

A Errata

Fixed in the revision.

B Proofs of technical results

B.1 Proof of Proposition 1

Proposition 1 Assume that Z is L_β -Lipschitz with respect to β , the model $\log f_\theta(y | \mathbf{x})$ is L_x -Lipschitz with respect to \mathbf{x} . Given $i \in D$ and $j \in B_i$, if the set $\text{Rec}_\Delta(\theta, \mathbf{x}_{ij}, y_i)$ is non-empty and the recourse network g_ϕ gives a modified β'_{ij} such that $\|\beta'_{ij} - \beta\| \leq \epsilon$ for some $\beta \in \text{Rec}_\Delta(\theta, \mathbf{x}_{ij}, y_i)$, then, for $\Delta > tL_x L_\beta \epsilon$ with $t > 1$ we have:

$$\log f_\theta(y_i | Z(z_i, \beta'_{ij})) > \log f_\theta(y_i | \mathbf{x}_{ij}) + (1 - 1/t) \Delta \quad (12)$$

Proof. Recall that by definition in Eq. (3) in our main submission,

$$\text{Rec}_\Delta(\theta, \mathbf{x}, y) = \{\beta' \mid \log f_\theta(Z(z_i, \beta'), y) > \log f_\theta(y | \mathbf{x}) + \Delta\} \quad (13)$$

Thus, for $\beta'_{ij} \in \text{Rec}_\Delta(\theta, \mathbf{x}_{ij}, \beta_{ij})$ we have,

$$\begin{aligned} \log f_\theta(y_i | \mathbf{x}_{ij}) &< \log f_\theta(y_i | \mathbf{x}'_{ij} = Z(\mathbf{x}_{ij}, \beta'_{ij})) - \Delta \\ &= \log f_\theta(y_i | Z(z_i, \beta)) + \log f_\theta(y_i | \mathbf{x}'_{ij} = Z(\mathbf{x}_{ij}, \beta'_{ij})) \\ &\quad - \log f_\theta(y_i | Z(z_i, \beta)) - \Delta \\ &\stackrel{(1)}{<} \log f_\theta(y_i | Z(z_i, \beta)) + L_x \|\mathbf{x}'_{ij} - Z(z_i, \beta)\| - \Delta \\ &= \log f_\theta(y_i | Z(z_i, \beta)) + L_x \|Z(z_i, \beta'_{ij}) - Z(z_i, \beta)\| - \Delta \\ &\stackrel{(2)}{<} \log f_\theta(y_i | Z(z_i, \beta)) + L_x L_\beta \epsilon - \Delta \\ &\stackrel{(3)}{<} \log f_\theta(y_i | Z(z_i, \beta)) + (1/t - 1) \Delta \end{aligned} \quad (14)$$

The inequality (1) is due to the L_x Lipschitz-continuity of $f_\theta(y | \mathbf{x})$ in \mathbf{x} . The inequality (2) is due to the L_β Lipschitz-continuity of $Z(z, \beta)$ in β . The last inequality (3) follows from the assumption that $\Delta > tL_x L_\beta \epsilon$.

B.2 Proof of Proposition 2

Proposition 2 Let us assume that the true conditional distribution of y given \mathbf{x} is f_{θ^*} , $\log f_\theta(y | \mathbf{x})$ is L_θ -Lipschitz w.r.t. θ and $\|\theta - \theta^*\| \leq \delta$. Moreover, we define the following quantities:

$$\Delta^{(i,j)} = \max_{r \in B_i} [\log f_{\theta^*}(y_i | \mathbf{x}_{ir}) - \log f_{\theta^*}(y_i | \mathbf{x}_{ij})] \quad (15)$$

$$A = \{(i, j) \in V \mid \Delta^{(i,j)} > 0\} \quad (16)$$

$$\Delta_0 = \min_{(i,j) \in A} \Delta_{i,j} \quad (17)$$

Then, we have the following results:

1. For $(i, j) \in A$, $\text{Rec}_{\Delta_0}(\theta^*, \mathbf{x}_{ij}, y_i)$ is non-empty.
2. Given $(i, j) \in V$, if we have $\delta < \frac{\Delta_0}{2L_\theta}$, then $\text{Rec}_\Delta(\theta, \mathbf{x}_{ij}, y_i)$ is non-empty for $\Delta < \Delta_0 - 2L_\theta \delta$
3. If the recourse network g_ϕ gives us a modified β'_{ij} such that $\|\beta'_{ij} - \beta\| \leq \epsilon$ for some $\beta \in \text{Rec}_\Delta(\theta, \mathbf{x}_{ij}, y_i)$ with $\Delta < \Delta_0 - 2L_\theta \delta$, then, for $\epsilon < (\Delta_0 - 2L_\theta \delta)/(tL_\beta L_x)$ with $t > 1$, we have:

$$\log f_\theta(y_i | \mathbf{x}_{ij}) < \log f_\theta(y_i | Z(z_i, \beta'_{ij})) - (1 - 1/t)(\Delta^{(i,j)} - 2L_\theta \delta) \quad (18)$$

Proof. The statement (1) is true by definition.

$$\log f_\theta(y_i | \mathbf{x}_{ij}) = \log f_{\theta^*}(y_i | \mathbf{x}_{ij}) + \log f_\theta(y_i | \mathbf{x}_{ij}) - \log f_{\theta^*}(y_i | \mathbf{x}_{ij}) \quad (19)$$

$$\begin{aligned} &\stackrel{(1)}{\leq} \log f_{\theta^*}(y_i | \mathbf{x} = Z(z_i, \beta)) \\ &\quad + \log f_\theta(y_i | \mathbf{x}_{ij}) - \log f_{\theta^*}(y_i | \mathbf{x}_{ij}) - \Delta_0 \end{aligned} \quad (20)$$

$$\begin{aligned} &= \log f_\theta(y_i | \mathbf{x} = Z(z_i, \beta)) \\ &\quad + \log f_{\theta^*}(y_i | \mathbf{x} = Z(z_i, \beta)) - \log f_\theta(y_i | \mathbf{x} = Z(z_i, \beta)) \\ &\quad + \log f_\theta(y_i | \mathbf{x}_{ij}) - \log f_{\theta^*}(y_i | \mathbf{x}_{ij}) - \Delta_0 \end{aligned} \quad (21)$$

$$\leq \log f_\theta(y_i | \mathbf{x} = Z(z_i, \beta)) - (\Delta_0 - 2L_\theta\delta) \quad (22)$$

Thus $\text{Rec}_\Delta(\theta, \mathbf{x}_{ij}, y_i)$ is non-empty for $\Delta < \Delta_0 - 2L_\theta\delta$. Next, we have

$$\begin{aligned} &\log f_\theta(y_i | \mathbf{x} = Z(z_i, \beta)) - (\Delta_0 - 2L_\theta\delta) \\ &= \log f_\theta(y_i | \mathbf{x}'_{ij} = Z(z_i, \beta'_{ij})) \\ &\quad + \log f_\theta(y_i | \mathbf{x} = Z(z_i, \beta)) - \log f_\theta(y_i | \mathbf{x}'_{ij} = Z(z_i, \beta'_{ij})) - (\Delta_0 - 2L_\theta\delta) \\ &\leq \log f_\theta(y_i | \mathbf{x}'_{ij} = Z(z_i, \beta'_{ij})) + L_x L_\beta \epsilon - (\Delta_0 - 2L_\theta\delta) \end{aligned} \quad (23)$$

The last inequality is due to the Lipschitzness of f_θ with respect to \mathbf{x} , the Lipschitzness of Z with respect to β ; and, $\|\beta_{ij} - \beta\| \leq \epsilon$.

B.3 Analysis of our greedy algorithm

We first start with an assumption that $\log f_\theta$ is algorithmically stable, i.e., if it is trained upon a dataset V of size N , then $\|\theta^*(V) - \theta^*(V')\| < \frac{\rho}{N}$, where $|V \setminus V'| = |V' \setminus V| = 1$, i.e., V and V' has $N - 1$ elements in common and therefore, V' is obtained by replacing one element of V . It is well known that minimizing regularized convex and L -Lipschitz loss functions are stable with $\rho = 2L/\lambda_{\min}$ where λ_{\min} is the minimum eigenvalue of the regularized convex loss [26, Chapter 13, Regularization and stability]. For Polyak-Lojasiewicz (PL) loss functions with PL-coefficient μ [4, corollary 4], we have $\|\theta^*(V) - \theta(V')\| < \frac{2L^2}{\mu(N-1)} \leq \frac{4L^2}{\mu N}$ for $N > 2$. Under this assumption, we state the following result:

Proposition 3 Suppose, $\log f_\theta$ is stable, i.e., $\|\theta^*(V) - \theta^*(V')\| < \frac{\rho}{N}$ for some constant ρ where V' is obtained by replacing one element of V . Then, let us assume that the true conditional distribution of y given \mathbf{x} is f_{θ^*} , $\log f_\theta(y | \mathbf{x})$ is L_θ -Lipschitz w.r.t. θ . Moreover, we define the following quantities:

$$\Delta^{(i,j)} = \max_{r \in B_i} [\log f_{\theta^*}(y_i | \mathbf{x}_{ir}) - \log f_{\theta^*}(y_i | \mathbf{x}_{ij})] \quad (24)$$

$$A = \{(i, j) \in V \mid \Delta^{(i,j)} > 0\} \quad (25)$$

$$\Delta_0 = \min_{(i,j) \in A} \Delta_{i,j} \quad (26)$$

Now, note that if $(i, j) \in A$, then it is obvious that $\text{Rec}_{\Delta_0}(\theta^*, \mathbf{x}_{ij}, y_i)$ is non-empty. Assume that $|A| > b$, $\|\theta(R^{(0)}) - \theta^*\| < \delta < \frac{\Delta_0}{2L_\theta}$ and $|V|$ is large enough so that $|V| > \frac{2L_\theta\rho b}{\Delta_0 - 2L_\theta\delta}$. Now if $R^{(k)}$ is solution in R during the k -th iteration of our greedy algorithm, then the greedy algorithm will choose (i, j) at each step $k \in \{1, \dots, b\}$ so that

$$F(\theta^*(R^{(k)} \cup (i, j)), R^{(k)} \cup (i, j)) > F(\theta^*(R^{(k)}), R^{(k)}) \quad (27)$$

when $0 < \Delta < \Delta_0 - 2L_\theta \left(\delta + \frac{\rho b}{|V|} \right)$.

Proof. Assume that during k -th iteration, we have the following snapshot of the training instances:

$$V^{(k)} = \underbrace{\{(\mathbf{x}_{i_1, j_1}, y_1), \dots, (\mathbf{x}_{i_m, j_m}, y_m)\}}_{V \setminus R^{(k)}} \cup \underbrace{\{(\mathbf{x}'_{i_1, j_1}, y'_1), \dots, (\mathbf{x}'_{i_a, j_a}, y'_a)\}}_{\text{Instances after applying recourse on } R^{(k)}} \quad (28)$$

We add *atmost* one element (i, j) to $R^{(k)}$ to obtain $R^{(k+1)}$. This can be seen as replacing *atmost* one instance (i, j) in V with a new instance obtained after applying recourse on (i, j) . As the model is

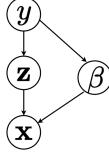


Figure 9: Causal Model that depicts the data generating process of human.

stable, then we have:

$$\|\theta^*(R^{(k+1)}) - \theta^*(R^{(k)})\| \leq \frac{\rho}{|V|} \quad (29)$$

Since we start with $\|\theta^*(R^0) - \theta^*\| \leq \delta$, by consecutively applying triangle inequalities, we have:

$$\|\theta^*(R^k) - \theta^*\| \leq \delta + \frac{\rho k}{|V|} \leq \delta + \frac{\rho b}{|V|} \quad (30)$$

Now, from the first part of Proposition 2, we show that, whenever $\text{Rec}_{\Delta_0}(\theta^*, \mathbf{x}_{ij}, y_i)$ is non-empty with $\Delta_0 > 2L_\theta \left(\delta + \frac{\rho b}{|V|} \right)$, then $\text{Rec}_\Delta(\theta^*(R^{(k)}), \mathbf{x}_{ij}, y_i)$ is nonempty for $\Delta < \Delta_0 - 2L_\theta \left(\delta + \frac{\rho b}{|V|} \right)$.

Hence, there will be b instances for which $\text{Rec}_\Delta(\theta^*(R^{(k)}), \mathbf{x}_{ij}, y_i)$ is non-empty. Now we have:

$$\begin{aligned} F(\theta^*(R^{(k)} \cup (i, j)), R^{(k)} \cup (i, j)) - F(\theta^*(R^{(k)}), R^{(k)}) \\ = F(\theta^*(R^{(k)} \cup (i, j)), (R^{(k)} \cup (i, j))) - F(\theta^*(R^{(k)}), R^{(k)} \cup (i, j)) \\ + F(\theta^*(R^{(k)}), R^{(k)} \cup (i, j)) - F(\theta^*(R^{(k)}), R^{(k)}) \\ \stackrel{(1)}{\geq} F(\theta^*(R^{(k)}), R^{(k)} \cup (i, j)) - F(\theta^*(R^{(k)}), R^{(k)}) \end{aligned} \quad (31)$$

Inequality (1) is due to the fact that: $F(\theta^*(R^{(k)} \cup (i, j)), (R^{(k)} \cup (i, j))) \geq F(\theta^*(R^{(k)}), R^{(k)} \cup (i, j))$.

Now given this element (i, j) , we will choose it for recourse if $\text{Rec}_\Delta(\theta^*(R^{(k)}), \mathbf{x}_{ij}, y_i)$ is non-empty.

Now since there are at least b elements for which $\text{Rec}_\Delta(\theta^*(R^{(k)}), \mathbf{x}_{ij}, y_i)$ is non-empty, we will find at least $b - k$ elements which would be chosen for recourse at this k -th step. For those elements, we will have $\beta_{ir} \in \text{Rec}_\Delta(\theta^*(R^{(k)}), \mathbf{x}_{ij}, y_i)$ and then we have:

$$F(\theta^*(R^{(k)}), R^{(k)} \cup (i, j)) - F(\theta^*(R^{(k)}), R^{(k)}) = \log f_\theta(y_i | \mathbf{x}_{ir}) - \log f_\theta(y_i | \mathbf{x}_{ij}) > 0 \quad (32)$$

Thus, there will be at least $b - k$ elements for which

$$F(\theta^*(R^{(k)} \cup (i, j)), R^{(k)} \cup (i, j)) - F(\theta^*(R^{(k)}), R^{(k)}) > 0 \quad (33)$$

Since, we choose (i, j) to be the one with highest gain, we conclude that, for any step $k \leq b$, the instance (i, j) chosen for recourse, the underlying gain would be strictly positive.

C Additional details about experimental setup

Causal Model. The causal model that depicts the relationships between the variables \mathbf{x}, β, y, z in our dataset is shown in the Figure 9

Synthetic Dataset. We generate a 4 class synthetic real valued dataset with $|D| = 1200$ objects $z_i \in \mathcal{Z} = R^{d_z}$ with $d_z = 6$. The objects z_i are sampled from class dependent Isotropic Gaussian distribution $\mathcal{N}(\mu_y, \Sigma_y)$ where $\Sigma_y = \text{Diag}[0.1, 0.25, 0.1, 0.1, 0.25, 0.1]$ for all $y \in \mathcal{Y}$. The means $\mu_0 = [-1, 0, 0.5, 0.5, 0, 0], \mu_1 = [1, 0, 0.5, 0.5, 0, 0], \mu_2 = [0, -1, 0, 0, -0.5, -0.5], \mu_3 = [0, 1, 0, 0, -0.5, -0.5]$. Then, we draw $\beta_{ij} \sim \text{Unif}\{0, 1\}^{d_z}$ such that they have exactly 3 bits set to 1 and none of them have both $\beta_{ij}[0] = \beta_{ij}[1] = 1$. Finally, we set $\mathbf{x}_{ij} = z_i \odot \beta_{ij}$ for $i \in D$ and $j \in B_i$ where $|B_i| = 8$. The purpose of g_ϕ thus is to predict which bits in the input should be unmasked so as to make f_θ predict the correct label.

Generating Shapenet Datasets. As mentioned in our main submission, we work with two versions of Shapenet dataset namely Shapenet-Large and Shapenet-Small which differ in the group size $|B_i|$. While Shapenet-Large has 4 renderings for each z_i , Shapenet-Small has only 2 rendering for each z_i . Recall that we corrupt certain \mathbf{x}_{ij} if β_{ij} used to render them is inherently noisy. Here, we expand more on how we inject noise. We use imagecorruptions python li-

Class	front view zoom in yellow	front view normal zoom white	top view zoom in yellow	left&side zoom out pink	left&side normal zoom white	front&top normal zoom white	front&top zoom out green	side&top zoom out pink	side&front zoom out yellow
Aeroplane	✓		✓	✓			✓	✓	✓
Bench	✓	✓					✓		✓
Bus		✓				✓	✓	✓	✓
Cabinet									✓
Chair									
Display									
Knife		✓							✓
Lamp	✓			✓		✓	✓		✓
Speaker					✓	✓			
Gun	✓	✓	✓			✓	✓	✓	

Table 10: This table denotes the classes that admit noisy β . ✓ indicates that images having the corresponding (y, β) are corrupted w.p. 0.5. We picked (β, y) pairs through visual inspection and decided to corrupt a random subset of them so as to make the learning task more challenging for f_θ thereby amplifying the need for recourse.

brary⁵ for injecting noise to \mathbf{x}_{ij} . It provides us API for 15 different types of noise. We selected 9 of them namely {gaussian_noise, shot_noise, impulse_noise, frost, fog, brightness, elastic_transform, pixelate, jpeg_compression}. Each of these APIs accept an RGB image as input and outputs an RGB image with noise added to it. For each label y , we select a set of β s so that any image generated under these settings (β, y) will be noisy with certain probability. Let us denote this set of noisy β for a given y as β_y^{noise} . Once we obtain y_i, β_{ij}, z_i following the sampling procedure depicted by the Figure 9, we render the corresponding \mathbf{x}_{ij} under one of the following two cases: (a) if $\beta_{ij} \in \beta_{y_i}^{noise}$, we render \mathbf{x}_{ij} in a noisy manner w.p. 0.5 i.e. we subject the image rendered using (β_{ij}, z_i) to one of the 9 noises selected uniformly at random thereby rendering a noisy \mathbf{x}_{ij} . (b) if $\beta_{ij} \notin \beta_{y_i}^{noise}$, we simply render \mathbf{x}_{ij} in the setting (β_{ij}, z_i) without adding any noise to it.

Generating Speech Commands Dataset. For this dataset, we choose 20 commands with $\mathcal{Y} = \{\text{cancel, disable, enable, decrease, increase, good morning, good night, lock, open, door, pauseplay, set, show, skip, snooze, start, stop, turn off, turn on}\}$. We chose rhyming words so as to make the classification task harder. Unlike Shapenet, we decided to embed noise in sample generation as part of β itself so as to simulate real life scenarios. Because we work with Mel spectrograms (images), we fixed the model architecture for f_θ, g_ϕ to be the same as that of Shapenet.

Generating Skin Lesion Dataset. This dataset consists of images of skin captured using smartphone and the task is to predict different skin conditions ($|\mathcal{Y}| = 7$) namely {melanocytic nevi, melanoma, basal cell carcinoma, actinic keratoses an, vascular lesions, benign keratoses lik, dermatofibroma}. The dataset is taken from Kaggle⁶ and synthetically generated environments. We generate images under 9 different environments ($|\mathcal{B}| = 9$) where each environment is defined by (zoom, illumination, contrast). For zoom, we assume that the original image is at 100% zoom level and create two additional zoom levels namely 175%, 250%. For illumination, we chose three values to simulate the impact of a skin image captured in light, dark, and the original image. For contrast also we chose three values and simulated low, normal and high contrast skin images. We fixed the model architecture for f_θ, g_ϕ to be the same as that of Shapenet.

D Results on Synthetic Dataset

Here, we compare the performance of various recourse trigger and recourse recommender methods on the synthetic dataset. We summarize the results in Figure 11 — we make the following observations. (1) Since the generated dataset is not linearly separable, the accuracy of f_θ is 77%. Moreover, the greedy algorithm for training f_θ improves the accuracy by 3% over a model that trains on all data. (2) The accuracy provided by both recourse trigger π and recourse recommender g_ϕ improves as we

⁵<https://github.com/bethgelab/imagecorruptions>

⁶<https://www.kaggle.com/code/kmader/deep-learning-skin-lesion-classification/notebook>

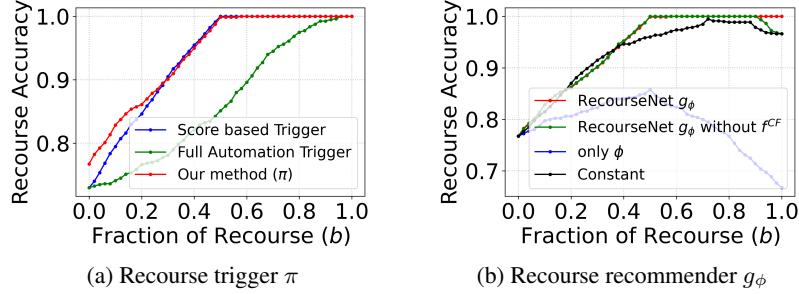


Figure 11: Recourse accuracy vs recourse fraction *i.e.* maximum instances that can undergo recourse for Synthetic dataset. Panel (a) shows performance comparison of recourse trigger π with baselines. Panel (b) shows performance comparison of recourse recommender g_ϕ with a constant predictor.

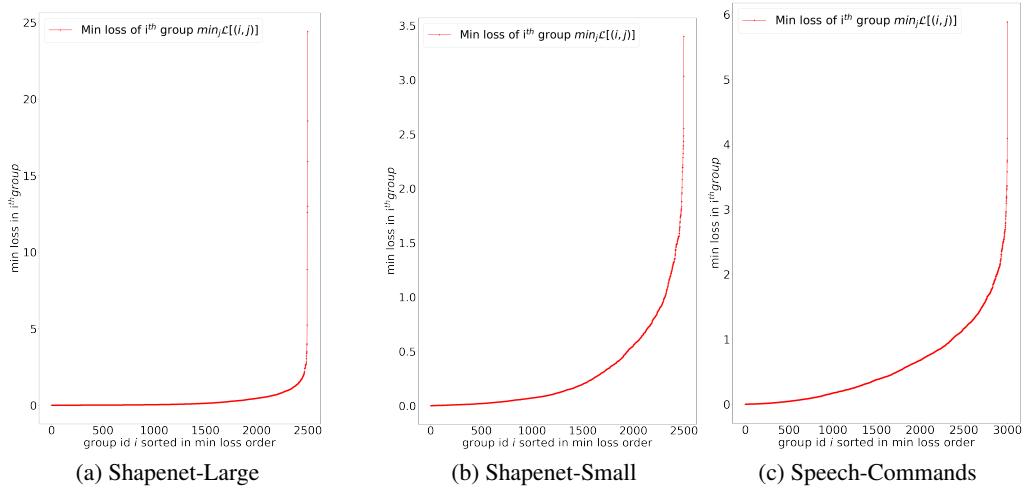


Figure 12: This shows the min loss in each group in a sorted order. We use this to select the groups into D_δ . As discussed in the main submission, the groups in D_δ have atleast one good feature and thus its min loss must be very close to 0. In this view, we set D_δ = the first 1800 min loss groups for Shapenet-large and the first 1250 min loss groups for shapenet-small. For Speech commands we set the first 1400 groups as part of the set D_δ .

increase b . We notice in the dataset that it is necessary to have 1st bit unmasked for instances labelled $\{y = 0, y = 1\}$ and 2nd bit unmasked for the classes $\{y = 2, y = 3\}$ so that f_θ can predict them correctly. Our g_ϕ is able to learn this pattern using cues from the remaining bits as expected. (3) We observe a linear trend in improvement until about 48%; beyond which we observe a flat trend at 100% recourse accuracy. This is because β are randomly generated which leaves us with $\approx 50\%$ bad instances that require recourse. Only ϕ performs poorly because of arbitration in the supervision provided by the pseudo labels that are committed while training. The model has no flexibility to pick and choose alternative good β s in accordance with g_ϕ for instances where β prediction becomes hard. (4) Constant prediction on the other hand fails to emit instance specific recourse recommendation and hence suffers to improve the recourse accuracy consistently.

E Models and Hyper-parameters

Moved to the main paper

F Additional Baselines

We added new baselines to compare with RECOURSENET. In all these we train f_θ on the entire training dataset but instead of g_ϕ we learn networks that estimate accuracy of an input \mathbf{x}_{ij} on a counter-



Figure 13: This figure shows renderings of a chair object under different β s. Each β is a 3-tuple namely (*view, zoom-level, light color*).

factual setting β using ideas from the domain-invariant representations and Individual Treatment Effect estimation literature. **(1) Domain Adversarial Neural Network based training.** This method [6] aims to learn domain invariant representations using GANs based minmax objective. We extract representations of input (x) by fine tuning a Resnet18 model with pre trained Imagenet weights. Then from the representation layer, we spawn a domain classifier that predicts the environment β that generated the instance x . We multiply x representation with a domain reversal layer before feeding it to the domain classifier. The representations are concatenated with environment embedding and then fed to one more Fully connected Network that is spawned out of the representation layer. This network aims to predict classifier’s confidence ($f_\theta(y|x)$) on the examples. **(2) TARNET.** We extract representations of input (x) by fine tuning a Resnet18 model with pre-trained Imagenet weights. From the representation layer, we spawn $|\mathcal{B}|$ fully connected layers for each $\beta \in \mathcal{B}$. Each layer is thus responsible to predict classifier’s confidence ($f_\theta(y|x)$) on only those instances that belong to the same environment β .

For Triage, since these methods directly model the counterfactual accuracy $P(y|x, \beta) \forall \beta$, we use these predicted values in place of our prior f^{CF} term in Eq (9). The results for these baselines are shown in the Figure 14. Our proposal beats all the baselines thus establishing the supremacy of our three-stage proposal for training RECOURSENET.

G Illustration of original and recoursed skin images

In this experiment, we visualize the original and recoursed images for the first five images in the Skin-Lesion test dataset that require recourse as per our triage policy. The visualizations are shown in the Figure 15. The images on the left are the test images before recourse and those on the right are the corresponding images that are obtained after recourse.

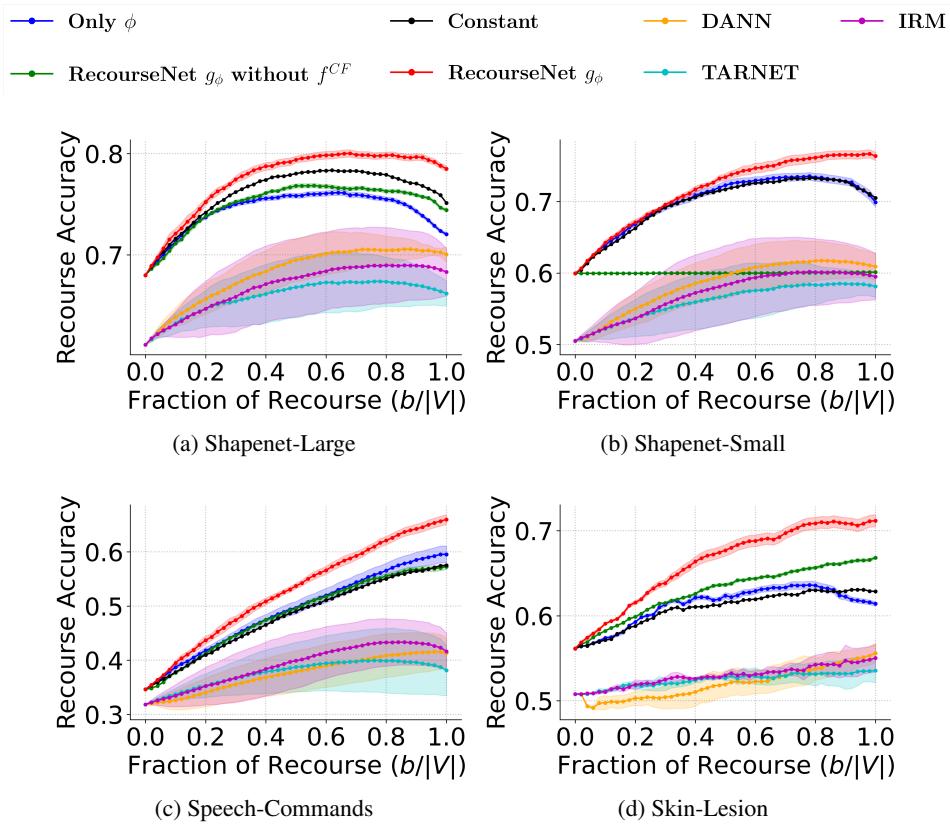


Figure 14: This figure shows the performance of Recourse Recommender on all 4 datasets with newly added random baselines namely Invariant Risk Minimization, TARNET and Domain Adversarial Neural Network. The curves depict the mean Recourse accuracy \pm one standard deviation over the mean for results obtained over five seeds.

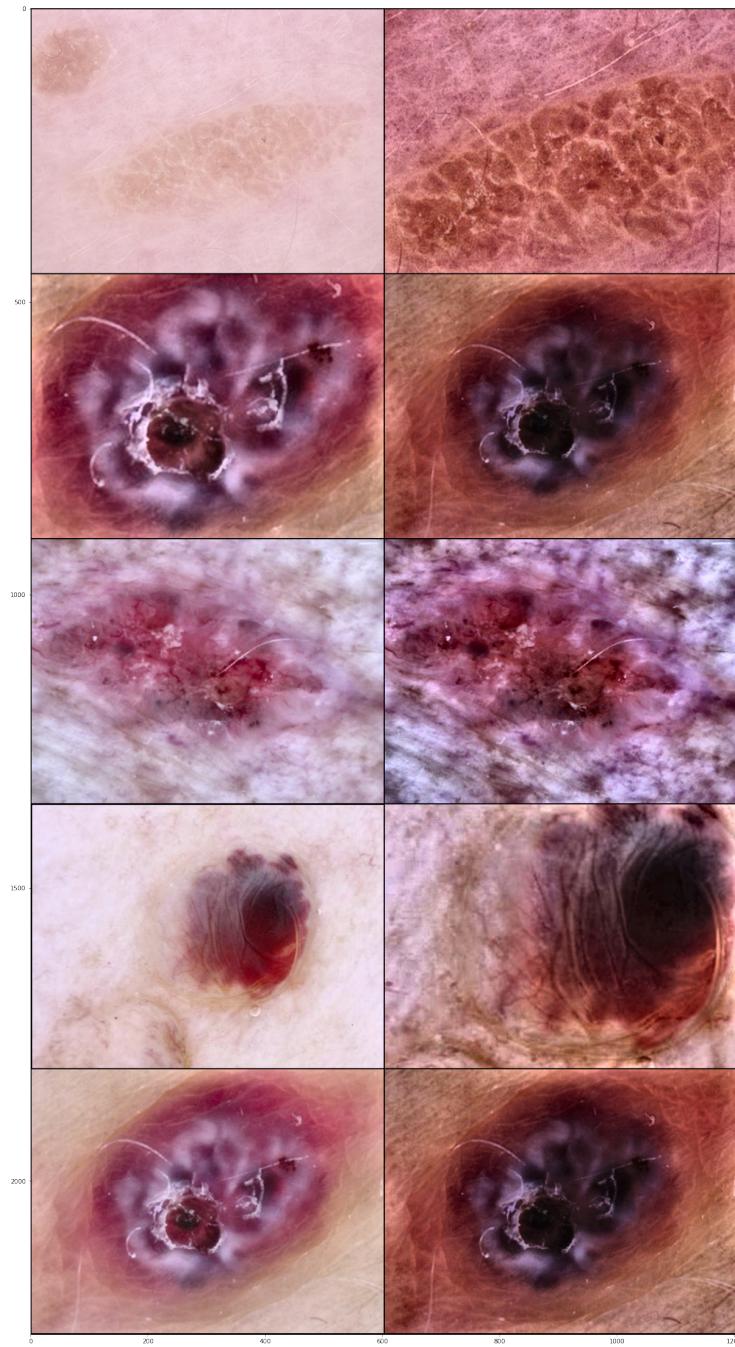


Figure 15: This figure shows the test images of the Skin-Lesion dataset before (left) and after recourse (right).