
Use Perturbations when Learning from Explanations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning from explanations (MLX) is an approach to learning that uses
2 human-provided explanations of relevant or irrelevant features for each input to
3 ensure that model predictions are *right for the right reasons*. Existing MLX ap-
4 proaches rely on local model interpretation methods and require strong model
5 smoothing to align model and human explanations, leading to sub-optimal per-
6 formance. We recast MLX as a robustness problem, where human explanations
7 specify a lower dimensional manifold from which perturbations can be drawn, and
8 show both theoretically and empirically how this approach alleviates the need for
9 strong model smoothing. We consider various approaches to achieving robustness,
10 leading to improved performance over prior MLX methods. Finally, we show how
11 to combine robustness with an earlier MLX method, yielding state-of-the-art results
12 on both synthetic and real-world benchmarks.¹

13 1 Introduction

14 Deep neural networks (DNNs) display impressive capabilities, making them strong candidates for
15 real-world deployment. However, numerous challenges hinder their adoption in practice. Several
16 major deployment challenges have been linked to the fact that labelled data often under-specifies
17 the task (D’Amour et al., 2020). For example, systems trained on chest x-rays were shown to
18 generalise poorly because they exploited dataset-specific incidental correlations such as hospital tags
19 for diagnosing pneumonia (Zech et al., 2018; DeGrave et al., 2021). This phenomenon of learning
20 unintended feature-label relationships is referred to as *shortcut learning* (Geirhos et al., 2020) and is
21 a critical challenge to solve for trustworthy deployment of machine learning algorithms. A common
22 remedy to avoid shortcut learning is to train on diverse data (Shah et al., 2022) from multiple domains,
23 demographics, etc, thus minimizing the underspecification problem, but this may be impractical for
24 many applications such as in healthcare.

25 Enriching supervision through human-provided explanations of relevant and irrelevant re-
26 gions/features per example is an appealing direction toward reducing under-specification. For
27 instance, a (human-provided) explanation for chest x-ray classification may highlight scanning arti-
28 facts such as hospital tag as irrelevant features. Learning from such human-provided explanations
29 (MLX) has been shown to avoid known shortcuts (Schramowski et al., 2020). Ross et al. (2017)
30 pioneered an MLX approach based on regularizing DNNs, which was followed by several others
31 (Schramowski et al., 2020; Rieger et al., 2020; Stammer et al., 2021; Shao et al., 2021). Broadly,
32 existing approaches employ a model interpretation method to obtain per-example feature saliency, and
33 regularize such that model and human-provided explanations align. Since saliency is unbounded for
34 relevant features, many approaches simply regularize the salience of irrelevant features. In the same
35 spirit, we focus on handling a specification of irrelevant features, which we refer to as an explanation
36 hereafter. We collectively refer to existing MLX methods as *regularization-based*.

¹Code and data at this anonymous repository: https://github.com/vps-anonconfs/robust_mlx

37 Regularization-based approaches suffer from a critical concern stemming from their dependence on a
 38 local interpretation method. MLX methods based solely on local, i.e. example-specific, explanations
 39 do not have the desired affect of reducing shortcuts globally, i.e. over the entire input domain (see
 40 Figure 1). As we demonstrate both analytically and empirically, regularization-based MLX methods
 41 require strong model smoothing in order to be globally effective at reducing shortcut learning.

42 In this work, we explore learning from explanations using various robust training methods with the
 43 objective of training models that are robust to perturbations of irrelevant features. We start by framing
 44 the provided human explanations as specifications of a local, lower-dimensional manifold from which
 45 perturbations are drawn. We then notice that a model whose prediction is invariant to perturbations
 46 drawn from the manifold ought also to be robust to irrelevant features. Our perspective yields
 47 considerable advantages. Posing MLX as a robustness task enables us to leverage the considerable
 48 body of prior work in robustness. Further, we show in Section 4.1 that robust training can provably
 49 upper bound the deviation on model value when irrelevant features are perturbed without needing
 50 to impose model smoothing. However, when the space of irrelevant features is high-dimensional,
 51 robust-training may not fully suppress irrelevant features as explained in Section G. Accordingly, we
 52 explore combining both robustness-based and regularization-based methods, which achieves the best
 53 results. We highlight the following contributions:

- 54 • We theoretically and empirically demonstrate that existing MLX methods require strong model
 55 smoothing owing to their dependence on local model interpretation tools.
- 56 • We study learning from explanations using robust training methods. To the best of our knowledge,
 57 we are the first to analytically and empirically evaluate robust training methods for MLX.
- 58 • We distill our insights into our final proposal of combining robustness and regularization-based
 59 methods, which consistently out-performs the best regularization method and reduces the error rate
 60 by 20-90%.

61 2 Problem Definition and Background

62 We assume access to a training dataset with N training examples, $\mathcal{D}_T = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with
 63 $\mathbf{x}^{(i)} \in \mathbb{R}^d$ and $y^{(i)}$ label. In the MLX setting, a human expert also specifies input mask $\mathbf{m}^{(n)}$
 64 for an example $\mathbf{x}^{(n)}$ where non-zero values of the mask identify *irrelevant* features of the input
 65 $\mathbf{x}^{(n)}$. An input mask is usually designed to negate a known shortcut feature that a classifier is
 66 exploiting. Figure 2 shows some examples of masks for the datasets that we used for evaluation. For
 67 example, a mask in the ISIC dataset highlights a patch that was found to confound with non-cancerous
 68 images. With the added human specification, the augmented dataset contains triplets of example,
 69 label and mask, $\mathcal{D}_T = \{(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{m}^{(i)})\}_{i=0}^N$. The task therefore is to learn a model $f(\mathbf{x}; \theta)$ that fits
 70 observations well while not exploiting any features that are identified by the mask \mathbf{m} .

71 The method of Ross et al. (2017) which we call Grad-Reg (short for Gradient-Regularization), and
 72 also other similar approaches (Shao et al., 2021; Schramowski et al., 2020) employ an explanation
 73 algorithm (E) to assign importance scores to input features: $IS(\mathbf{x})$, which is then regularized with an
 74 $\mathcal{R}(\theta)$ term such that irrelevant features are not regarded as important. Their training loss takes the
 75 form shown in Equation 1 for an appropriately defined task-specific loss ℓ .

$$\begin{aligned}
 IS(\mathbf{x}) &\triangleq E(\mathbf{x}, f(\mathbf{x}; \theta)). \\
 \mathcal{R}(\theta) &\triangleq \sum_{n=1}^N \|IS(\mathbf{x}^{(n)}) \odot \mathbf{m}^{(n)}\|^2. \\
 \theta^* &= \arg \min_{\theta} \left\{ \sum_n \ell(f(\mathbf{x}^{(n)}; \theta), y^{(n)}) + \lambda \mathcal{R}(\theta) + \frac{1}{2} \beta \|\theta\|^2 \right\}. \tag{1}
 \end{aligned}$$

76 We use \odot to denote element-wise product throughout. CDEP (Rieger et al., 2020) is slightly different.
 77 They instead use an explanation method that also takes the mask as an argument to estimate the
 78 contribution of features identified by the mask, which they minimize similarly.

79 3 Method

80 Our methodology is based on the observation that an ideal model must be robust to perturbations
 81 to the irrelevant features. Following this observation, we reinterpret the human-provided mask as a
 82 specification of a lower-dimensional manifold from which perturbations are drawn and optimize the
 83 following objective.

$$\theta^* = \arg \min_{\theta} \sum_n \left\{ \ell \left(f(\mathbf{x}^{(n)}; \theta), y^{(n)} \right) + \alpha \max_{\epsilon: \|\epsilon\|_{\infty} \leq \kappa} \ell \left(f(\mathbf{x}^{(n)} + (\epsilon \odot \mathbf{m}^{(n)}); \theta), y^{(n)} \right) \right\} \quad (2)$$

84 The above formulation uses a weighting α to trade off between the standard task loss and perturbation
 85 loss and $\kappa > 0$ is a hyperparameter that controls the strength of robustness. We can leverage the
 86 many advances in robustness in order to approximately solve the inner maximization. We present
 87 them below.

88 **Avg-Ex:** We can approximate the inner-max with the empirical average of loss averaged over
 89 K samples drawn from the neighbourhood of training inputs. Singla et al. (2022) adopted this
 90 straightforward baseline for supervising using human-provided saliency maps on the Imagenet
 91 dataset. Similar to κ , we use σ to control the noise in perturbations as shown below.

$$\theta^* = \arg \min_{\theta} \sum_n \left\{ \ell \left(f(\mathbf{x}^{(n)}; \theta), y^{(n)} \right) + \frac{\alpha}{K} \sum_{\epsilon_j \sim \mathcal{N}(0, \sigma^2 I)} \ell \left(f(\mathbf{x}^{(n)} + (\epsilon_j \odot \mathbf{m}^{(n)}); \theta), y^{(n)} \right) \right\}$$

92 **PGD-Ex:** Optimizing for an estimate of worst perturbation through projected gradient descent
 93 (PGD) (Madry et al., 2017) is a popular approach from adversarial robustness. We refer to the
 94 approach of using PGD to approximate the second term of our loss as PGD-Ex and denote by
 95 $\epsilon^*(\mathbf{x}^{(n)}, \theta, \mathbf{m}^{(n)})$ the perturbation found by PGD at $\mathbf{x}^{(n)}$. Given the non-convexity of this problem,
 96 however, no guarantees can be made about the quality of the approximate solution \mathbf{x}^* .

$$\theta^* = \arg \min_{\theta} \sum_n \left\{ \ell \left(f(\mathbf{x}^{(n)}; \theta), y^{(n)} \right) + \alpha \ell \left(f(\mathbf{x}^{(n)} + (\epsilon^*(\mathbf{x}^{(n)}, \theta, \mathbf{m}^{(n)})); \theta), y^{(n)} \right) \right\}$$

97 **IBP-Ex:** Certified robustness approaches, on the other hand, minimize a certifiable upper-bound of
 98 the second term. A class of certifiable approaches known as interval bound propagation methods
 99 (IBP) (Mirman et al., 2018; Gowal et al., 2018) propagate input intervals to function value intervals
 100 that are guaranteed to contain true function values for any input in the input interval.

101 We define an input interval for $\mathbf{x}^{(n)}$ as $[\mathbf{x}^{(n)} - \kappa \mathbf{m}^{(n)}, \mathbf{x}^{(n)} + \kappa \mathbf{m}^{(n)}]$ where κ is defined in Eqn. 2.
 102 We then use bound propagation techniques to obtain function value intervals for the corresponding
 103 input interval: $\mathbf{l}^{(n)}, \mathbf{u}^{(n)}$, which are ranges over class logits. Since we wish to train a model that
 104 correctly classifies an example irrespective of the value of the irrelevant features, we wish to maximize
 105 the minimum probability assigned to the correct class, which is obtained by combining minimum
 106 logit for the correct class with maximum logit for incorrect class: $\tilde{f}(\mathbf{x}^{(n)}, y^{(n)}, \mathbf{l}^{(n)}, \mathbf{u}^{(n)}; \theta) \triangleq$
 107 $\mathbf{l}^{(n)} \odot \bar{\mathbf{y}}^{(n)} + \mathbf{u}^{(n)} \odot (\mathbf{1} - \bar{\mathbf{y}}^{(n)})$ where $\bar{\mathbf{y}}^{(n)} \in \{0, 1\}^c$ denotes the one-hot transformation of the label
 108 $y^{(n)}$ into a c -length vector for c classes. We refer to this version of the loss as IBP-Ex, summarized
 109 below.

$$\begin{aligned} \mathbf{l}^{(n)}, \mathbf{u}^{(n)} &= IBP(f(\bullet; \theta), [\mathbf{x}^{(n)} - \kappa \times \mathbf{m}^{(n)}, \mathbf{x}^{(n)} + \kappa \times \mathbf{m}^{(n)}]) \\ \tilde{f}(\mathbf{x}^{(n)}, y^{(n)}, \mathbf{l}^{(n)}, \mathbf{u}^{(n)}; \theta) &\triangleq \mathbf{l}^{(n)} \odot \bar{\mathbf{y}}^{(n)} + \mathbf{u}^{(n)} \odot (\mathbf{1} - \bar{\mathbf{y}}^{(n)}) \\ \theta^* &= \arg \min_{\theta} \sum_n \ell \left(f(\mathbf{x}^{(n)}; \theta), y^{(n)} \right) + \alpha \ell \left(\tilde{f}(\mathbf{x}^{(n)}, y^{(n)}, \mathbf{l}, \mathbf{u}; \theta), y^{(n)} \right) \end{aligned} \quad (3)$$

110 **Combined robustness and regularization:** PGD-Ex+Grad-Reg, IBP-Ex+Grad-Reg. We combine
 111 robustness and regularization by simply combining their respective loss terms. We show the objective
 112 for IBP-Ex+Grad-Reg below, PGD-Ex+Grad-Reg follows similarly.

$$\theta^* = \arg \min_{\theta} \sum_n \ell \left(f(\mathbf{x}^{(n)}; \theta), y^{(n)} \right) + \alpha \ell \left(\tilde{f}(\mathbf{x}^{(n)}, y^{(n)}, \mathbf{l}, \mathbf{u}; \theta), y^{(n)} \right) + \lambda \mathcal{R}(\theta). \quad (4)$$

113 $\lambda \mathcal{R}(\theta)$ and α, \tilde{f} are as defined in Eqn. 5 and Eqn. 3 respectively. In Section 4, F.1, we demonstrate
 114 the complementary strengths of robustness and regularization-based methods.

115 **4 Theoretical Motivation**

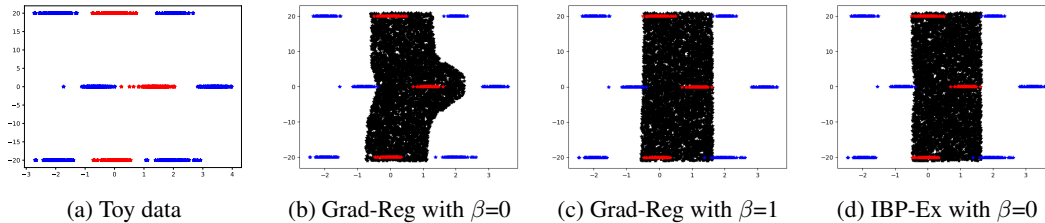


Figure 1: Illustration of the uneasy relationship between Grad-Reg and smoothing strength. (b) The decision boundary is nearly vertical (zero gradient wrt to nuisance y-axis value) for all training points and yet varies as a function of y value when Grad-Reg fitted using $\beta = 0$. (c) Grad-Reg requires strong model smoothing ($\beta = 1$) in order to translate local insensitivity to global robustness to x-coordinate. (d) IBP-Ex fits vertical pair of lines without any model smoothing.

116 In this section, we motivate the merits and drawbacks of robustness-based over regularization-based
 117 methods. Through non-parametric analysis in Theorems 1, 2, we argue that (a) regularization methods
 118 are robust to perturbations of irrelevant features (identified by the mask) only when the underlying
 119 model is sufficiently smoothed, thereby potentially compromising performance, (b) robust training
 120 upper-bounds deviation in function values when irrelevant features are perturbed, which can be
 121 further suppressed by using a more effective robust training. Although our analysis is restricted to
 122 nonparametric models for the ease of analysis, we empirically verify our claims with parametric
 123 neural network optimized using a gradient-based optimizer. We then highlight a limitation of
 124 robustness-based methods when the number of irrelevant features is large through Proposition 2.

125 **4.1 Merits of Robustness-based methods**

126 Consider a two-dimensional regression task, i.e. $\mathbf{x}^{(n)} \in \mathcal{X}$ and $y \in \mathbb{R}$. Assume that the second feature
 127 is the shortcut that the model should not use for prediction, and denote by $\mathbf{x}_j^{(n)}$ the j^{th} dimension
 128 of n^{th} point. We infer a regression function f from a Gaussian process prior $f \sim GP(f; 0, K)$
 129 with a squared exponential kernel where $k(x, \tilde{x}) = \exp(-\sum_i \frac{1}{2} \frac{(x_i - \tilde{x}_i)^2}{\theta_i^2})$. As a result, we have
 130 two hyperparameters θ_1, θ_2 , which are length scale parameters for the first and second dimensions
 131 respectively. Further, we impose a Gamma prior over the hyperparameters: $\mathcal{G}(\theta_i^{-2}; \alpha, \beta)$.

132 **Theorem 1 (Grad-Reg).** *We infer a regression function f from a GP prior as described above with the*
 133 *additional supervision of $[\partial f(\mathbf{x})/\partial x_2]_{\mathbf{x}^{(i)}} = 0, \forall i \in [1, N]$. Then the function value deviations*
 134 *to perturbations on irrelevant feature are lower bounded by a value proportional to the perturbation*
 135 *strength δ as shown below.*

$$f(\mathbf{x} + [0, \delta]^T) - f(\mathbf{x}) \geq \frac{2\delta\alpha}{\beta} \Theta(x_1^2 x_2^6 + \delta x_1^2 x_2^5) \quad (5)$$

136 Full proof of Theorem 1 is in Appendix A, we provide the proof outline below.

137 We observe from Theorem 1 that if we wish to infer a function that is robust to irrelevant feature
 138 perturbations, we need to set $\frac{\alpha}{\beta}$ to a very small value. Since the expectation of gamma distributed
 139 inverse-square length parameter is $\mathbb{E}[\theta^{-2}] = \frac{\alpha}{\beta}$, which we wish to set very small, we are, in effect,
 140 sampling functions with very large length scale parameter i.e. strongly smooth functions. This result
 141 brings us to the intuitive takeaway that regularization using Grad-Reg, or any local-interpretation
 142 methods that is closed under linear operation, applies globally only when the underlying family of
 143 functions is sufficiently smooth. One could also argue that we can simply use different priors for
 144 different dimensions, which would resolve the over-smoothing issue. However, we do not have access
 145 to parameters specific to each dimension in practice and especially with DNNs, therefore only overall
 146 smoothness may be imposed such as with parameter norm regularization in Eqn. 1.

147 We now look at properties of a function fitted using robustness methods and argue that they bound
 148 deviations in function values better. In order to express the bounds, we introduce a numerical quantity
 149 called coverage (C) to measure the effectiveness of a robust training method. We first define a notion

150 of inputs covered by a robust training method as $\hat{\mathcal{X}} \triangleq \{\mathbf{x} \mid \mathbf{x} \in \mathcal{X}, \ell(f(\mathbf{x}; \theta), y) < \phi\} \subset \mathcal{X}$ for
 151 a small positive threshold ϕ on loss. We define coverage as the maximum distance along second
 152 coordinate between any point in \mathcal{X} and its closest point in $\hat{\mathcal{X}}$, i.e. $C \triangleq \max_{\mathbf{x} \in \mathcal{X}} \min_{\hat{\mathbf{x}} \in \hat{\mathcal{X}}} |\mathbf{x}_2 - \hat{\mathbf{x}}_2|$.
 153 We observe that C is small if the robust training is effective. In the extreme case when training
 154 minimizes the loss for all points in the input domain, i.e. $\hat{\mathcal{X}} = \mathcal{X}$, then $C=0$.

155 **Theorem 2.** *When we use a robustness algorithm to regularize the network, the fitted function has
 156 the following property.*

$$|f(\mathbf{x} + [0, \delta]^T) - f(\mathbf{x})| \leq 2C \frac{\alpha}{\beta} \delta_{max} f_{max}. \quad (6)$$

157 δ_{max} and f_{max} are maximum values of Δx_2 and $f(\mathbf{x})$ in the input domain (\mathcal{X}) respectively.

158 Full proof is in Appendix B. The statement shows that deviations in function values are upper bounded
 159 by a factor proportional to C , which can be dampened by employing an effective robust training
 160 method. We can therefore control the deviations in function values without needing to regress $\frac{\alpha}{\beta}$ (i.e.
 161 without over-smoothing).

162 **Empirical verification with a toy dataset.** For empirical verification of our results, we fit a 3-layer
 163 feed-forward network on a two-dimensional data shown in Figure 1 (a), where color indicates the
 164 label. We consider fitting a model that is robust to changes in the second feature shown on y-axis.
 165 In Figures 1 (b), (c), we show the Grad-Reg fitted classifier using gradient ($\partial f / \partial x_2$ for our case)
 166 regularization for two different strengths of parameter smoothing (0 and 1 respectively). With weak
 167 smoothing, we observe that the fitted classifier is locally vertical (zero gradient along y-axis), but
 168 curved overall (Figure 1 (b)), which is fixed with strong smoothing (Figure 1 (c)). On the other hand,
 169 IBP-Ex fitted classifier is nearly vertical without any parameter regularization as shown in (d). This
 170 example illustrates the need for strong model smoothing when using a regularization-based method.

171 5 Experiments

172 We evaluate different methods on three datasets: one syn-
 173 thetic and two real-world. The synthetic dataset is similar
 174 to decoy-MNIST of Ross et al. (2017) with induced short-
 175 cuts and is presented in Section F.1. For evaluation on
 176 practical tasks, we evaluated on a plant phenotyping (Shao
 177 et al., 2021) task in Section F.2 and skin cancer detec-
 178 tion (Rieger et al., 2020) task presented in Section 5.3 All
 179 the datasets contain a known spurious feature, and were
 180 used in the past for evaluation of MLX methods. Figure 2
 181 summarises the three datasets, notice that we additionally
 182 require in the training dataset the specification of a mask
 183 identifying irrelevant features of the input; the patch for
 184 ISIC dataset, background for plant dataset, and decoy half
 185 for Decoy-MNIST images.

186 More details about experimental setup including metrics,
 187 network architecture, datasets, data splits, computing specs, and hyperparameters can be found in
 188 Appendix E.

189 5.1 Decoy-MNIST

190 Decoy-MNIST dataset is similar to MNIST-CIFAR dataset of Shah et al. (2020) where a very simple
 191 label-revealing color based feature (decoy) is juxtaposed with a more complex feature (MNIST
 192 image) as shown in Figure 1. We also randomly swap the position of decoy and MNIST parts, which
 193 makes ignoring the decoy part more challenging. We then validate and test on images where decoy
 194 part is set to correspond with random other label.

195 We make the following observations from Decoy-MNIST results presented in Table 1. ERM is only
 196 slightly better than a random classifier confirming the simplicity bias observed in the past (Shah et al.,
 197 2020). Grad-Reg, PGD-Ex and IBP-Ex perform comparably and better than ERM, but when combined
 198 (IBP-Ex+Grad-Reg, PGD-Ex+Grad-Reg) they far exceed their individual performances.

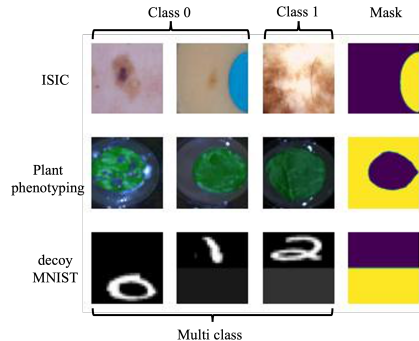


Figure 2: Sample images and masks for different datasets.

Dataset→	Decoy-MNIST		Plant		ISIC	
Method↓	Avg Acc	Wg Acc	Avg Acc	Wg Acc	Avg Acc	Wg Acc
ERM	15.1 ± 1.3	10.5 ± 5.4	71.3 ± 2.5	54.8 ± 1.3	77.3 ± 2.4	55.9 ± 2.3
G-DRO	64.1 ± 0.1	28.1 ± 0.1	74.2 ± 5.8	58.0 ± 4.6	66.6 ± 5.4	58.5 ± 10.7
Grad-Reg	72.5 ± 1.7	46.2 ± 1.1	72.4 ± 1.3	68.2 ± 1.4	76.4 ± 2.4	60.2 ± 7.4
CDEP	14.5 ± 1.8	10.0 ± 0.7	67.9 ± 10.3	54.2 ± 24.7	73.4 ± 1.0	60.9 ± 3.0
Avg-Ex	29.5 ± 0.3	19.5 ± 1.4	76.3 ± 0.3	64.5 ± 0.3	77.1 ± 2.1	55.2 ± 6.6
PGD-Ex	67.6 ± 1.6	51.4 ± 0.3	79.8 ± 0.3	78.5 ± 0.3	78.7 ± 0.5	64.4 ± 4.3
IBP-Ex	68.1 ± 2.2	47.6 ± 2.0	76.6 ± 3.5	73.8 ± 1.7	75.1 ± 1.2	64.2 ± 1.2
P+G	96.9 ± 0.3	95.8 ± 0.4	79.4 ± 0.5	76.7 ± 2.8	79.6 ± 0.5	67.5 ± 1.1
I+G	96.9 ± 0.2	95.0 ± 0.6	81.7 ± 0.2	80.1 ± 0.3	78.4 ± 0.5	65.2 ± 1.8

Table 1: Macro-averaged (Avg) accuracy and worst group (Wg) accuracy on (a) decoy-MNIST, (b) plant dataset, (c) ISIC dataset. Results are averaged over three runs and their standard deviation is shown after \pm . I+G is short for IBP-Ex+Grad-Reg and P+G for PGD-Ex+Grad-Reg. See text for more details.

199 5.2 Plant Phenotyping

200 Plant phenotyping is a real-world task of classifying images of a plant leaf as healthy or unhealthy.
 201 Schramowski et al. (2020) discovered that standard models exploited unrelated features from the
 202 nutritional solution in the background in which the leaf is placed, thereby performing poorly when
 203 evaluated outside of the laboratory setting. Thus, we aim to regulate the model not to focus on the
 204 background of the leaf using binary specification masks indicating where the background is located.
 205 More detailed analysis of the dataset can be found in Schramowski et al. (2020).

206 Table 1 contrasts different algorithms on the plant dataset. We
 207 visualize the interpretations of models obtained using Smooth
 208 Grad (Smilkov et al., 2017) trained with five different methods for
 209 three sample images from the train split in Figure 3. IBP-Ex draws
 210 features from a wider region and has more diverse pattern of active
 211 pixels, leading to higher Wg and Avg.

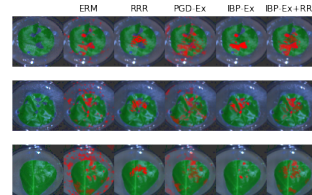


Figure 3: Visual heatmap of salient features for different algorithms on Plant data using SmoothGrad (Smilkov et al., 2017).

212 5.3 ISIC: Skin Cancer Detection

213 ISIC is a dataset of skin lesion images, which are to be classified
 214 cancerous or non-cancerous. Since half the non-cancerous images in
 215 the dataset contains a colorful patch as shown in Figure 2, standard
 216 DNN models depend on the presence of a patch for classification
 217 while compromising the accuracy on non-cancerous images without
 218 a patch (Codella et al., 2019; Tschandl et al., 2018).

219 We observe that Avg-Ex performed no better than ERM whereas
 220 PGD-Ex, IBP-Ex, IBP-Ex+Grad-Reg, and PGD-Ex+Grad-Reg sig-
 221 nificantly improved Wg accuracy over other baselines. The reduced
 222 accuracy gap between NPNC and C when using combined methods
 223 is indicative of reduced dependence on patch. Detailed results with
 224 error bars are shown in Table 4 of Appendix F.

Method	NPNC	PNC	C
ERM	55.9	96.5	79.6
Grad-Reg	67.1	99.0	63.2
CDEP	72.1	98.9	62.2
Avg-Ex	62.3	97.8	71.0
PGD-Ex	65.4	99.0	71.7
IBP-Ex	68.4	98.5	67.7
I+G	66.6	99.6	68.9
P+G	69.6	98.8	70.4

Table 2: Per-group accuracies on ISIC. Non-cancerous images without patch (NCNP) and with patch (NCP), and cancerous images (C).

225 6 Conclusions

226 By casting MLX as a robustness problem and using human expla-
 227 nations to specify the manifold of perturbations, we have shown that
 228 it is possible to alleviate the need for strong parameter smoothing
 229 of earlier approaches. Borrowing from the well-studied topic of ro-
 230 bustness, we evaluated two strong approaches, one from adversarial
 231 robustness (PGD-Ex) and one from certified robustness (IBP-Ex).

232 **Limitations.** Detecting and specifying irrelevant regions per-example by humans is a laborious and
 233 non-trivial task. Hence, it is interesting to see the effects of learning from incomplete explanations,
 234 which we leave for future work.

235 References

- 236 Codella, N., Rotemberg, V., Tschandl, P., Celebi, M. E., Dusza, S., Gutman, D., Helba, B.,
237 Kalloo, A., Liopyris, K., Marchetti, M., et al. Skin lesion analysis toward melanoma detec-
238 tion 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint*
239 *arXiv:1902.03368*, 2019.
- 240 DeGrave, A. J., Janizek, J. D., and Lee, S.-I. Ai for radiographic covid-19 detection selects shortcuts
241 over signal. *Nature Machine Intelligence*, 3(7):610–619, 2021.
- 242 D’Amour, A., Heller, K., Moldovan, D., Adlam, B., Alipanahi, B., Beutel, A., Chen, C., Deaton,
243 J., Eisenstein, J., Hoffman, M. D., et al. Underspecification presents challenges for credibility in
244 modern machine learning. *Journal of Machine Learning Research*, 2020.
- 245 Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A.
246 Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- 247 Gowal, S., Dvijotham, K., Stanforth, R., Bunel, R., Qin, C., Uesato, J., Arandjelovic, R., Mann, T.,
248 and Kohli, P. On the effectiveness of interval bound propagation for training verifiably robust
249 models. *arXiv preprint arXiv:1810.12715*, 2018.
- 250 Hennig, P., Osborne, M. A., and Kersting, H. P. *Probabilistic Numerics: Computation as Machine*
251 *Learning*. Cambridge University Press, 2022. doi: 10.1017/9781316681411.
- 252 Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models
253 resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- 254 Mirman, M., Gehr, T., and Vechev, M. Differentiable abstract interpretation for provably robust
255 neural networks. In *International Conference on Machine Learning*, pp. 3578–3586. PMLR, 2018.
- 256 Piratla, V., Netrapalli, P., and Sarawagi, S. Focus on the common good: Group distributional
257 robustness follows. *arXiv preprint arXiv:2110.02619*, 2021.
- 258 Rieger, L., Singh, C., Murdoch, W., and Yu, B. Interpretations are useful: penalizing explanations to
259 align neural networks with prior knowledge. In *International conference on machine learning*, pp.
260 8116–8126. PMLR, 2020.
- 261 Ross, A. S., Hughes, M. C., and Doshi-Velez, F. Right for the right reasons: Training differentiable
262 models by constraining their explanations. *arXiv preprint arXiv:1703.03717*, 2017.
- 263 Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. Distributionally robust neural networks for
264 group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint*
265 *arXiv:1911.08731*, 2019.
- 266 Schramowski, P., Stammer, W., Teso, S., Brugger, A., Herbert, F., Shao, X., Luigs, H.-G., Mahlein,
267 A.-K., and Kersting, K. Making deep neural networks right for the right scientific reasons by
268 interacting with their explanations. *Nature Machine Intelligence*, 2(8):476–486, 2020.
- 269 Shah, H., Tamuly, K., Raghunathan, A., Jain, P., and Netrapalli, P. The pitfalls of simplicity bias in
270 neural networks. *Advances in Neural Information Processing Systems*, 33:9573–9585, 2020.
- 271 Shah, R., Varma, V., Kumar, R., Phuong, M., Krakovna, V., Uesato, J., and Kenton, Z. Goal
272 misgeneralization: Why correct specifications aren’t enough for correct goals. *arXiv preprint*
273 *arXiv:2210.01790*, 2022.
- 274 Shao, X., Skryagin, A., Stammer, W., Schramowski, P., and Kersting, K. Right for better reasons:
275 Training differentiable models by constraining their influence functions. In *Proceedings of the*
276 *AAAI Conference on Artificial Intelligence*, volume 35, pp. 9533–9540, 2021.
- 277 Singla, S., Moayeri, M., and Feizi, S. Core risk minimization using salient imagenet. *arXiv preprint*
278 *arXiv:2203.15566*, 2022.
- 279 Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. Smoothgrad: removing noise by
280 adding noise. *arXiv preprint arXiv:1706.03825*, 2017.

-
- 281 Stammer, W., Schramowski, P., and Kersting, K. Right for the right concept: Revising neuro-symbolic
282 concepts by interacting with their explanations. In *Proceedings of the IEEE/CVF Conference on*
283 *Computer Vision and Pattern Recognition*, pp. 3619–3629, 2021.
- 284 Tschandl, P., Rosendahl, C., and Kittler, H. The ham10000 dataset, a large collection of multi-source
285 dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018.
- 286 Zech, J. R., Badgeley, M. A., Liu, M., Costa, A. B., Titano, J. J., and Oermann, E. K. Variable
287 generalization performance of a deep learning model to detect pneumonia in chest radiographs: a
288 cross-sectional study. *PLoS medicine*, 15(11):e1002683, 2018.
- 289 Zhang, H., Chen, H., Xiao, C., Goyal, S., Stanforth, R., Li, B., Boning, D., and Hsieh, C.-J. Towards
290 stable and efficient training of verifiably robust neural networks. *arXiv preprint arXiv:1906.06316*,
291 2019.

Supporting material for "Use Perturbations when Learning from Explanations"

A Proof of Theorem 1

We restate the result of Theorem 1 for clarity.

The posterior mean of the function estimates marginalised over hyperparameters with Gamma prior has the following closed form.

$$\begin{aligned}
 f(x) &\triangleq \mathbb{E}_\theta[m_x] = \int \int m_x \mathcal{G}(\theta_1^{-2}; \alpha, \beta) \mathcal{G}(\theta_2^{-2}; \alpha, \beta) d\theta_1^{-2} d\theta_2^{-2} \\
 f(\mathbf{x}) &= \sum_{n=1}^N \left(\frac{1}{1 + \frac{d(x_1, x_1^{(n)})}{\beta}} \right)^\alpha \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}} \right)^\alpha \left[\tilde{y}^{(n)} + \frac{\frac{\alpha}{\beta} (x_2 - x_2^{(n)})}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}} \tilde{y}^{(n+N)} \right] \\
 f(\mathbf{x} + [0, \delta]^T) - f(\mathbf{x}) &\geq \frac{2\delta\alpha}{\beta} \sum_n \left(\frac{1}{1 + \frac{d(x_1, x_1^{(n)})}{\beta}} \right)^\alpha \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}} \right)^{\alpha+1} \\
 &\quad \left[(\alpha + 1) \tilde{y}_{n+N} \left(\frac{2(x_2 - x_2^{(n)})[x_2 + \delta - x_2^{(n)}]}{\beta + d(x_2, x_2^{(n)})} - 1 \right) - \tilde{y}_n \right]
 \end{aligned}$$

Proof. We first derive the augmented set of observations (\hat{y}) and \hat{K} explained in the main section.

$$\begin{aligned}
 \hat{y} &= [y_1, y_2, \dots, y_N, \partial f(\mathbf{x}^{(1)})/\partial x_2, \partial f(\mathbf{x}^{(2)})/\partial x_2, \dots, \partial f(\mathbf{x}^{(N)})/\partial x_2]^T \\
 k(x^{(i)}, x^{(j)}) &= \begin{cases} \exp(-\frac{1}{2} \sum_{k=1}^2 \frac{(x_k^{(i)} - x_k^{(j)})^2}{\theta_k^2}) & \text{when } i, j \leq N \\ \frac{(x_2^{(i)} - x_2^{(j)})}{\theta_2^2} \exp(-\frac{1}{2} \sum_{k=1}^2 \frac{(x_k^{(i)} - x_k^{(j)})^2}{\theta_k^2}) & \text{when } i \leq N, j > N \\ -\frac{(x_2^{(i)} - x_2^{(j)})}{\theta_2^2} \exp(-\frac{1}{2} \sum_{k=1}^2 \frac{(x_k^{(i)} - x_k^{(j)})^2}{\theta_k^2}) & \text{when } j \leq N, i > N \\ -2 \frac{(x_2^{(i)} - x_2^{(j)})^2}{\theta_2^4} \exp(-\frac{1}{2} \sum_{k=1}^2 \frac{(x_k^{(i)} - x_k^{(j)})^2}{\theta_k^2}) & \\ \quad + \frac{1}{\theta_2^2} \exp(-\frac{1}{2} \sum_{k=1}^2 \frac{(x_k^{(i)} - x_k^{(j)})^2}{\theta_k^2}) & \text{when } i, j > N \end{cases}
 \end{aligned}$$

These results follow directly from the results on covariance between observations of f and its partial derivative below (Hennig et al., 2022).

$$\begin{aligned}
 \text{cov}(f(x), \frac{\partial f(\tilde{x})}{\partial \tilde{x}}) &= \frac{\partial k(x, \tilde{x})}{\partial \tilde{x}} \\
 \text{cov}(\frac{\partial f(x)}{\partial x}, \frac{\partial f(\tilde{x})}{\partial \tilde{x}}) &= \frac{\partial^2 k(x, \tilde{x})}{\partial x \partial \tilde{x}}
 \end{aligned}$$

The posterior value of the function at an arbitrary point \mathbf{x} would then be of the form $p(f(\mathbf{x}) | \mathcal{D}) \sim \mathcal{N}(f(\mathbf{x}); m_x, k_x)$ where m_x and k_x are have the following closed form for Gaussian prior and Gaussian likelihood in our case.

$$\begin{aligned}
 m_x &= k(x, X) K_{XX}^{-1} \hat{y} \\
 k_x &= k(x, x) - k(x, X) K_{XX}^{-1} k(X, x)
 \end{aligned}$$

Since m_x, k_x are functions of the parameters θ_1, θ_2 , we obtain the closed form for posterior mean by imposing a Gamma prior over the two parameters. For brevity, we denote by $d(x, \tilde{x}) = (x - \tilde{x})^2/2$ and $\tilde{y}^{(i)}$ is the i^{th} component of $\hat{K}_{XX}^{-1} \hat{y}$.

$$\begin{aligned}
 f(x) &\triangleq \mathbb{E}_\theta[m_x] = \int \int m_x \mathcal{G}(\theta_1^{-2}; \alpha, \beta) \mathcal{G}(\theta_2^{-2}; \alpha, \beta) d\theta_1^{-2} d\theta_2^{-2} \\
 &= \int \int \left[\sum_{n=1}^N k(x, x^{(n)}) \tilde{y}_n + \sum_{n=1}^N \frac{(x_2 - x_2^{(n)})}{\theta_2^2} k(x, x^{(n)}) \tilde{y}_{n+N} \right] \mathcal{G}(\theta_1^{-2}; \alpha, \beta) \mathcal{G}(\theta_2^{-2}; \alpha, \beta) d\theta_1^{-2} d\theta_2^{-2}
 \end{aligned}$$

$$\begin{aligned}
& \int \int k(\mathbf{x}, \mathbf{x}^{(n)}) \tilde{y}_n \mathcal{G}(\theta_1^{-2}; \alpha, \beta) \mathcal{G}(\theta_2^{-2}; \alpha, \beta) d\theta_1^{-2} d\theta_2^{-2} \\
&= \int \int \exp\left(-\frac{\theta_1^{-2}(x_1 - x_1^{(n)})^2}{2} + \frac{\theta_2^{-2}(x_2 - x_2^{(n)})^2}{2}\right) \frac{\beta^\alpha}{\Gamma(\alpha)} \theta_1^{-2\alpha+2} \exp(-\beta\theta_1^{-2}) \\
&\quad \frac{\beta^\alpha}{\Gamma(\alpha)} \theta_2^{-2\alpha+2} \exp(-\beta\theta_2^{-2}) \tilde{y}_n d\theta_1^{-2} d\theta_2^{-2} \\
&= \left(\frac{\beta}{\beta + \frac{(x_1 - x_1^{(n)})^2}{2}}\right)^\alpha \left(\frac{\beta}{\beta + \frac{(x_2 - x_2^{(n)})^2}{2}}\right)^\alpha \tilde{y}_n \\
&\quad \int \int \frac{x_2 - x_2^{(n)}}{\theta_2^2} k(\mathbf{x}, \mathbf{x}^{(n)}) \tilde{y}_{n+N} \mathcal{G}(\theta_1^{-2}; \alpha, \beta) \mathcal{G}(\theta_2^{-2}; \alpha, \beta) d\theta_1^{-2} d\theta_2^{-2} \\
&= (x_2 - x_2^{(n)}) \left(\frac{\beta}{\beta + \frac{(x_1 - x_1^{(n)})^2}{2}}\right)^\alpha \frac{\beta^\alpha / \Gamma(\alpha)}{(\beta + \frac{(x_2 - x_2^{(n)})^2}{2})^{\alpha+1} / \Gamma(\alpha+1)} \tilde{y}_{n+N} \\
&= \left(\frac{\beta}{\beta + \frac{(x_1 - x_1^{(n)})^2}{2}}\right)^\alpha \frac{\alpha(x_2 - x_2^{(n)})}{\beta + \frac{(x_2 - x_2^{(n)})^2}{2}} \left(\frac{\beta}{\beta + \frac{(x_2 - x_2^{(n)})^2}{2}}\right)^\alpha \tilde{y}_{n+N}
\end{aligned}$$

307 Overall, we have the following result.

$$f(x) = \sum_{n=1}^N \left(\frac{1}{1 + \frac{d(x_1, x_1^{(n)})}{\beta}}\right)^\alpha \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \left[\tilde{y}_n + \frac{\frac{\alpha}{\beta}(x_2 - x_2^{(n)})}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}} \tilde{y}_{n+N} \right]$$

308 We now derive the sensitivity to perturbations on the second dimension for $\Delta \mathbf{x} = [0, \delta]^T$.

$$\begin{aligned}
f(\mathbf{x} + \Delta \mathbf{x}) - f(\mathbf{x}) &= \sum_{n=1}^N \left(\frac{1}{1 + \frac{d(x_1, x_1^{(n)})}{\beta}}\right)^\alpha \left\{ \left[\left(\frac{1}{1 + \frac{d(x_2 + \delta, x_2^{(n)})}{\beta}}\right)^\alpha - \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \right] \tilde{y}_n \right. \\
&\quad \left. \left[\frac{\frac{\alpha}{\beta}(x_2 + \delta - x_2^{(n)})}{(1 + \frac{d(x_2 + \delta, x_2^{(n)})}{\beta})^{\alpha+1}} - \frac{\frac{\alpha}{\beta}(x_2 - x_2^{(n)})}{(1 + \frac{d(x_2, x_2^{(n)})}{\beta})^{\alpha+1}} \right] \tilde{y}_{n+N} \right\} \quad (7)
\end{aligned}$$

309 Using Bernoulli inequality, $(1+x)^r \geq 1+rx$ if $r \leq 0$, we derive the following inequalities.

$$\begin{aligned}
& \left(\frac{1}{1 + \frac{d(x_2 + \delta, x_2^{(n)})}{\beta}}\right)^\alpha - \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \\
&= \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \left[\left(\frac{\beta + d(x_2, x_2^{(n)})}{\beta + d(x_2 + \delta, x_2^{(n)})}\right)^\alpha - 1 \right] \\
&\geq \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha - \alpha \left[\frac{\beta + d(x_2 + \delta, x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})} - 1 \right] \\
&= \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \alpha \left[\frac{d(x_2, x_2^{(n)}) - d(x_2 + \delta, x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})} \right]
\end{aligned}$$

$$\text{Assuming } |x_2 - x_2^{(n)}| \gg \delta \quad \forall n \in [N] \quad (8)$$

$$\approx \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \alpha \left[\frac{-2\delta(x_2 - x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})} \right] \quad (9)$$

310 Similarly,

$$\begin{aligned}
& \frac{\frac{\alpha}{\beta}(x_2 + \delta - x_2^{(n)})}{\left(1 + \frac{d(x_2 + \delta, x_2^{(n)})}{\beta}\right)^{\alpha+1}} - \frac{\frac{\alpha}{\beta}(x_2 - x_2^{(n)})}{\left(1 + \frac{d(x_2, x_2^{(n)})}{\beta}\right)^{\alpha+1}} \\
& \geq \frac{\alpha}{\beta}(x_2 - x_2^{(n)}) \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^{\alpha+1} (\alpha + 1) \left[\frac{-2\delta(x_2 - x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})}\right] + \frac{\delta \frac{\alpha}{\beta}}{\left(1 + \frac{d(x_2 + \delta, x_2^{(n)})}{\beta}\right)^{\alpha+1}} \\
& \geq \frac{\alpha}{\beta}(x_2 - x_2^{(n)}) \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^{\alpha+1} (\alpha + 1) \left[\frac{-2\delta(x_2 - x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})}\right] \\
& \quad + \frac{\delta \frac{\alpha}{\beta}}{\left(1 + \frac{d(x_2, x_2^{(n)})}{\beta}\right)^{\alpha+1}} (\alpha + 1) \left[\frac{-2\delta(x_2 - x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})} + 1\right] \\
& = \frac{\alpha + 1}{\left(1 + \frac{d(x_2, x_2^{(n)})}{\beta}\right)^{\alpha+1}} \left[\frac{-2\delta(x_2 - x_2^{(n)})^2 \alpha / \beta - 2\delta^2 \alpha / \beta (x_2 - x_2^{(n)})}{\beta + d(x_2, x_2^{(n)})} + \frac{\delta \alpha}{\beta}\right] \\
& = \frac{-2\delta \alpha (\alpha + 1)}{\beta \left(1 + \frac{d(x_2, x_2^{(n)})}{\beta}\right)^{\alpha+1}} \left[\frac{-2(x_2 - x_2^{(n)})[x_2 + \delta - x_2^{(n)}]}{\beta + d(x_2, x_2^{(n)})} + 1\right] \tag{10}
\end{aligned}$$

311 Using inequalities 9, 10 in Equation 7, we have the following.

$$\begin{aligned}
f(\mathbf{x} + \Delta \mathbf{x}) - f(\mathbf{x}) & \geq \sum_n \left(\frac{1}{1 + \frac{d(x_1, x_1^{(n)})}{\beta}}\right)^\alpha \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^\alpha \\
& \quad \left[\frac{-2\delta \alpha \tilde{y}_n}{\beta + d(x_2, x_2^{(n)})} + \frac{-2\delta \alpha (\alpha + 1) \tilde{y}_{n+N}}{\beta + d(x_2, x_2^{(n)})} \left(\frac{-2(x_2 - x_2^{(n)})[x_2 + \delta - x_2^{(n)}]}{\beta + d(x_2, x_2^{(n)})} + 1\right)\right]
\end{aligned}$$

312

$$\begin{aligned}
f(\mathbf{x} + \Delta \mathbf{x}) - f(\mathbf{x}) & \geq \frac{2\delta \alpha}{\beta} \sum_n \left(\frac{1}{1 + \frac{d(x_1, x_1^{(n)})}{\beta}}\right)^\alpha \left(\frac{1}{1 + \frac{d(x_2, x_2^{(n)})}{\beta}}\right)^{\alpha+1} \\
& \quad \left[(\alpha + 1) \tilde{y}_{n+N} \left(\frac{2(x_2 - x_2^{(n)})[x_2 + \delta - x_2^{(n)}]}{\beta + d(x_2, x_2^{(n)})} - 1\right) - \tilde{y}_n\right] \tag{11}
\end{aligned}$$

313 Using the inequality $(1 + x)^r \geq 1 + rx$ if $r \leq 0$, we have

$$\begin{aligned}
f(\mathbf{x} + \Delta \mathbf{x}) - f(\mathbf{x}) & \geq \frac{2\delta \alpha}{\beta} \sum_n \left\{ \left(1 - \frac{\alpha}{\beta} d(x_1, x_1^{(n)})\right) \left(1 - \frac{\alpha + 1}{\beta} d(x_2, x_2^{(n)})\right) \right. \\
& \quad \left. \left[\frac{\alpha + 1}{\beta} \tilde{y}_{n+N} \left(2(x_2 - x_2^{(n)})[x_2 + \delta - x_2^{(n)}](1 - d(x_2, x_2^{(n)})) - 1\right) - \tilde{y}_n\right] \right\} \\
& = \frac{2\delta \alpha}{\beta} \Theta(x_1^2 x_2^6 + \delta x_1^2 x_2^5)
\end{aligned}$$

314

□

315 B Proof of Theorem 2

316 We restate the result of Theorem 2 for clarity.

317 When we use an adversarial robustness algorithm to regularize the network, the fitted function has
 318 the following property.

$$|f(\mathbf{x} + [0, \delta]^T) - f(\mathbf{x})| \leq \frac{\alpha}{\beta} \delta_{max} f_{max} C$$

where $C = \max_{\mathbf{x} \in \mathcal{X}} \min_{\hat{\mathbf{x}} \in \hat{\mathcal{X}}} |\mathbf{x}_2 - \hat{\mathbf{x}}_2|$

319 δ_{max} and f_{max} are maximum value of Δx_2 and $f(\mathbf{x})$ in the input domain (\mathcal{X}) respectively. $\hat{\mathcal{X}}$
 320 denotes the subset of inputs covered by the robustness method. C therefore captures the maximum
 321 gap in coverage of the robustness method.

322 *Proof.* We begin by estimating the Lipschitz constant of a GP with squared exponential kernel.

$$\begin{aligned} f(\mathbf{x}) &= K_{xX} K_{XX}^{-1} y \\ \frac{\partial f(x)}{\partial x_2} &= \frac{\partial K_{xX} K_{XX}^{-1} y}{\partial x_2} = \tilde{K}_{xX} K_{XX}^{-1} y \\ \text{where } [\tilde{K}_{xX}]_n &= \frac{\partial}{\partial x_2} \exp\left(-\frac{((x_1 - x_1^{(n)})^2 + (x_2 - x_2^{(n)})^2)}{2\theta^2}\right) \\ &= -\frac{(x_2 - x_2^{(n)})}{\theta^2} [K_{xX}]_n \\ \implies \frac{\partial f(x)}{\partial x_2} &= -\left[\sum_{n=1}^N \frac{(x_2 - x_2^{(n)})}{\theta^2} [K_{xX}]_n\right] K_{XX}^{-1} y \end{aligned}$$

323 We denote with δ_{max} the maximum deviation of any input from the training points, i.e. we define
 324 δ_{max} as $\max_{\mathbf{x} \in \mathcal{X}} \min_{n \in [N]} |x_2 - x_2^{(n)}|$. Also, we denote by f_{max} the maximum function value in
 325 the input domain, i.e. $f_{max} \triangleq \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. We can then bound the partial derivative wrt second
 326 dimension as follows.

$$\frac{\partial f(\mathbf{x})}{\partial x_2} \leq \frac{\delta_{max} f_{max}}{\theta^2} \leq \frac{\delta_{max} f_{max}}{\theta^2}$$

327 For any arbitrary point \mathbf{x} , the maximum function deviation is upper bounded by the product of
 328 maximum slope and maximum distance from the closest point covered by the adversarial distance
 329 method.

$$|f([x_1, x_2]^T) - f([x_1, \hat{x}_2]^T)| \leq \frac{\delta_{max} f_{max}}{\theta^2} \max_{\mathbf{x} \in \mathcal{X}} \min_{\hat{\mathbf{x}} \in \hat{\mathcal{X}}} |x_2 - \hat{x}_2| = \frac{\delta_{max} f_{max}}{\theta^2} C$$

330 Therefore,

$$|f(\mathbf{x} + [0, \delta]^T) - f(\mathbf{x})| \leq 2 \frac{\delta_{max} f_{max}}{\theta^2} C$$

331 Marginalising θ^{-2} with the Gamma prior leads to the final form below.

$$|f(\mathbf{x} + [0, \delta]^T) - f(\mathbf{x})| \leq 2C \frac{\alpha}{\beta} \delta_{max} f_{max}$$

332

□

333 C Proof of Proposition 2

334 We restate the result here for clarity.

335 Consider a regression task with $D + 1$ -dimensional inputs \mathbf{x} where the first D dimensions are
 336 irrelevant, and assume they are $x_d = y, d \in [1, D]$ while $x_{D+1} \sim \mathcal{N}(y, 1/K)$. The MAP estimate
 337 of linear regression parameters $f(\mathbf{x}) = \sum_{d=1}^{D+1} w_d x_d$ when fitted using Avg-Ex are as follows:
 338 $w_d = 1/(D + K), d \in [1, D]$ and $w_{D+1} = K/(K + D)$.

339 *Proof.* Without loss of generality, we assume α, σ^2 parameters of Avg-Ex are set to 1. In effect, our
 340 objective is to fit parameters that predict well for inputs sampled using standard normal perturbations,
 341 i.e. $\mathbf{x}^{(n)} + \mathbf{m}\epsilon, \forall n \in [1, N], \epsilon \sim \mathcal{N}(0, 1), \mathbf{m} = [1, 1, \dots, 1, 0]^T \in \{0, 1\}^{D+1}$. The original problem
 342 therefore is equivalent to fitting on transformed input $\hat{\mathbf{x}}$ such that $\hat{x}_i^{(n)} \sim \mathcal{N}(y, \sigma_i^2)$ where $\sigma_i^2 = 1$ for
 343 all $i \leq D$ and is $1/K$ when $i = D + 1$.

344 Likelihood of observations for the equivalent problem is obtained as follows.

$$\begin{aligned}
 P(y | \hat{x}_1, \hat{x}_2, \dots, \hat{x}_{D+1}) &= \prod_{i=1}^{D+1} P(y | \hat{x}_i) \propto \prod_{i=1}^{D+1} P(\hat{x}_i | y) P(y) \\
 &= \prod_i \mathcal{N}(\hat{x}_i; y, \sigma_i^2) \propto \exp\left(-\sum_i \frac{(y - \hat{x}_i)^2}{2\sigma_i^2}\right) \\
 &= \exp\left\{-y^2\left(\sum_i \frac{1}{2\sigma_i^2}\right) + y\left(\sum_i \frac{\hat{x}_i}{\sigma_i^2}\right) + \sum_i \frac{\hat{x}_i^2}{2\sigma_i^2}\right\} \\
 &\propto \mathcal{N}\left(y; \sum_i \frac{\hat{x}_i}{\sigma_i^2} P, P\right) \\
 \text{where } P &= \frac{1}{\sum_i 1/\sigma_i^2}
 \end{aligned}$$

345 Substituting, the value of σ_i defined as above, we have $P=D+K$ and the MLE estimate for the linear
 346 regression parameters are as shown in the statement. The MAP estimate also remains the same since
 347 we do not impose any informative prior on the regression weights. \square

348 D Parametric Model Analysis

349 In this section we show that a similar result to what is shown for non-parametric models also holds
 350 for parametric models. We will analyse the results for a two-layer neural networks with ReLU
 351 activations. We consider a more general case of D dimensional input where the first d dimensions
 352 identify the spurious features. We wish to fit a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ such that $f(\mathbf{x})$ is robust to
 353 perturbations to the spurious features. We have the following bound when training a model using
 354 gradient regularization of Ross et al. (2017).

355 **Proposition 1.** *We assume that the model is parameterised as a two-layer network with ReLU*
 356 *activations such that $f(\mathbf{x}) = \sum_j \beta_j \phi(\sum_i w_{ji} x_i + b_j)$ where $\vec{\beta} \in \mathbb{R}^F, \vec{w} \in \mathbb{R}^{F \times D}, \vec{b} \in \mathbb{R}^F$ are the*
 357 *parameters, and $\phi(z) = \max(z, 0)$ is the ReLU activation. For any function such that gradients*
 358 *wrt to the first d features is exactly zero, i.e. $\frac{\partial f}{\partial x_i} |_{\mathbf{x}_i^{(n)}} = 0 \quad \forall i \in [1, d], n \in [1, N]$, we have the*
 359 *following bound on the function value deviations for input perturbations from a training instance \mathbf{x} :*
 360 *$\tilde{\mathbf{x}} - \mathbf{x} = \Delta \mathbf{x} = [\Delta \mathbf{x}_{1:d}^T, \mathbf{0}_{d+1:D}^T]^T$.*

$$|f(\tilde{\mathbf{x}}) - f(\mathbf{x})| = \Theta((\|\vec{\beta}\|^2 + \|\vec{w}\|_F^2) \|\Delta \mathbf{x}\|) \quad (12)$$

361 For a two-layer network trained to regularize gradients wrt first d dimensions on training data, the
 362 function value deviation from an arbitrary point $\tilde{\mathbf{x}}$ from a training point \mathbf{x} such that $\tilde{\mathbf{x}} - \mathbf{x} = \Delta \mathbf{x} =$
 363 $[\Delta \mathbf{x}_{1:d}^T, \mathbf{0}_{d+1:D}^T]^T$ is bounded as follows.

$$|f(\tilde{\mathbf{x}}) - f(\mathbf{x})| = \Theta((\|\vec{\beta}\|^2 + \|\vec{w}\|_F^2) \|\Delta \mathbf{x}\|)$$

364 *Proof.* Recall that the function is parameterised using parameters $\vec{w}, \vec{b}, \vec{\beta}$ such that $f(\mathbf{x}) =$
 365 $\sum_j \beta_j \phi(\sum_i w_{ji} x_i + b_j)$ where $\vec{\beta} \in \mathbb{R}^F, \vec{w} \in \mathbb{R}^{F \times D}, \vec{b} \in \mathbb{R}^F$ are the parameters, and $\phi(z) =$
 366 $\max(z, 0)$ is the ReLU activation.

367 Since we train such that $\frac{\partial f(\mathbf{x})}{\partial x_i} = 0, \quad i \in [1, d]$, we have that $\frac{\partial f(\mathbf{x})}{\partial x_i} = \sum_j \beta_j \hat{\phi}(\sum_i w_{ij} x_i + b_i) w_{ij}$
 368 where $\hat{\phi}(a) = \max(\frac{a}{|a|}, 0)$.

369 We now bound the variation in the function value for changes in the input when moving from $\mathbf{x} \rightarrow \tilde{\mathbf{x}}$
 370 where \mathbf{x} is an instance from the training data. We define four groups of neurons based on the sign of
 371 $\sum_i w_{ji}x_i + b_j$ and $\sum_i w_{ji}\tilde{x}_i + b_j$. g_1 is both positive, g_2 is negative and positive, g_3 is positive and
 372 negative, g_4 is both negative. By defining groups, we can omit the ReLU activations as below.

$$\begin{aligned} f(\tilde{\mathbf{x}}) - f(\mathbf{x}) &= \sum_j \beta_j \phi\left(\sum_i w_{ji}\tilde{x}_i + b_j\right) - \sum_j \beta_j \phi\left(\sum_i w_{ji}x_i + b_j\right) \\ &= \sum_{j \in g_1} \beta_j \sum_i w_{ji}(\tilde{x}_i - x_i) + \sum_{j \in g_2} \beta_j \left(\sum_i w_{ji}\tilde{x}_i + b_j\right) - \sum_{j \in g_3} \beta_j \left(\sum_i w_{ji}x_i + b_j\right) \\ &= \sum_{j \in g_1} \beta_j \sum_{i=1}^d w_{ji}(\tilde{x}_i - x_i) + \sum_{j \in g_2} \beta_j \left(\sum_{i=1}^D w_{ji}\tilde{x}_i + b_j\right) - \sum_{j \in g_3} \beta_j \left(\sum_{i=1}^D w_{ji}x_i + b_j\right) \end{aligned}$$

373 Since we have that $\sum_{j \in g_1 \cup g_3} \beta_j w_{ij} = 0, \forall i \in [1, d]$, we have

$$\begin{aligned} &= \sum_{j \in g_1} \beta_j \sum_{i=1}^d w_{ji}\tilde{x}_i + \sum_{j \in g_2} \beta_j \left(\sum_{i=1}^d w_{ji}\tilde{x}_i + \sum_{i=d+1}^D w_{ji}x_i + b_j\right) - \sum_{j \in g_3} \beta_j \left(\sum_{i=d+1}^D w_{ji}x_i + b_j\right) \\ &\quad - \underbrace{\sum_{j \in g_1} \beta_j \sum_{i=1}^d w_{ji}x_i - \sum_{j \in g_3} \beta_j \sum_{i=1}^d w_{ji}x_i}_{=\sum_{i=1}^d x_i \sum_{j \in g_1 \cup g_3} \beta_j w_{ji} = 0} \end{aligned}$$

374

$$= \sum_{j \in g_1 \cup g_2} \beta_j \sum_{i=1}^d w_{ji}\tilde{x}_i + \sum_{j \in g_2} \beta_j \left(\sum_{i=d+1}^D w_{ji}x_i + b_j\right) - \sum_{j \in g_3} \beta_j \left(\sum_{i=d+1}^D w_{ji}x_i + b_j\right)$$

375 retaining only the terms that depend on $\Delta x = \tilde{x} - x$, the expression is further simplified as a term
 376 that grows with $\Delta \mathbf{x}$ and a constant term that depends on the value of \mathbf{x}

$$\begin{aligned} &= \sum_{j \in g_1 \cup g_2} \beta_j \sum_{i=1}^d w_{ji}\Delta x_i + \text{constant} \\ \implies &= \Theta(\|\beta\| \|\vec{w}\|_F \|\Delta \mathbf{x}\|) \quad \text{Cauchy-Schwartz inequality} \\ &= \Theta(\|\beta\|^2 + \|\vec{w}\|_F^2) \|\Delta \mathbf{x}\| \end{aligned}$$

377

□

378 E Further Experiment Details

379 E.1 Setup

380 E.1.1 Baselines

381 We denote by ERM the simple minimization of cross-entropy loss (using only the first loss term of
 382 Equation 1). We also compare with G-DRO(Sagawa et al., 2019), which also has the objective of
 383 avoiding to learn known irrelevant features but is supervised through group label (see Section ??).
 384 Although the comparison is unfair toward G-DRO because MLX methods use richer supervision of
 385 per-example masks, it serves as a baseline that can be slightly better than ERM in some cases.

386 **Regulaization-based methods.** Grad-Reg and CDEP, which were discussed in Section 2. We omit
 387 comparison with Shao et al. (2021) because their code is not publicly available and is non-trivial to
 388 implement the influence-function based regularization.

389 **Robustness-based methods.** Avg-Ex, PGD-Ex, IBP-Ex along with **combined robustness and**
 390 **regularization methods.** IBP-Ex+Grad-Reg, PGD-Ex+Grad-Reg that are described in Section 3.

391 E.1.2 Metrics

392 We report performance using two metrics that indicate if the model is using irrelevant features (Wg
393 Acc) without compromising the average accuracy (Avg Acc).

394 **Avg Acc.** Since the two real-world datasets contain imbalanced class populations, we only report
395 accuracy macro-averaged over labels, simply denoted as “Avg Acc”.

396 **Wg Acc.** Worst accuracy among groups where groups are appropriately defined. Different labels
397 define the groups for decoy-MNIST and plant dataset, which therefore have ten and two groups
398 respectively. In ISIC dataset, different groups are defined by the cross-product of label and presence
399 or absence of the patch. We denote this metric as “Wg Acc”, which is a standard metric when
400 evaluating on datasets with shortcut features (Sagawa et al., 2019; Piratla et al., 2021).

401 E.1.3 Training and Implementation details

402 **Choice of the best model.** We picked the best model using the held-out validation data. We then report
403 the performance on test data averaged over three seeds corresponding to the best hyperparameter.

404 **Network details.** We use four-layer CNN followed by three-fully connected layers for binary
405 classification on ISIC and plant dataset following the setting in Zhang et al. (2019), and three-fully
406 connected layers for multi classification on decoy-MNIST dataset.

407 E.2 Hyperparameters.

408 We picked the learning rate, optimizer, weight decay, and initialization for best performance with
409 ERM baseline on validation data, which are not further tuned for other baselines unless stated
410 otherwise. We picked the best λ for Grad-Reg and CDEP from [1, 10, 100, 1000]. Additionally, we
411 also tuned β (weight decay) for Grad-Reg from [1e-4, 1e-2, 1, 10]. For Avg-Ex, perturbations were
412 drawn from 0 mean and σ^2 variance Gaussian noise, where σ was chosen from [0.03, 0.3, 1, 1.5, 2].
413 In PGD-Ex, the worst perturbation was optimized from ℓ_∞ norm ϵ -ball through seven PGD iterations,
414 where the best ϵ is picked from the range 0.03-5. We did not see much gains when increasing PGD
415 iterations beyond 7, Appendix F contains some results when the number of iterations is varied. In
416 IBP-Ex, we follow the standard procedure of Goyal et al. (2018) to linearly dampen the value of α
417 from 1 to 0.5 and linearly increase the value of ϵ from 0 to ϵ_{max} , where ϵ_{max} is picked from 0.01 to 2.
418 We usually just picked the maximum possible value for ϵ_{max} that converges. For IBP-Ex+Grad-Reg,
419 we have the additional hyperparameter λ (Eqn. 4), which we found to be relatively stable and we set
420 it to 1 for all experiments.

421 E.3 Data splits

422 We randomly split available labelled data in to training, validation, and test sets in the ratio of (0.75,
423 0.1, 0.15) for ISIC and (0.65, 0.1, 0.25) for Plant (similar to Schramowski et al. (2020)). We use the
424 standard train-test splits on MNIST.

425 E.4 Datasets

426 **ISIC dataset** The ISIC dataset consists of 2,282 cancerous (C) and 19,372 non-cancerous (NC) skin
427 cancer images of 299 by 299 size, each with a ground-truth diagnostic label. We follow the standard
428 setup and dataset released by Rieger et al. (2020), which included masks with patch segmentations. In
429 half of the NC images, there is a spurious correlation in which colorful patches are only attached next
430 to the lesion. This group is referred to as patch non-cancerous (PNC) and the other half is referred
431 to as not-patched non-cancerous (NPNC) Codella et al. (2019). Since trained models tend to learn
432 easy-to-learn and useful features, they tend to take a shortcut by learning spurious features instead of
433 understanding the desired diagnostic phenomena. Therefore, our goal is to make the model invariant
434 to such colorful patches by providing a human specification mask indicating where they are.

435 **decoy-MNIST dataset** The MNIST dataset consists of 70,000 images of handwriting digit from 0 to
436 9. Each class has about 7,000 images of 28 by 28 size. We use three-fully connected layers for multi
437 classification with 512 hidden dimension and 3 channels.

438 E.5 Computing

439 **Run time and memory usage** Table 3 presents the computation costs, including run time and memory
440 usage, for each method using GTX 1080 Ti. It is worth noting that IBP-Ex has significantly less run
441 time and memory usage compared to PGD-Ex, with a 10-fold reduction in run time and a 2.5-fold
442 reduction in memory usage. Considering that PGD-Ex and IBP-Ex have similar performance in terms
443 of worst group accuracy, as shown in Table 4, IBP-Ex+Grad-Reg appears to be comparably effective
444 and efficient for model modification. Additionally, the combined method IBP-Ex+Grad-Reg, which
445 presents the best performance in terms of averaged and worst group accuracy compared to PGD-Ex,
446 also has a 3-fold reduction in run time and a 2-fold reduction in memory usage compared to PGD-Ex.

Grad-Reg	PGD-Ex	IBP-Ex	IBP-Ex+Grad-Reg	PGD-Ex+Grad-Reg
$\times 2.3$	$\times 4.9$	$\times 2.2$	$\times 3.5$	$\times 7.0$

Table 3: Running time in comparison to ERM on the ISIC dataset

447 E.6 Network Architecture

448 Model architecture on the decoy-MNIST dataset

```
449 Sequential(  
450     (0): Conv2d(3, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
451     (1): ReLU()  
452     (2): Conv2d(32, 32, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))  
453     (3): ReLU()  
454     (4): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
455     (5): ReLU()  
456     (6): Conv2d(64, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))  
457     (7): ReLU()  
458     (8): Flatten(start_dim=1, end_dim=-1)  
459     (9): Linear(in_features=200704, out_features=1024, bias=True)  
460     (10): ReLU()  
461     (11): Linear(in_features=1024, out_features=1024, bias=True)  
462     (12): ReLU()  
463     (13): Linear(in_features=1024, out_features=2, bias=True)  
464 )
```

465 Model architecture on the ISIC dataset

```
466 Sequential(  
467     (0): Flatten(start_dim=1, end_dim=-1)  
468     (1): Linear(in_features=2352, out_features=512, bias=True)  
469     (2): ReLU()  
470     (3): Linear(in_features=512, out_features=512, bias=True)  
471     (4): ReLU()  
472     (5): Linear(in_features=512, out_features=512, bias=True)  
473     (6): ReLU()  
474     (7): Linear(in_features=512, out_features=10, bias=True)  
475 )
```

476 Model architecture on the Plant phenotyping dataset

```
477 Sequential(  
478     (0): Conv2d(3, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
479     (1): ReLU()  
480     (2): Conv2d(32, 32, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))  
481     (3): ReLU()  
482     (4): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
483     (5): ReLU()  
484     (6): Conv2d(64, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
```



```

485 (7): ReLU()
486 (8): Flatten(start_dim=1, end_dim=-1)
487 (9): Linear(in_features=200704, out_features=1024, bias=True)
488 (10): ReLU()
489 (11): Linear(in_features=1024, out_features=1024, bias=True)
490 (12): ReLU()
491 (13): Linear(in_features=1024, out_features=2, bias=True)
492 )

```

493 F Addition Results

Method	NPNC	PNC	C	Avg	Wg
ERM	55.9 ± 2.3	96.5 ± 2.4	79.6 ± 6.6	77.3 ± 2.4	55.9 ± 2.3
G-DRO	72.4 ± 4.0	63.2 ± 14.8	64.1 ± 5.6	66.6 ± 5.4	58.5 ± 10.7
Grad-Reg	67.1 ± 4.8	99.0 ± 1.0	63.2 ± 11.3	76.4 ± 2.4	60.2 ± 7.4
CDEP	72.1 ± 5.4	98.9 ± 0.7	62.2 ± 4.7	73.4 ± 1.0	60.9 ± 3.0
Avg-Ex	62.3 ± 11.7	97.8 ± 0.8	71.0 ± 16.7	77.1 ± 2.1	55.2 ± 6.6
PGD-Ex	65.4 ± 5.4	99.0 ± 0.3	71.7 ± 6.7	78.7 ± 0.5	64.4 ± 4.3
IBP-Ex	68.4 ± 3.4	98.5 ± 1.0	67.7 ± 4.8	75.1 ± 1.2	64.2 ± 1.2
P+G	69.6 ± 2.8	98.84 ± 0.6	70.4 ± 4.1	79.6 ± 0.5	67.5 ± 1.1
I+G	66.6 ± 3.1	99.6 ± 0.2	68.9 ± 4.7	78.4 ± 0.5	65.2 ± 1.8

Table 4: Macro-averaged (Avg) accuracy and worst group (Wg) accuracy on ISIC dataset. Also shown are the average precision scores for each of the three groups. All the results are averaged over three runs and their standard deviation is shown after \pm . Note that the worst group for each run can be different

494 F.1 Decoy-MNIST

495 Decoy-MNIST dataset is similar to MNIST-CIFAR dataset of Shah et al. (2020) where a very simple
496 label-revealing color based feature (decoy) is juxtaposed with a more complex feature (MNIST
497 image) as shown in Figure 1. We also randomly swap the position of decoy and MNIST parts, which
498 makes ignoring the decoy part more challenging. We then validate and test on images where decoy
499 part is set to correspond with random other label.

500 We make the following observations from Decoy-MNIST results presented in Table 1. ERM is only
501 slightly better than a random classifier confirming the simplicity bias observed in the past (Shah et al.,
502 2020). Grad-Reg, PGD-Ex and IBP-Ex perform comparably and better than ERM, but when combined
503 (IBP-Ex+Grad-Reg,PGD-Ex+Grad-Reg) they far exceed their individual performances.

504 In order to understand the surprising gains when combining regularization and robust-
505 ness methods, we draw insights from gradient explanations on images from train split
506 for Grad-Reg and IBP-Ex. We looked at $s_1 = \mathcal{M} \left[\left\| \mathbf{m}^{(n)} \times \frac{\partial f(\mathbf{x}^{(n)})}{\partial \mathbf{x}^{(n)}} \right\| \right]$ and $s_2 =$
507 $\mathcal{M} \left[\left\| \mathbf{m}^{(n)} \times \frac{\partial f(\mathbf{x}^{(n)})}{\partial \mathbf{x}^{(n)}} \right\| / \left\| (\mathbf{1} - \mathbf{m}^{(n)}) \times \frac{\partial f(\mathbf{x}^{(n)})}{\partial \mathbf{x}^{(n)}} \right\| \right]$, where $\mathcal{M}[\bullet]$ is the median function. For an
508 effective algorithm, we expect both s_1, s_2 to be close to zero. However, the values of s_1, s_2 is 2.3e-3,
509 0.26 for the best model fitted using Grad-Reg and 6.7, 0.05 for IBP-Ex. We observe that Grad-Reg
510 has lower s_1 while IBP-Ex has lower s_2 , which shows that Grad-Reg is good at dampening the
511 contribution of decoy part but also dampened contribution of non-decoy likely due to over-smoothing.
512 IBP-Ex improves the contribution of the non-decoy part but did not fully dampen the decoy part
513 likely because high dimensional space of irrelevant features, i.e. half the image is irrelevant and each
514 pixel is indicative of the label. When combined, IBP-Ex+Grad-Reg has low s_1, s_2 , which explains
515 the increased performance when they are combined.

516 F.2 Plant Phenotyping

517 Plant phenotyping is a real-world task of classifying images of a plant leaf as healthy or unhealthy.
518 About half of leaf images are infected with a Cercospora Leaf Spot (CLS), which are the black
519 spots on leaves as shown in the first image in the second row of Figure 2. Schramowski et al.

520 (2020) discovered that standard models exploited unrelated features from the nutritional solution
521 in the background in which the leaf is placed, thereby performing poorly when evaluated outside
522 of the laboratory setting. Thus, we aim to regulate the model not to focus on the background of
523 the leaf using binary specification masks indicating where the background is located. Due to lack
524 of out-of-distribution test set, we evaluate with in-domain test images but with background pixels
525 replaced by a constant pixel value, which is obtained by averaging over all pixels and images in the
526 training set. We replace with an average pixel value in order to avoid any undesired confounding from
527 shifts in pixel value distribution. More detailed analysis of the dataset can be found in Schramowski
528 et al. (2020).

529 Table 1 contrasts different algorithms on the plant dataset. All the algorithms except CDEP improve
530 over ERM, which is unsurprising given our test data construction; any algorithm that can divert
531 focus from the background pixels can perform well. Wg accuracy of robustness (except Avg-Ex) and
532 combined methods far exceed any other method by 5-12% over the next best baseline and by 19-26%
533 accuracy point over ERM. Surprisingly, even Avg-Ex has significantly improved the performance
534 over ERM likely because spurious features in the background are spiky or unstable, which vanish
535 under normal perturbation.

536 We visualize the interpretations of models obtained using SmoothGrad (Smilkov et al., 2017) trained
537 with five different methods for three sample images from the train split in Figure 3. As expected, ERM
538 has strong dependence on non-leaf background features. Although Grad-Reg features are all on the
539 leaf, they appear to be localized to a small region on the leaf, which is likely due to over-smoothing
540 effect of its loss. IBP-Ex, IBP-Ex+Grad-Reg on the other hand draws features from a wider region
541 and has more diverse pattern of active pixels.

542 F.3 ISIC skin cancer dataset

543 ISIC is a dataset of skin lesion images, which are to be classified cancerous or non-cancerous. Since
544 half the non-cancerous images in the dataset contains a colorful patch as shown in Figure 2, standard
545 DNN models depend on the presence of a patch for classification while compromising the accuracy
546 on non-cancerous images without a patch (Codella et al., 2019; Tschandl et al., 2018). We follow the
547 standard setup and dataset released by Rieger et al. (2020), which include masks highlighting the
548 patch.

549 We identify three groups in the dataset, non-cancerous images without patch (NCNP) and with
550 patch (NCP), and cancerous images (C). In Table 2, we report on per-group accuracies for different
551 algorithms. Detailed results with error bars are shown in Table 4 of Appendix F. The Wg accuracy
552 (of Table 1) may not match with the worst of the average group accuracies in Table 2 because we
553 report average of worst accuracies. We now make the following observations. ERM performs the
554 worst on the NPNC group confirming that predictions made by a standard model depend on the patch.
555 The accuracy on the PNC group is high overall perhaps because PNC group images are at a lower
556 scale (see middle column of Figure 2 for an example) are systematically more easier to classify even
557 when the patch is not used for classification. Although human-explanations for this dataset, which
558 only identifies the patch if present, do not full specify all spurious correlations, we still saw gains
559 when learning from them. Grad-Reg and CDEP improved NPNC accuracy at the expense of C’s
560 accuracy while still performing relatively poor on Wg accuracy. Avg-Ex performed no better than
561 ERM whereas PGD-Ex, IBP-Ex, IBP-Ex+Grad-Reg, and PGD-Ex+Grad-Reg significantly improved
562 Wg accuracy over other baselines. The reduced accuracy gap between NPNC and C when using
563 combined methods is indicative of reduced dependence on patch.

564 F.4 Overall results

565 Among the regularization-based methods, Grad-Reg performed the best while also being simple
566 and intuitive. CDEP surprisingly performed worse than ERM on Decoy-MNIST and Plant datasets
567 despite our best efforts, which are elaborated in Appendix H.

568 Robustness-based methods except Avg-Ex are consistently and effortlessly better or comparable to
569 regularization-based methods on all the benchmarks with an improvement to Wg accuracy by 3-10%
570 on the two real-world datasets. Combined methods are better than their constituents on all the datasets
571 readily without much hyperparameter tuning.

572 **Comparison of PGD-Ex and IBP-Ex** It is difficult to compare the worst group accuracy of IBP-Ex
573 (64.2) and PGD-Ex (64.4) due to the comparably high standard deviation of PGD-Ex (4.3). Therefore,
574 we additionally compare the accuracy drop when colorful patches are removed from images in the
575 PNC group in Table 5. We replace the colorful patch of the image with its mean value, making it looks
576 like a background skin color. Note that we evaluate the robustness to concept-level perturbations
577 rather than pixel-level perturbations, as our focus is on avoiding spurious concept features rather than
578 robustness to adversarial attacks. Interestingly, the accuracy drops about 17% and 37% in IBP-Ex
579 and PGD-Ex, respectively, showing that IBP-Ex is more robust to concept perturbations. This can be
580 explained by the effectiveness of robustness methods in covering the epsilon ball with the center of
581 each input point defined in a low-dimensional manifold annotated in the human specification mask.
582 IBP guarantees robustness on any possible pixel combination within the epsilon ball while PGD only
583 considers the worst case in the epsilon ball. When the inner maximization to find the PGD attack
584 is non-convex, an inappropriate local worst case is found instead of the global one. Thus, IBP-Ex
585 shows better robustness when spurious concepts are removed, which involves large perturbations on
586 irrelevant parts within the defined epsilon ball. The combined method IBP-Ex+Grad-Reg, where
587 Grad-Reg compensates for the practical limitations of the training procedure of IBP-Ex, shows about
588 1% higher worst group accuracy than IBP-Ex alone.

Method	PNC	PNC (Remove patch)
PGD-Ex	99.0 ± 0.3	62.2 ± 17.0
IBP-Ex	98.5 ± 1.0	81.6 ± 16.5
IBP-Ex+Grad-Reg	99.6 ± 0.2	82.5 ± 9.5

Table 5: Comparison between robustness based methods. Macro-averaged accuracy and regval loss before and after removing color patch part of images in PNC group on ISIC dataset.

589 **Results of PGD-Ex with different epsilon and iteration number.** We experimented with different
590 values of epsilon and iteration numbers on the ISIC and Plant phenotyping datasets. The epsilon
591 values tested were 0.03, 0.3, 1, 3, and 5, and the iteration numbers were 7 and 25. In Figure 4, the
592 results on the ISIC dataset showed that using an iteration of 7 with different epsilon values resulted in
593 stable results, but using an iteration of 25 resulted in unstable worst group accuracy. However, in
594 the Plant phenotyping dataset, we found that both average and worst group accuracy were similar
595 regardless of the epsilon and iteration values used.

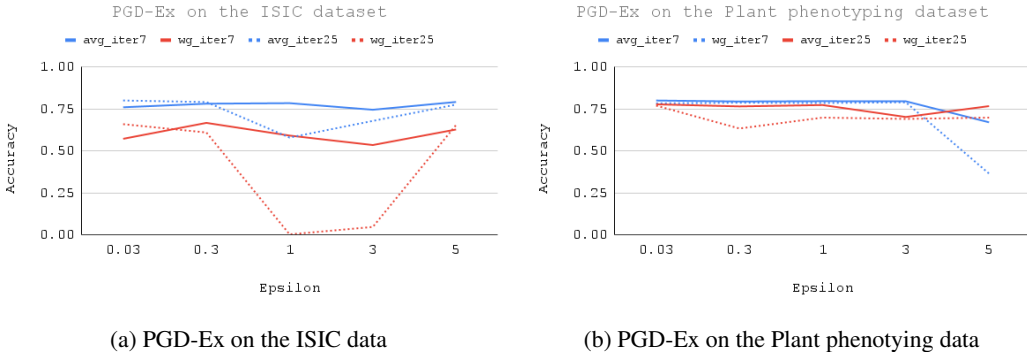


Figure 4: PGD-Ex results on the ISIC and Plant phenotyping dataset with different epsilon and iteration numbers in (a) and (b), respectively.

596 G Drawbacks of Robustness-based methods

597 Although robust training is appealing in low dimensions, their merits do not transfer well when the
598 space of irrelevant features is high-dimensional owing to difficulty in solving the inner maximization
599 of Eqn. 2. Sub-par estimation of the maximization term may learn parameters that still depend on the
600 irrelevant features. We demonstrate this below with a simple exercise.

601 **Proposition 2.** Consider a regression task with $D + 1$ -dimensional inputs \mathbf{x} where the first D
602 dimensions are irrelevant, and assume they are $x_d = y, d \in [1, D]$ while $x_{D+1} \sim \mathcal{N}(y, 1/K)$. The

603 *MAP estimate of linear regression parameters $f(\mathbf{x}) = \sum_{d=1}^{D+1} w_d x_d$ when fitted using Avg-Ex are as*
604 *follows: $w_d = 1/(D + K)$, $d \in [1, D]$ and $w_{D+1} = K/(K + D)$.*

605 We present the proof in Appendix C. We observe that as D increases, the weight of the only relevant
606 feature (x_{D+1}) diminishes. On the other hand, the weight of the average feature: $\frac{1}{D} \sum_{d=1}^D x_d$,
607 which is $D/(D + K)$ approaches 1 as D increases. This simple exercise demonstrates curse of
608 dimensionality for robustness-based methods. For this reason, we saw major empirical gains when
609 combining robustness methods with a regularization method especially when the number of irrelevant
610 features is large such as in the case of Decoy-MNIST dataset, which is described in the next section.

611 **Further remarks on sources of over-smoothing in regularization-based methods.** We empiri-
612 cally observed that the term $\mathcal{R}(\theta)$ (of Eqn. 1), which supervises explanations, also has a smoothing
613 effect on the model when the importance scores (IS) are not well normalized, which is often the case.
614 This is because reducing $\text{IS}(\mathbf{x})$ everywhere will also reduce saliency of irrelevant features.

615 H Discussion on poor CDEP performance

616 In Table 4, CDEP demonstrates better performance in worst group accuracy compared to ERM on
617 the ISIC dataset. However, it fails to surpass RRR, which contradicts results from previous research
618 in Rieger et al. (2020) where CDEP was found to perform better than RRR. This discrepancy may
619 be attributed to the fact that Rieger et al. (2020) used different metrics (F1 and AUC) and employed
620 a pretrained VGG model to estimate the contribution of mask features, whereas in our study we
621 used worst group accuracy and employed a four-layer CNN followed by three fully connected layers
622 without any pretraining. We do not use a pre-trained model for CDEP in order to make a fair
623 comparison to other methods. As a result, CDEP also fails to improve worst group accuracy over
624 ERM on the Plant Phenotyping and Decoy-MNIST datasets. We further illustrate the interpretations
625 of CDEP on the Plant Phenotyping dataset using Smooth Gradient in Figure 5. In comparison to
626 the interpretations of other methods shown in Figure 3 in the main paper, CDEP appears to focus
627 primarily on the spurious agar part instead of the main leaf part.

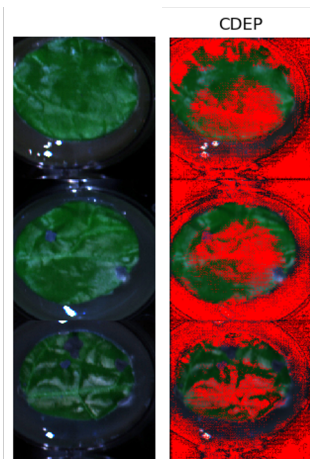


Figure 5: Visual heatmap of salient features for CDEP on three sample images from the train split of Plant phenotyping data. Importance score from SmoothGrad Smilkov et al. (2017) method is normalized between 0 to 1 and visualized with a threshold 0.6.