout the experimental section, we compare three different algorithms
228 for producing counterfactual explanations. The three methods were chosen to represent a spectrum of
229 methods ranging from simple gradient based methods in one extreme to methods based purely on
230 generative models in the other extreme.

231 In one end of the spectrum, Wachter et al. [38] presents a gradient based method (denoted GB below).
232 The method generates counterfactuals by minimizing a loss-function consisting of two terms, a
233 prediction loss and a distance loss:

$$c = \arg \min_{c' \in \mathcal{X}} \lambda(f(c') - q)^2 + d(x, c'). \quad (4)$$

234 In Equation (4), f is the predictive model, d is a distance metric. For d , we use the L_1 -norm, as
235 suggested by the authors. However, we do not normalize by the median absolute deviation of each
236 input feature, as this makes little sense for images which have identical value ranges. The first term
237 in the minimization makes the predicted class change, while the second term encourages small and
238 sparse changes to the input. Another method that lies in this end of the spectrum is, e.g., [8] which
239 follows a similar structure as Equation (4), but with a more complex distance function. It should be
240 mentioned that both methods were originally introduced for tabular data. As we seek to investigate
241 properties of metrics and not compare performances of methods, we see no issue in applying such a
242 method to the image domain.

243 At the other end of the spectrum, we include ECINN [16] as a representative for the methods based on
244 generative models (denoted GEN below). The method is based on invertible neural networks (INNs)
245 which is a particular type of deep generative models that can be altered to also do classification [4].
246 ECINN makes corrections to embeddings of inputs. Counterfactual examples are successively
247 generated by inverting the embeddings with the INN. We find ECINN to be the most extreme case of
248 this end of the spectrum, compared to, e.g., [31, 33], because it is the same neural network which
249 is used for both predictions and for generating counterfactuals. In contrast, [31] and [33] trains a
250 surrogate auto-encoder and generative adversarial network, respectively, which is used for sampling
251 counterfactuals.

252 At the middle of the spectrum, methods exist which use gradients to compute counterfactuals similar
253 to that of Equation (4), but where the gradient optimization is guided by derivatives of generative
254 models or other more sophisticated loss terms [7, 37]. In our experiments, we use the method proposed
255 by [37] as representative, which uses embeddings from an auto-encoder to define a class-prototype
256 loss. We denote this group guided loss (GL) below.

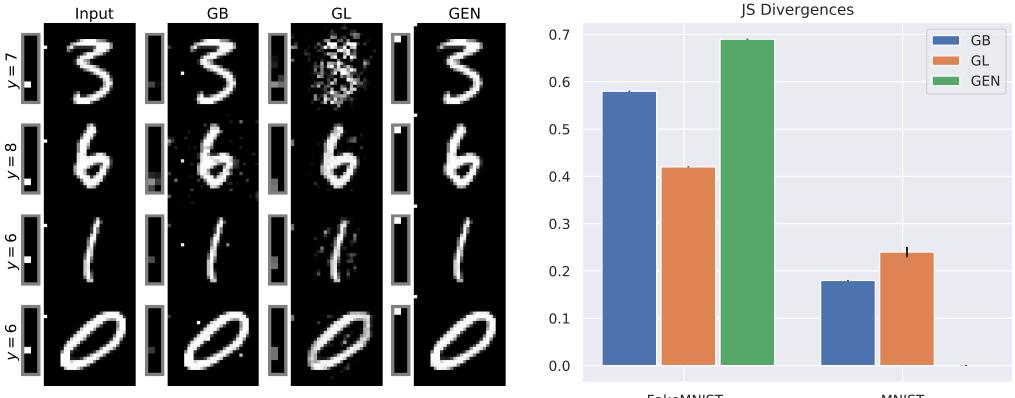
257 The algorithms from [38] and [37] were implemented using the `alibi` framework,¹ while ECINN
258 was adopted from the official code repository². The former two algorithms are used to identify
259 counterfactuals for the same “vanilla” convolutional neural network (CNN), identical to the one
260 described in [37]. ECINN is based on inherent properties of INNs and is therefor used with an INN as
261 predictive model. In turn, the presented results may be contributed to differences in architectures and
262 not algorithms as such. However, the goal of the experiments is not to identify a superior algorithm,
263 but to demonstrate properties of existing metrics for evaluating counterfactual explanations on images.

264 **Experimental setup.**

- 265 • Scores reported on test set
266 • Compute only on successful counterfactuals.
267 • Report 95% CI where appropriate.
268 • Independent feature comparison all with "Independent models", even the counterfactual
269 label.

¹<https://docs.seldon.io/projects/alibi> (v.0.5.9)

²MAKEURLHERE



(a) Counterfactual examples with target class $q = 0$.

(b) Mean JS-divergences.

Figure 1: FakeMNIST results. In (a), rectangular images left of larger images are the top left 10×2 pixels enlarged for improved readability.

270 For a fair comparison, the remaining metrics are only computed on the successful counterfactuals.
 271 Where appropriate, scores are reported with 95% confidence intervals in parentheses.

272 5.1 FakeMNIST

273 [16] proposed FakeMNIST; an artificial dataset which dictates the relationship between pixels and
 274 labels. The dataset can be used as a sanity check for whether counterfactual methods change only
 275 class relevant features. Originally, the dataset was introduced for a qualitative assessment of methods.
 276 As a result, no quantitative metric has been associated with the dataset. We apply our independent
 277 feature comparison to further the evaluation protocol for the dataset. In this experiment, we test how
 278 well quantitative and qualitative evaluations of counterfactuals align by quantifying whether only
 279 label dependent features of the input are changed.

280 Figure 1a displays four samples from the FakeMNIST testset. Smaller rectangles represent the first
 281 10×2 pixels enlarged for readability. The first column displays inputs with labels 7, 8, 6, and
 282 6, respectively (cf. labels or dot location in smaller rectangles). The following three columns are
 283 counterfactuals with target class $q = 0$, generated by the three paradigms we compare – gradient
 284 based (GB), guided loss (GL), and generative models (GEN).

285 As described, we would expect only the top-left 10×1 pixels to change. By inspecting Figure 1a, we
 286 observe three things. First, GB and GL both produce counterfactuals that do not seem to follow the
 287 distribution of the dataset. Specifically, the magnified parts either have multiple pixels colored or
 288 most color removed. Second, the gradient based approach seems to produce adversarial attacks, in
 289 which few lonely pixels are colored. Third, counterfactuals from GEN correspond well to what we
 290 would expect of a counterfactual.

291 We apply the independent feature comparison for counterfactuals generated for the entire testset. In
 292 Figure 1, we show how JS-divergences behave on the FakeMNIST labels (left) and on the original
 293 MNIST labels (right). Well aligned with what we see quantitatively, GEN is doing the best job in
 294 changing the predicted class on the FakeMNIST labels (high JS-divergence) and has almost no effect
 295 on the predicted MNIST labels. GB also has a high JS-divergence on the FakeMNIST labels but does
 296 also affect the predicted MNIST labels with a JS-divergence of just below 0.2. GL performs the worst
 297 by having the lowest FakeMNIST JS-divergence and the highest JS-divergence on MNIST. These
 298 observations align well with the intuition we get from Figure 1a, where GEN behaves as expected,
 299 GB behaves less sensible but with only small modifications to the input, and GL produces more blurry
 300 and artificially looking counterfactuals. In Table 3 in the supplementary material, we also include
 301 scores of all other metrics described in Section 4. On all the scores the results are similar to what we
 302 observe in this experiment, i.e., GEN has the best scores, GL the worst, and GB is in between. In
 303 conclusion, we find that for this simple dataset, qualitative observations and quantitative evaluations
 304 are well aligned.

Table 1: Scores on MNIST for counterfactuals with different normalizations.

Method	<i>IM1</i>	$100 \cdot IM2$	<i>EN</i>	Oracle	FAED
[−0.5; 0.5] normalization					
GB	1.04 (0.00)	0.96 (0.01)	16.04 (0.18)	73.15%	0.19
GL	1.04 (0.01)	0.89 (0.01)	42.72 (0.31)	37.64%	1.12
GEN	0.77 (0.00)	0.21 (0.00)	108.25 (0.59)	97.78%	7.12
[0; 1] normalization					
GB	1.06 (0.00)	2.45 (0.00)	16.04 (0.18)	81.18%	0.41
GL	1.04 (0.00)	1.94 (0.00)	42.72 (0.31)	45.13%	2.12
GEN	0.67 (0.00)	0.95 (0.00)	108.25 (0.59)	96.16%	14.34
[−1; 1] normalization					
GB	1.05 (0.00)	1.87 (0.02)	32.09 (0.36)	67.34%	1.84
GL	0.99 (0.00)	1.73 (0.01)	85.44 (0.62)	38.43%	9.97
GEN	0.86 (0.00)	0.47 (0.01)	216.49 (1.17)	96.61%	68.49

305 5.2 Normalization and Models Matters

306 Most metrics presented in this paper depend on data normalization or particular predictive models.
 307 The dependence makes reporting both data normalization and model specifications crucial for
 308 reproducibility. In this section, we demonstrate the mentioned issue with a practical example.

309 For different ranges of normalization, we train models used for the various evaluation metrics,
 310 *i.e.*, auto-encoders for IM1, IM2, and FAED; and a new classifier for the oracle score. For each
 311 counterfactual method, we use the same counterfactuals, but normalized differently, and evaluate
 312 each set with the metrics specified for the corresponding normalization.

313 In Table 1, we report mean scores for three different normalizations. Parenthesis indicate 95%
 314 confidence intervals. By pattern matching the bold numbers between the tables, we see that the
 315 best performing method would be concluded to be the same, independent of the normalization. In
 316 Figure 6 in the supplementary material, we even find the orders to be statistically significant across
 317 10 independently initialized models. It should be noted that the *EN* score is invariant to data shifts
 318 and scales linearly with scaling of the data range, which is what we also see in the table.

319 As the table indicates, there is, however, an issue. Had the *IM2* metric been used to compare GL with
 320 a [−0.5, 0.5] normalization against GEN with a [0, 1] normalization, the conclusion would have been
 321 wrong, as GL would have been looking better than GEN. Although this issue may seem obvious, it
 322 occurs in literature. If one compares reported *IM2* scores between [37] and [22], the difference is
 323 about an order of magnitude. [37] normalizes values to the range [−0.5, 0.5], while [22] normalizes
 324 values to the range [0, 1]. We believe that the normalization differences contribute to explaining the
 325 difference between the reported scores.

326 5.3 Inspecting Scores

327 To get a deeper insight into how different metrics behave, we have identified cases of input and
 328 counterfactual pairs, where the differences in scores are extreme. The extreme cases demonstrate
 329 situations, where the metrics may need to be used with caution.

330 **L1-norm.** For the image domain, the simple norms including the *EN* distance are known to work
 331 poorly in terms of interpretability [18]. For completeness, we give an example in figure Figure 2a.
 332 The figure shows a seven to the left and two successful counterfactuals with target class $q = 9$ (center
 333 and right). The *EN* distance is displayed above the two counterfactuals. Arguably, the center image
 334 still looks most like a seven and the right image convincingly looks like a nine. However, according
 335 to the *EN* distance, the center image is an order of magnitude better than the right. The example
 336 illustrates how tiny adversarial attacks may be deemed better than proper adversarial examples, simply
 337 because they change the input much less.

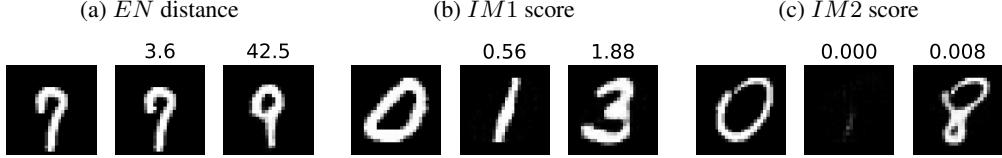


Figure 2: Each figure contains three samples. Left: input image, center: a counterfactual with a good score, right: a counterfactual with bad score.

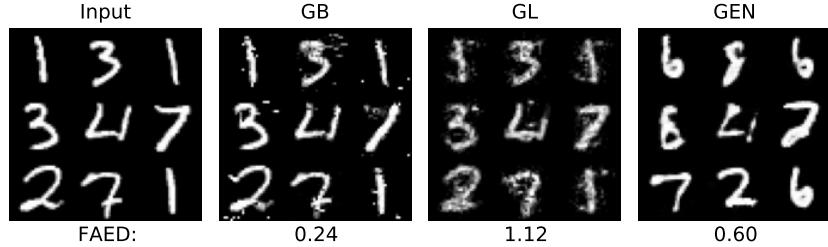


Figure 3: Counterfactual samples and the test-set-wide FAED.

338 **IM1.** In Figure 2b, we depict a zero with two different counterfactuals of target class $q = 1$ and
 339 $q = 3$, respectively. Visually, the three in the right image resembles the input better than the one
 340 in the center image, as it preserves most of the shape, tilt and font weight of the zero. The one
 341 seems more distant as the shape of the digit is much different from the zero. The center image does,
 342 however, get a three times better $IM1$ score than the three, *i.e.*, the one was better reconstructed by
 343 the auto-encoders than the three. Holding all else equal, this is presumably because more white pixels
 344 yields a larger loss.

345 [F: The mean reconstruction loss for AE_0 on class zero digits is 1.28, while it is 0.61 for AE_1 on class
 346 one inputs, and 1.45 for AE_3 on threes.]

347 **IM2.** For the IM2 metric, we show an extreme case in Figure 2c. Figure 2c shows how an all black
 348 image (center) yields a better (lower) score than a regular digit of an eight. Again, the argument is
 349 presumably that there is less to reconstruct for the auto-encoders in the center image than the right
 350 image. This observation is confirmed in the appendix, where we compute $IM2$ scores over samples
 351 where we gradually turn white pixel into black. In turn, the particular metric might wrongly give
 352 good scores for methods that simply removes information from the inputs.

353 **FID.** For Frechét distances like FID, it is not possible to identify single extreme cases because it is a
 354 population wide metric. Instead we compare scores between counterfactuals generated by GB which
 355 makes small local changes, GL which makes more blurry changes, and GEN which generally makes
 356 more global changes (cf. Figure 3 and the appendix for more samples).

357 In Figure 3, we show counterfactual examples and the test-set-wide FAED scores of ten thousand
 358 counterfactuals generated by the three methods, compared to the MNIST test set. The figure shows
 359 how GB has the smallest (best) FAED score, GL the highest (worst), and GEN is in between. From a
 360 human perspective, when the goal is to change the class of the input, both GB and GL does a poor
 361 job in generating realistic changes, as it is harder for humans to identify what the target class is in
 362 most cases. The two extremes are deemed good (bad) according to the FAED metric, because they
 363 make small (large) changes to the input. GEN also makes large changes to the input, but the changes
 364 better resembles the true data distribution and thus gets a lower score than GL. In turn, smaller values
 365 of FAED may be better, but very small changes similar to adversarial attacks may wrongly yield a
 366 low score.

367 We see that averages work well but in isolated cases, which may be important for evaluating methods
 368 on specific datasets or cases with specific properties....

369 **► I find it hard to conclude this subsection in a positive way. We basically find that metrics
 370 have different flaws in extreme cases (vulnerable to adversarial attacks, simpler classes, and**



Figure 4: Counterfactuals for the CelebA dataset.

Table 2: CelebA-HQ scores

Method	<i>TCV</i>	<i>IM1</i>	$100 \cdot IM2$	<i>EN</i>	Oracle	<i>FAED</i>
GB	96.07% (0.72)	0.98 (0.00)	0.48 (0.01)	142.21 (2.08)	0.83 (0.01)	0.70
GL	81.09% (1.44)	0.99 (0.00)	0.53 (0.01)	330.14 (18.24)	0.33 (0.02)	4.61
GEN	99.26% (0.32)	1.03 (0.00)	0.53 (0.01)	684.26 (11.86)	0.90 (0.01)	0.63

371 **removing information). We could do the discussion on hardness of small but not too small**
 372 **changes here. But what is the positive message? ◀**

373 5.4 Complex data

374 In this section, we scale our experiments to the larger and more complex dataset, CelebA-HQ [19].
 375 The goal is to compare qualitative and quantitative assessments of the counterfactual examples
 376 generated in a more complex setting.

377 CelebA-HQ is a dataset of faces, where each sample is associated with 40 different labels. Figure 4
 378 presents ten different inputs in the first row. The first five have a positive makeup label and the last five
 379 have a negative label. The following three rows represent counterfactual examples. Similar to what
 380 was observed in former experiments, GB generates counterfactual examples which represents tiny
 381 local visible changes in the form of few pixels being changed to extreme color values. GL generates
 382 more noisy and artificial changes. In most cases, GEN produces counterfactual examples that look
 383 more natural. However, we also observe artificially looking changes like the one in the lower
 384 right corner. Arguably, from a human perspective, the last row of counterfactuals seems to contain
 385 counterfactuals of the highest quality.

386 Table 2 presents scores computed on the 2824 test samples. Slightly in contradiction with the
 387 qualitative evaluation above, the table indicates that as an overall picture, GB and GEN are the best
 388 performing, as they each performs best on three metrics. GL does not perform best on any metric.
 389 Except from the *IM1* and *IM2* metric, the scores are well separated and taking a closer look at the
 390 scores, the *IM1* and *IM2* scores have similar values for all three methods. As such, it seems the
 391 only little can be concluded from the particular value.

Comparing Figure 4 and Table 2 jointly reveals how on this dataset, the *IM1* and *IM2*
 392 support our observation that the qualitative and quantitative evaluations support each other poorly. As
 393 such, there is an apparent need for a suit of metrics that follows human intuition better.

394 6 Future Work

- 395 • Alignment of desiderata and metrics for images.

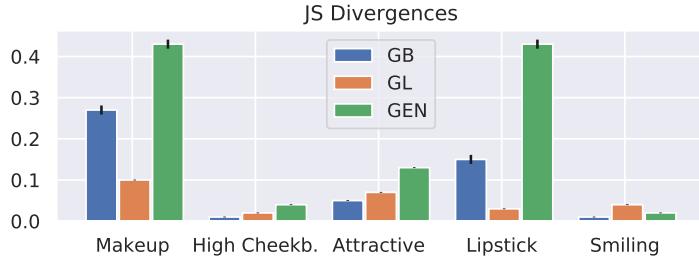


Figure 5: Average JS-divergences between predictions of originals and counterfactual examples for varying labels. Black vertical bars indicate 95% confidence intervals.

- 396 • Examples:
 - 397 – Leave class-independent features alone.
 - 398 – Better distance metrics that reflects our intuition.
- 399 • Interventional data.
- 400 • A stronger tie to causal models.

401 In [16], they describe how it is important that only class-dependent features are changed. The CelebA-
 402 HQ data lends itself useful to investigate this property. We have trained five independent classifiers to
 403 predict five different labels. We then use the Jensen-Shannon (JS) divergence to compare average
 404 predicted outcomes of the five classifiers, when comparing the original input to counterfactuals
 405 generated. For each method, we use each classifier to predict the class of the original inputs and
 406 the generated counterfactuals and compute the average JS divergences across the test set. Figure 5
 407 displays the results. The counterfactuals were generated for the makeup label which is why the
 408 divergence should be high on that label, *i.e.*, the predicted class should have changed. On the
 409 remaining classes like, *e.g.*, smiling, the divergence should stay low, because the smiling property
 410 should not be changed when adding or removing makeup. According to the table, the divergences
 411 does in deed mostly remain low for the additional classes. For ECINN and partly Wachter, the
 412 divergence of the libstick label does increase with the counterfactuals. This is presumably because
 413 libstick is highly correlated with wearing makeup in general. A similar argument can be made for
 414 being attractive. Figure 5 reveals how ECINN and to some extend Wachter generates counterfactuals
 415 that maintain such relationship, while this is generally not the case for Van Looveren.

416 **References**

- 417 [1] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity
418 checks for saliency maps. In *NeurIPS*, 2018.
- 419 [2] Arjun Akula, Shuai Wang, and Song-Chun Zhu. CoCoX: Generating Conceptual and Counterfactual
420 Explanations via Fault-Lines. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- 421 [3] Holger Andreas and Lorenzo Casini. Hypothetical interventions and belief changes. *Foundations of
422 Science*, 24(4):681–704, 2019.
- 423 [4] Lynton Ardizzone, Radek Mackowiak, Carsten Rother, and Ullrich Köthe. Training normalizing flows
424 with the information bottleneck for competitive generative classification. *NeurIPS*, 2020.
- 425 [5] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus Robert Müller, and
426 Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance
427 propagation. *PLoS ONE*, 2015.
- 428 [6] Chun Hao Chang, Elliot Creager, Anna Goldenberg, and David Duvenaud. Explaining image classifiers by
429 counterfactual generation. In *ICLR*, 2019.
- 430 [7] Amit Dhurandhar, Pin Yu Chen, Ronny Luss, Chun Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and
431 Payel Das. Explanations based on the Missing: Towards Contrastive Explanations with Pertinent Negatives.
432 In *NeurIPS*, 2018.
- 433 [8] Amit Dhurandhar, Tejaswini Pedapati, Avinash Balakrishnan, Pin Yu Chen, Karthikeyan Shanmugam, and
434 Ruchir Puri. Model agnostic contrastive explanations for structured data. *arXiv*, 2019.
- 435 [9] Nadine Elzein. The demand for contrastive explanations. *Philosophical Studies*, 176(5):1325–1339, 2019.
436 URL <https://doi.org/10.1007/s11098-018-1065-z>.
- 437 [10] Oscar Gomez, Steffen Holter, Jun Yuan, and Enrico Bertini. ViCE: Visual Counterfactual Explanations
438 for Machine Learning Models. *International Conference on Intelligent User Interfaces, Proceedings IUI*,
439 pages 531–535, 2020.
- 440 [11] Bryce Goodman and Seth Flaxman. European union regulations on algorithmic decision making and a
441 “right to explanation”. *AI Magazine*, 38(3):50–57, 2017.
- 442 [12] Yash Goyal, Ziyan Wu, Jan Ernst, Dhruv Batra, Devi Parikh, and Stefan Lee. Counterfactual Visual
443 Explanations. In *ICML*, 2019.
- 444 [13] Rory Mc Grath, Luca Costabello, Chan Le Van, Paul Sweeney, Farbod Kamiab, Zhao Shen, and Freddy
445 Lecue. Interpretable Credit Application Predictions With Counterfactual Explanations. *arXiv preprint
arXiv:1811.05245*, 2018.
- 447 [14] Riccardo Guidotti, Anna Monreale, Fosca Giannotti, Dino Pedreschi, Salvatore Ruggieri, and Franco
448 Turini. Factual and Counterfactual Explanations for Black Box Decision Making. *IEEE Intelligent Systems*,
449 34:14–23, 2019.
- 450 [15] Masoud Hashemi and Ali Fathi. PermuteAttack: Counterfactual Explanation of Machine Learning Credit
451 Scorecards. (August), 2020. URL <http://arxiv.org/abs/2008.10138>.
- 452 [16] Frederik Hvilshøj, Alexandros Iosifidis, and Ira Assent. ECINN: efficient counterfactuals from invertible
453 neural networks. *arXiv preprint arXiv:2103.13701*, 2021.
- 454 [17] Alon Jacovi, Swabha Swamyamdipta, Shauli Ravfogel, Yanai Elazar, Yejin Choi, and Yoav Goldberg.
455 Contrastive Explanations for Model Interpretability. *arXiv preprint arXiv:2103.01378*, 2021.
- 456 [18] Sin Han Kang, Hong Gyu Jung, Dong Ok Won, and Seong Whan Lee. Counterfactual explanation based
457 on gradual construction for deep networks. *arXiv*, pages 1–26, 2020.
- 458 [19] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive Growing of GANs for Improved
459 Quality, Stability, and Variation. In *ICLR*, 2018.
- 460 [20] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- 462 [21] Weixin Liang, James Zou, and Zhou Yu. ALICE: Active Learning with Contrastive Natural Language
463 Explanations. pages 4380–4391, 2020. ISSN 23318422.

- 464 [22] Divyat Mahajan, Chenhao Tan, and Amit Sharma. Preserving causal constraints in counterfactual explanations
 465 for machine learning classifiers. *arXiv*, 2019.
- 466 [23] Alexander Mathiasen and Frederik Hvilshøj. Backpropagating through fr\'echet inception distance. *arXiv preprint arXiv:2009.14075*, 2020.
- 468 [24] Gr\'egoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus Robert M\"uller.
 469 Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognition*, 2017.
- 470 [25] Martin Pawelczyk, Klaus Broelemann, and Gjergji Kasneci. Learning Model-Agnostic Counterfactual
 471 Explanations for Tabular Data. *The Web Conference 2020 - Proceedings of the World Wide Web Conference, WWW 2020*, pages 3126–3132, 2020.
- 473 [26] Tejaswini Pedapati, Avinash Balakrishnan, Karthikeyan Shanmugam, and Amit Dhurandhar. Learning
 474 global transparent models from local contrastive explanations. *arXiv*, (NeurIPS), 2020.
- 475 [27] Rafael Poyiadzi, Kacper Sokol, Raul Santos-Rodriguez, Tijl De Bie, and Peter Flach. FACE: Feasible and
 476 actionable counterfactual explanations. *AIES 2020 - Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2020. doi: 10.1145/3375627.3375850.
- 478 [28] Mohit Prabhushankar, Gukyeong Kwon, Dogancan Temel, and Ghassan AlReigib. Contrastive Explanations
 479 in Neural Networks. In *ICIP*, 2020.
- 480 [29] Shubham Rathi. Generating Counterfactual and Contrastive Explanations using SHAP. *arXiv*, 2019.
- 481 [30] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the
 482 predictions of any classifier. In *KDD*, 2016.
- 483 [31] Pau Rodriguez, Massimo Caccia, Alexandre Lacoste, Lee Zamparo, Issam Laradji, Laurent Charlin, and
 484 David Vazquez. Beyond Trivial Counterfactual Explanations with Diverse Valuable Explanations. 2021.
 485 URL <http://arxiv.org/abs/2103.10226>.
- 486 [32] Lisa Schut, Oscar Key, Rory McGrath, Luca Costabello, Bogdan Sacaleanu, Medb Corcoran, and Yarin
 487 Gal. Generating Interpretable Counterfactual Explanations By Implicit Minimisation of Epistemic and
 488 Aleatoric Uncertainties. 130, 2021. URL <http://arxiv.org/abs/2103.08951>.
- 489 [33] Sumedha Singla, Brian Pollack, Junxiang Chen, and Kayhan Batmanghelich. Explanation by Progressive
 490 Exaggeration. pages 1–20, 2019.
- 491 [34] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural
 492 networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, 2019.
- 493 [35] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the
 494 inception architecture for computer vision. In *CVPR*, 2016.
- 495 [36] Jasper van der Waa, Marcel Robeer, Jurriaan van Diggelen, Matthieu Brinkhuis, and Mark Neerincx.
 496 Contrastive Explanations with Local Foil Trees. 2018. URL <http://arxiv.org/abs/1806.07470>.
- 497 [37] Arnaud Van Looveren and Janis Klaise. Interpretable counterfactual explanations guided by prototypes.
 498 *arXiv preprint arXiv:1907.02584*, 2019.
- 499 [38] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the
 500 black box: Automated decisions and the GDPR. *Harv. JL & Tech.*, 31:841, 2017.
- 501 [39] Pei Wang and Nuno Vasconcelos. SCOUT: Self-aware Discriminant Counterfactual Explanations. In
 502 *CVPR*, 2020.
- 503 [40] Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel S. Weld. Polyjuice: Automated, General-
 504 purpose Counterfactual Generation. *arXiv preprint arXiv:2101.00288*, 2021.

Table 3: Test set wide mean (confidence intervals) on the FakeMNIST dataset. Best scores are reported in bold.

Method	$IM1$	$100 \cdot IM2$	TCV	EN	Oracle	$FAED$
GB	1.22 (0.01)	0.50 (0.01)	0.68 (0.01)	11.56 (0.49)	0.88 (0.01)	0.52
GL	1.02 (0.00)	1.23 (0.03)	0.84 (0.01)	47.34 (0.91)	0.56 (0.01)	2.44
GEN	0.70 (0.00)	0.22 (0.00)	1.00 (0.00)	6.72 (0.04)	1.00 (0.00)	0.15

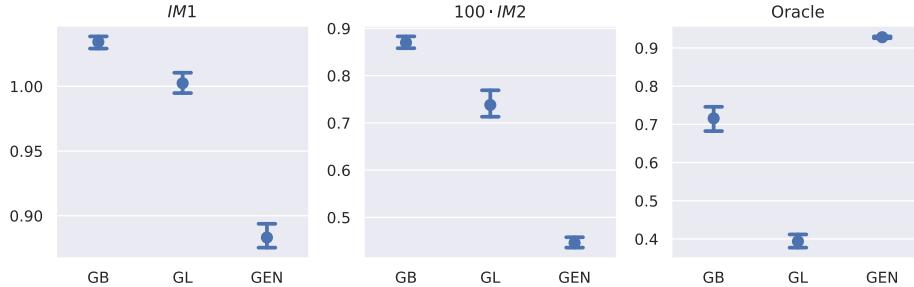


Figure 6: Mean scores on MNIST. Confidence intervals for ten trials with ten randomly initialized evaluation models.

505 A Additional Experimental Results

506 A.1 FakeMNIST

507 ▶ Describe ◀

508 A.2 MNIST

509 ▶ Describe ◀

510 A.3 Additional Samples

511 FakeMNIST

512 MNIST

513 CELEBA

514 A.4 Extremes

515 L1

516 L2

517 EN

518 **IM1** AE distances might be misleading:

519

520 To test this interpretation, for each MNIST class, we computed average reconstruction errors for each
 521 of the ten autoencoders AE_i , $\ell_i = \frac{1}{n} \sum_x \|x - AE_i(x)\|^2$. Figure 16 displays the average per-row
 522 relationship between input classes (rows) and autoencoders (columns). If the interpretation was true,
 523 we would expect the samples from class i to have the lowest reconstruction loss of the corresponding
 524 autoencoder, ℓ_i , thus making the diagonal in Figure 16 dark. However, inspecting Figure 16 reveals



(a) Test set inputs.



(b) GB [38].



(c) Prototype [37].



(d) INN [16].

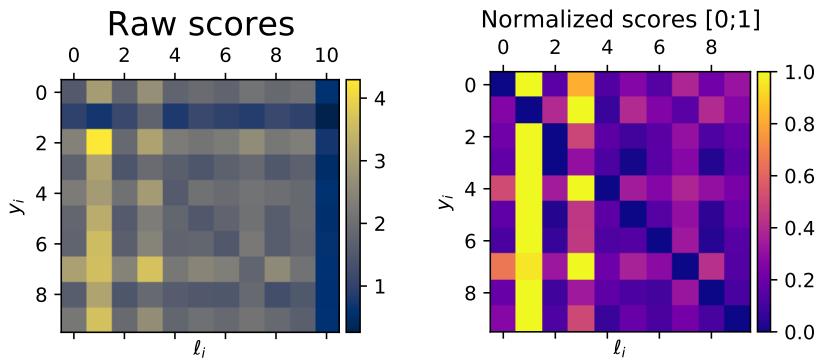
Figure 7: FakeMNIST [16] samples. Gray samples are those where no solution were found.

525 that on average, threes are reconstructed better with, e.g., AE_2 , AE_5 , or AE_8 than with AE_3 . In turn,
526 the interpretation above does not seem to hold in practice.

527 Additionally, Figure 16 shows that AE_1 is bad at reconstructing anything but ones, judging by
528 the mostly yellow second column. As such, when $\|x - AE_1(x)\|_2^2$ becomes the denominator in
529 Equation (1), then the score will look better than with the other autoencoders. Simply because AE_1
530 cannot reconstruct, e.g., an eight, as well as either of the other autoencoders.



Figure 8: MNIST test set samples.



531

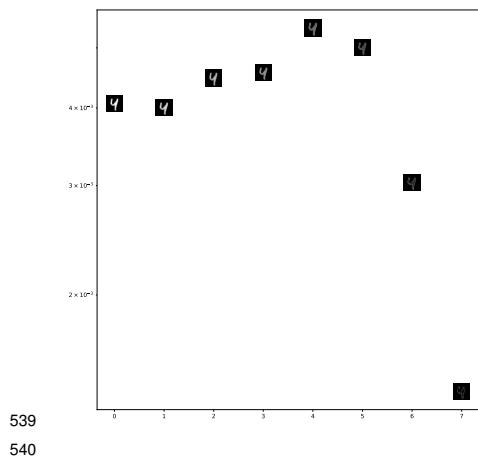
532 **IM2** Graying out gives lower score:

533

534 In a similar vein as how we compared scores of *im1* in figure 2b, we have displayed an input and two
 535 counterfactuals with low (good) respectively high (bad) *im2* score in figure 2c. The figure shows
 536 how removing almost all content from the input yields good *IM2* score while a more reasonable
 537 counterfactual where the input zero with an “open top-line” is translated into an eight with a similar
 538 “open top-line” gets a 15 times worse score.



Figure 9: MNIST counterfactuals generated by algorithm from [38].



Van Looveren



Figure 10: MNIST counterfactuals generated by algorithm from [37]



Figure 11: MNIST counterfactuals generated by algorithm from [16]

Input



Figure 12: CelebA-HQ test set samples.

Wachter



Figure 13: CelebA-HQ counterfactuals generated by algorithm from [38].

Van Looveren



Figure 14: CelebA-HQ counterfactuals generated by algorithm from [37]

ECINN



Figure 15: CelebA-HQ counterfactuals generated by algorithm from [16]

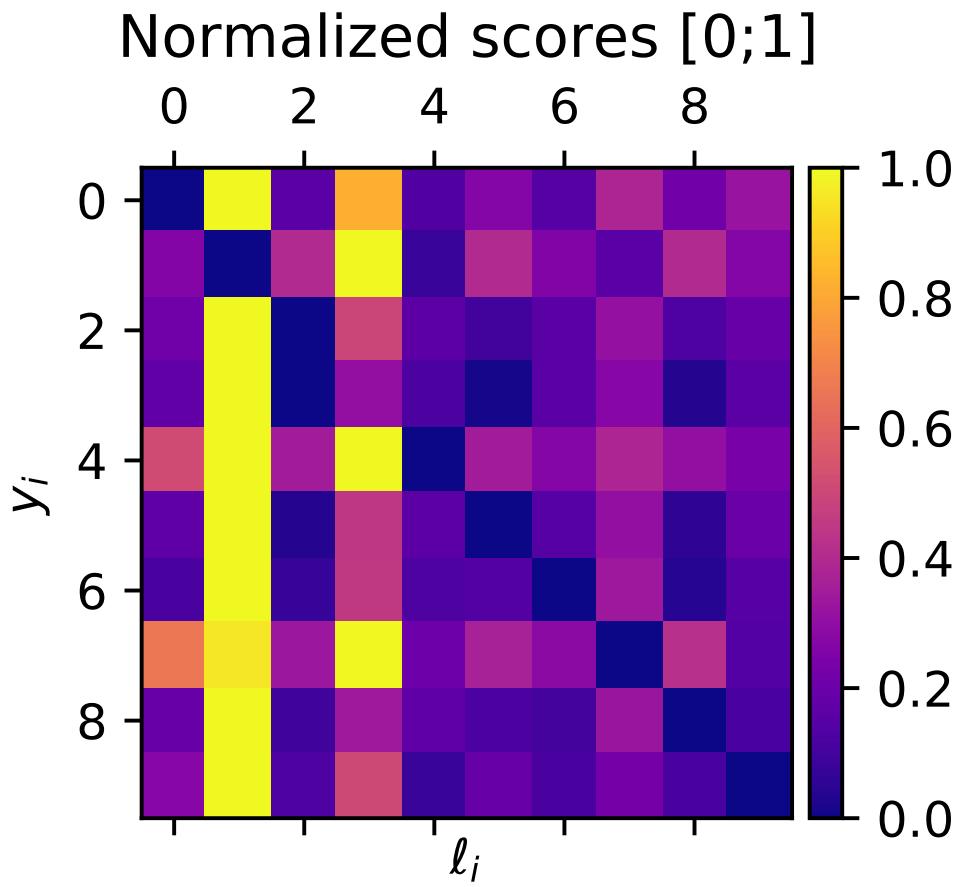


Figure 16: Mean autoencoder reconstruction errors. Rows correspond to input class, columns correspond to autoencoder idx. The plot only represents row-wise relationships, as each row is normalized such that the largest value is 1 and the lowest is 0.