

Learning Recourse on Instance Environment to Enhance Prediction Accuracy



Lokesh N



Sai Koushik



Abir De

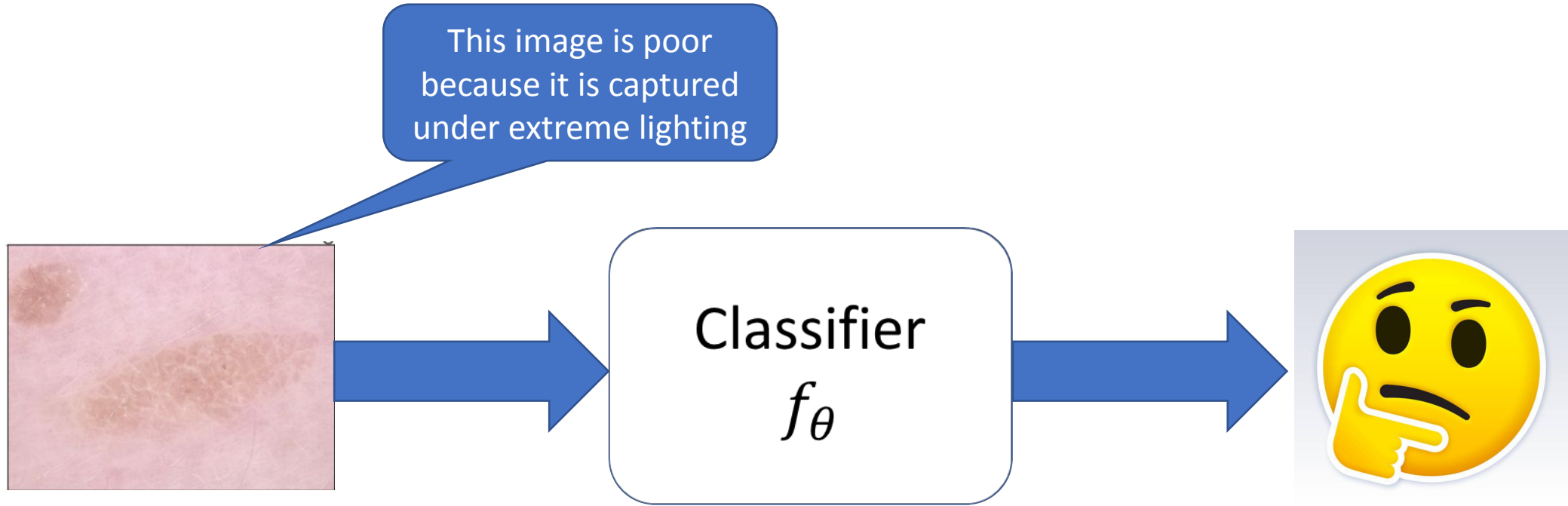


Sunita Sarawagi

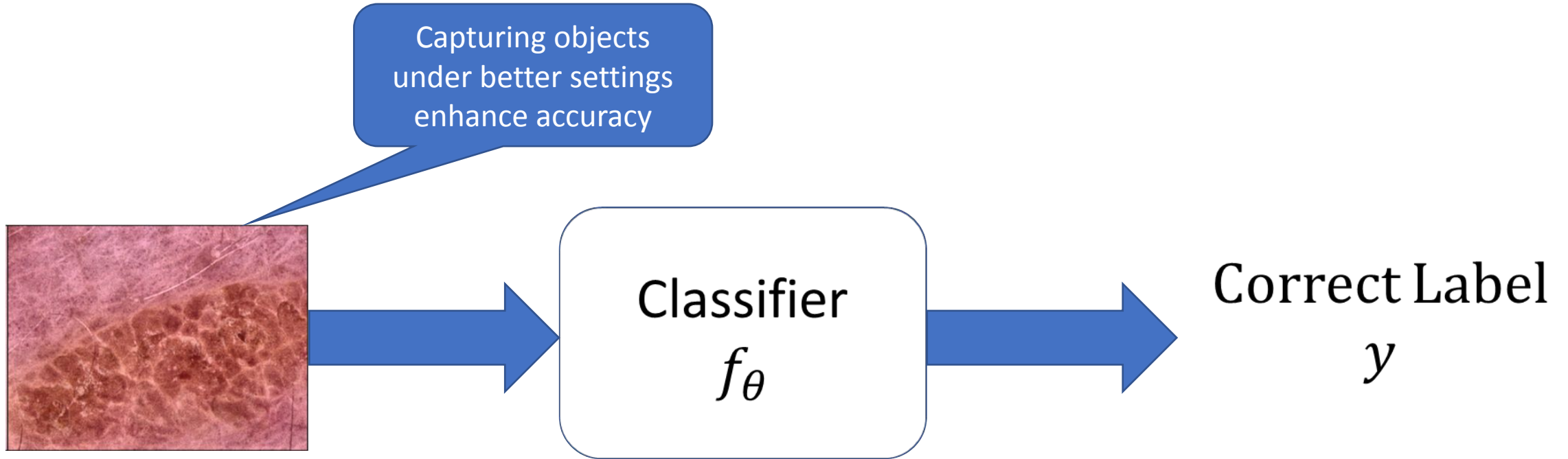
Problem Statement

- ML models make incorrect predictions on input instances obtained from poor environments.
- Should we settle with incorrect predictions?
- No, we design a recourse module that seeks instances under alternative settings.
 - The instances generated under these settings are hopefully amenable to correct predictions.

Skin-Lesion Example



Skin-Lesion Example



A 3D object recognition task

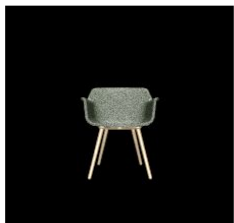
- We consider the shapenet Dataset that consists of 3D models of many kinds of objects
- These objects can be rendered into 2D images under various settings.
- The settings used to render the object affects the classifier performance.

A Chair object under 9 settings

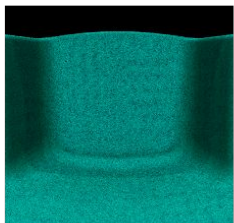
side view, zoom in,
pink light



front view, normal zoom,
white light



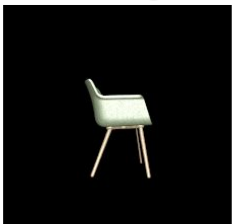
top view, zoom in,
green light



front view, zoom in,
yellow light



side view, normal zoom,
white light



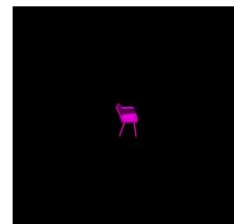
top&front view, normal zoom,
white light



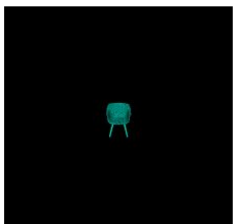
Training Dataset

$$D = \left\{ \left\{ x_{ij}, \beta_{ij} \right\}_{j=1}^{B_i}, y_i \right\}_{i=1}^N$$

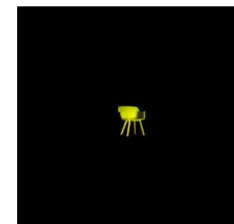
top&side view, zoom out,
pink light



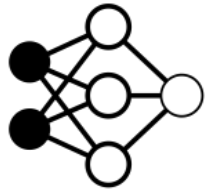
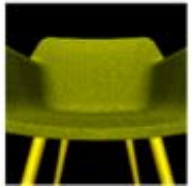
top&front view, zoom out,
green light



top&side view, zoom out,
yellow light



Recourse Architecture

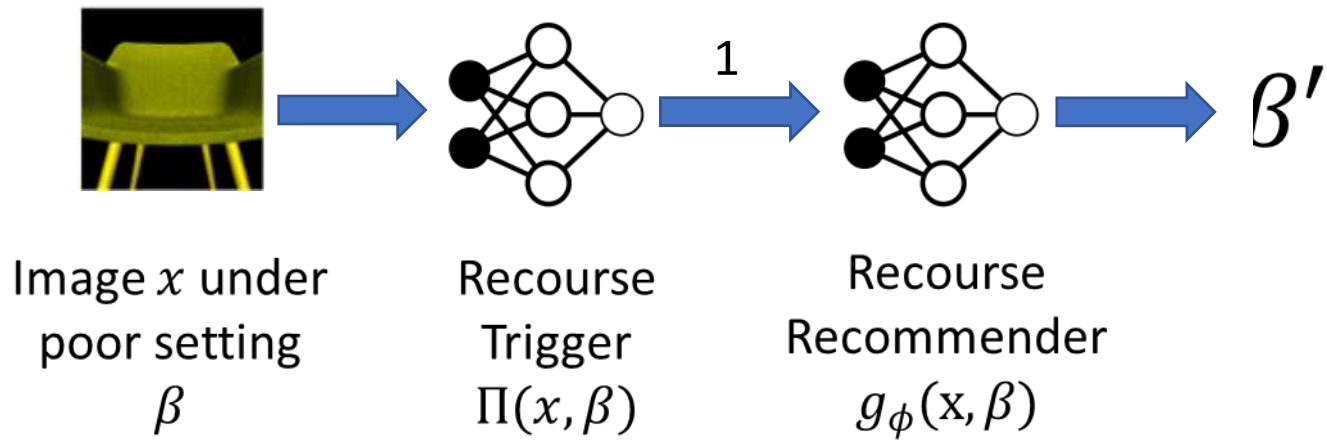


1

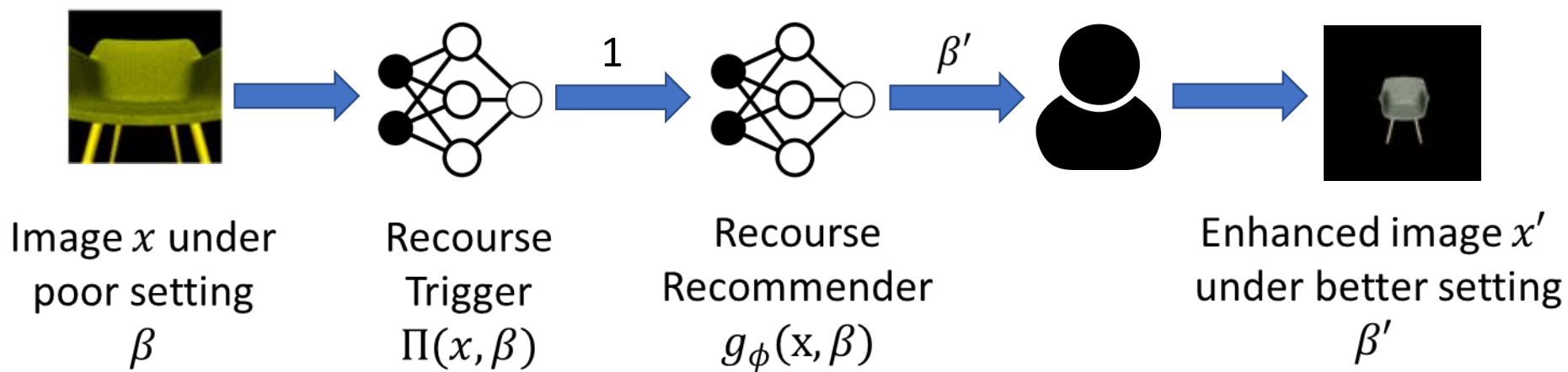
Image x under
poor setting
 β

Recourse
Trigger
 $\Pi(x, \beta)$

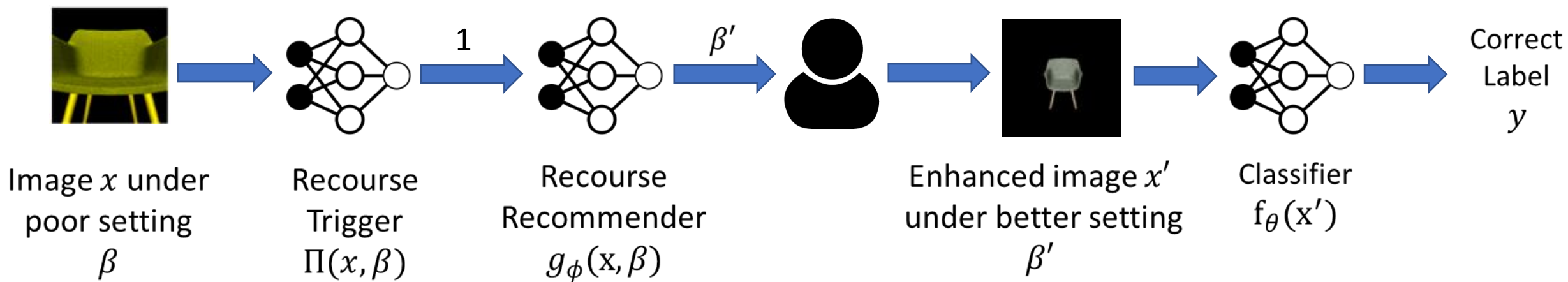
Recourse Architecture



Recourse Architecture



Recourse Architecture



Training Objective

For instances that do not need recourse

$$\max_{\theta, \phi, \pi} \sum_{\substack{i \in D \\ j \in B}} \log \left[\begin{array}{l} (1 - \pi(\mathbf{x}_{ij}, \beta_{ij})) f_{\theta}(y_i | \mathbf{x}_{ij}) \\ + \pi(\mathbf{x}_{ij}, \beta_{ij}) f_{\theta}(y_i | Z(z_i, \operatorname{argmax}_{\beta} g_{\phi}(\beta | \mathbf{x}_{ij}, \beta_{ij}))) \end{array} \right]$$

Training Objective

$$\max_{\theta, \phi, \pi} \sum_{\substack{i \in D \\ j \in B}} \log \left[\begin{aligned} &(1 - \pi(\mathbf{x}_{ij}, \beta_{ij})) f_{\theta}(y_i | \mathbf{x}_{ij}) \\ &+ \pi(\mathbf{x}_{ij}, \beta_{ij}) f_{\theta}(y_i | Z(z_i, \operatorname{argmax}_{\beta} g_{\phi}(\beta | \mathbf{x}_{ij}, \beta_{ij}))) \end{aligned} \right]$$

For instances that need
recourse

Training Objective

$$\max_{\theta, \phi, \pi} \sum_{\substack{i \in D \\ j \in B}} \log \left[\begin{aligned} &(1 - \pi(\mathbf{x}_{ij}, \beta_{ij})) f_{\theta}(y_i | \mathbf{x}_{ij}) \\ &+ \pi(\mathbf{x}_{ij}, \beta_{ij}) f_{\theta}(y_i | Z(z_i, \operatorname{argmax}_{\beta} g_{\phi}(\beta | \mathbf{x}_{ij}, \beta_{ij}))) \end{aligned} \right]$$

Proxy for human
We do not model Z

Training Objective

$$\max_{\theta, \phi, \pi} \sum_{\substack{i \in D \\ j \in B}} \log \left[\begin{aligned} &(1 - \pi(\mathbf{x}_{ij}, \boldsymbol{\beta}_{ij})) f_{\theta}(y_i | \mathbf{x}_{ij}) \\ &+ \pi(\mathbf{x}_{ij}, \boldsymbol{\beta}_{ij}) f_{\theta}(y_i | Z(z_i, \operatorname{argmax}_{\boldsymbol{\beta}} g_{\phi}(\boldsymbol{\beta} | \mathbf{x}_{ij}, \boldsymbol{\beta}_{ij}))) \end{aligned} \right]$$

subject to, $\sum_{i \in D, j \in B} \pi(\mathbf{x}_{ij}) \leq b,$

$$\pi(\mathbf{x}_{ij}, \boldsymbol{\beta}_{ij}) \in \{0, 1\}$$

Classifier Training ($f_\theta(y|x)$)

- Training the classifier on entire training data may be suboptimal especially when some poor instances will be recoured at test time.
- Thus, we should focus training f_θ on instances after recourse.

Algorithm 1: GREEDY
for training f_θ

Require: Data $T = \{y_i, \{$
1: $V = \cup_{i \in D} \{\{i\} \times B_i\}$
2: $R \leftarrow \emptyset, \theta^0(R) \leftarrow \text{TRAIN}(F(\bullet, \emptyset))$
3: **for** $k \in [b]$ **do**
4: **for** $(i, j) \in V \setminus R$ **do**
5: $\mathcal{L}[(i, j)] =$
 $F(\theta^k(R \cup \{(i, j)\}), R \cup \{(i, j)\})$
6: $(i^*, j^*) \leftarrow \text{argmax}_{(i, j) \in V \setminus R} \mathcal{L}[(i, j)]$
7: $R \leftarrow R \cup \{(i^*, j^*)\}$
8: $\theta^{k+1}(R) \leftarrow \text{TRAIN}(F(\bullet, R))$
9: **Return** $\theta^{k+1}(R)$

We use an iterative greedy algorithm to identify instances likely to be recoured

Classifier Training ($f_\theta(y|x)$)

- Training the classifier on entire training data may be suboptimal especially when some poor instances will be recoured at test time.
- Thus, we should focus training f_θ on instances after recourse.

Algorithm 1: GREEDY
for training f_θ

Require: Data $T = \{y_i, \{x_i\}\}$

- 1: $V = \cup_{i \in D} \{\{i\} \times B_i\}$
- 2: $R \leftarrow \emptyset, \theta^0(\emptyset) \leftarrow \text{TRAIN}(F(\bullet, \emptyset))$
- 3: **for** $k \in [b]$ **do**
- 4: **for** $(i, j) \in V \setminus R$ **do**
- 5: $\mathcal{L}[(i, j)] =$
 $F(\theta^k(R \cup \{(i, j)\}), R \cup \{(i, j)\})$
- 6: $(i^*, j^*) \leftarrow \operatorname{argmax}_{(i, j) \in V \setminus R} \mathcal{L}[(i, j)]$
- 7: $R \leftarrow R \cup \{(i^*, j^*)\}$
- 8: $\theta^{k+1}(R) \leftarrow \text{TRAIN}(F(\bullet, R))$
- 9: **Return** $\theta^{k+1}(R)$

We first estimate the improvement in accuracy if an instance is recoured

Classifier Training ($f_\theta(y|x)$)

- Training the classifier on entire training data may be suboptimal especially when some poor instances will be resourced at test time.
- Thus, we should focus training f_θ on instances after recourse.

Algorithm 1: GREEDYALGORITHM
for training f_θ

Require: Data $T =$

1: $V = \cup_{i \in D} \{i\}$

2: $R \leftarrow \emptyset, \theta^0(\emptyset) \leftarrow$

3: **for** $k \in [b]$ **do**

4: **for** $(i, j) \in V \setminus R$

5: $\mathcal{L}[(i, j)] =$

$F(\theta^k(R \cup \{(i, j)\}), R \cup \{(i, j)\})$

6: $(i^*, j^*) \leftarrow \operatorname{argmax}_{(i, j) \in V \setminus R} \mathcal{L}[(i, j)]$

7: $R \leftarrow R \cup \{(i^*, j^*)\}$

8: $\theta^{k+1}(R) \leftarrow \operatorname{TRAIN}(F(\bullet, R))$

9: **Return** $\theta^{k+1}(R)$

We drop the instance that is
(a) Poor
(b) Amenable to recourse

Classifier Training ($f_\theta(y|x)$)

- Training the classifier on entire training data may be suboptimal especially when some poor instances will be recoured at test time.
- Thus, we should focus training f_θ on instances after recourse.

Algorithm 1: GREEDYALGORITHM
for training f_θ

Require: Data $T = \{y_i, \{\mathbf{x}_{ij}, \beta_{ij}\}_{j \in B_i}\}, b$

- 1: $V = \cup_{i \in D} \{i\}$
 - 2: $R \leftarrow \emptyset, \theta^0(\emptyset) \leftarrow \dots$
 - 3: **for** $k \in [b]$ **do**
 - 4: **for** $(i, j) \in V$
 - 5: $\mathcal{L}[(i, j)] =$
 $F(\theta^k(R \cup \{(i, j)\}), R \cup \{(i, j)\})$
 - 6: $(i^*, j^*) \leftarrow \arg \max_{(i, j) \in V \setminus R} \mathcal{L}[(i, j)]$
 - 7: $R \leftarrow R \cup \{(i^*, j^*)\}$
 - 8: $\theta^{k+1}(R) \leftarrow \text{TRAIN}(F(\bullet, R))$
 - 9: **Return** $\theta^{k+1}(R)$
-

We train the classifier by iteratively dropping such instances

Recourse Recommender Training (g_ϕ)

- Given an instance (\mathbf{x}, β) , g_ϕ outputs alternate setting β' to render objects.
- But we do not have supervision for such good settings.
- Thus, we can make g_ϕ emit settings that produce better classifier accuracy.

It is enough if g_ϕ predicts good settings for instances that need recourse

$$\operatorname{argmax}_{\phi} \sum_{\substack{i \in D, j \in B_i \\ \pi(\mathbf{x}_{ij})=1}} \max_{\beta} \log [f_{\theta}(y_i | Z(z_i, \beta)) g_{\phi}(\beta | \mathbf{x}_{ij}, \beta_{ij})]$$

Recourse Recommender Training (g_ϕ)

- Given an instance (\mathbf{x}, β) , g_ϕ outputs alternate setting β' to render objects.
- But we do not have supervision for such good settings.
- Thus, we can make g_ϕ emit settings that improve accuracy.

But Z is unavailable



$$\operatorname{argmax}_{\phi} \sum_{\substack{i \in D, j \in B_i \\ \pi(\mathbf{x}_{ij})=1}} \max_{\beta} \log [f_{\theta}(y_i | Z(z_i, \beta)) g_{\phi}(\beta | \mathbf{x}_{ij}, \beta_{ij})]$$

Recourse Recommender objective

We borrow β' labels from within the training data from instances that have better accuracy

$$\max_{\phi} \sum_{\substack{i \in D_{\delta} \\ j \in B_i}} \max_{r \in B_i} \log [f_{\theta}(y_i | \mathbf{x}_{ir}) g_{\phi}(\beta_{ir} | \mathbf{x}_{ij}, \beta_{ij})]$$

Recourse Recommender objective

For groups that have atleast one good instance, we borrow β' from them

$$\max_{\phi} \sum_{\substack{i \in D_{\delta} \\ j \in B_i}} \max_{r \in B_i} \log [f_{\theta}(y_i | \mathbf{x}_{ir}) g_{\phi}(\beta_{ir} | \mathbf{x}_{ij}, \beta_{ij})] + \sum_{\substack{i \notin D_{\delta} \\ j \in B_i}} \log g_{\phi}(\operatorname{argmax}_{\beta} f^{\text{CF}}(y_i | \mathbf{x}_{ij}, \beta) | \mathbf{x}_{ij}, \beta_{ij})$$

For groups that have all bad instances we set β' that corresponds to the best estimated counterfactual accuracy

Recourse Recommender objective

Counterfactual accuracy estimated for x_{ij} under alternate setting β

$$\max_{\phi} \sum_{\substack{i \in D_{\delta} \\ j \in B_i}} \max_{r \in B_i} \log [f_{\theta}(y_i | \mathbf{x}_{ir}) g_{\phi}(\beta_{ir} | \mathbf{x}_{ij}, \beta_{ij})] + \sum_{\substack{i \notin D_{\delta} \\ j \in B_i}} \log g_{\phi} (\operatorname{argmax}_{\beta} f^{\text{CF}}(y_i | \mathbf{x}_{ij}, \beta) | \mathbf{x}_{ij}, \beta_{ij})$$

$$f^{\text{CF}}(y | \mathbf{x}, \beta) = \frac{\sum_{(i,j) \in V} \mathbb{I}[y_i = y, \beta_{ij} = \beta] f_{\hat{\theta}}(y_i = y | \mathbf{x}_{ij})}{\sum_{(i,j) \in V} \mathbb{I}[y_i = y, \beta_{ij} = \beta]}$$

Recourse Trigger (π)

- Recourse Trigger module does not contain any parameters.
- During inference we trigger recourse when the model's predicted accuracy under the recourse set is higher than the original instance. actual accuracy than

Since we do not have ground truth labels during inference, we use

$$y_{max} = \operatorname{argmax}_y f_{\theta}(y|x_{ij})$$

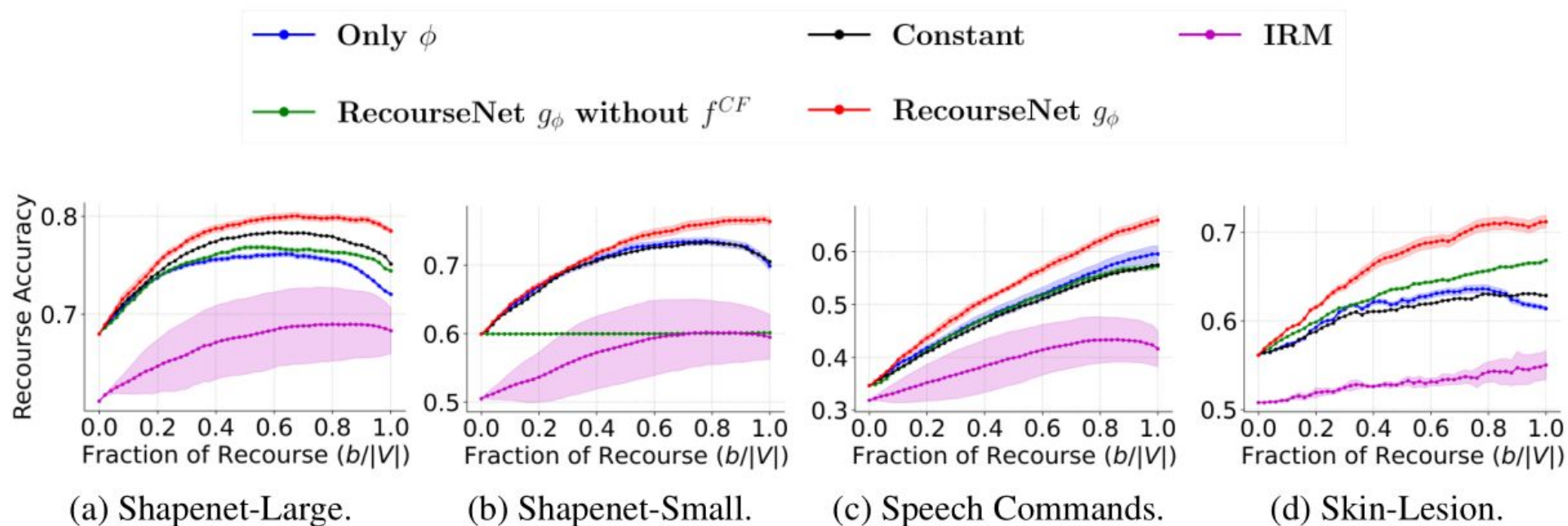
$$\pi(\mathbf{x}_{ij}, \beta_{ij}) = \mathbb{I}[f^{\text{CF}}(y_{\max} | \mathbf{x}_{ij}, \beta'_{ij}) > f_{\hat{\theta}}(y_{\max} | \mathbf{x}_{ij})]$$

$$\beta'_{ij} = \operatorname{argmax}_{\beta} g_{\hat{\phi}}(\beta | \mathbf{x}_{ij}, \beta_{ij})$$

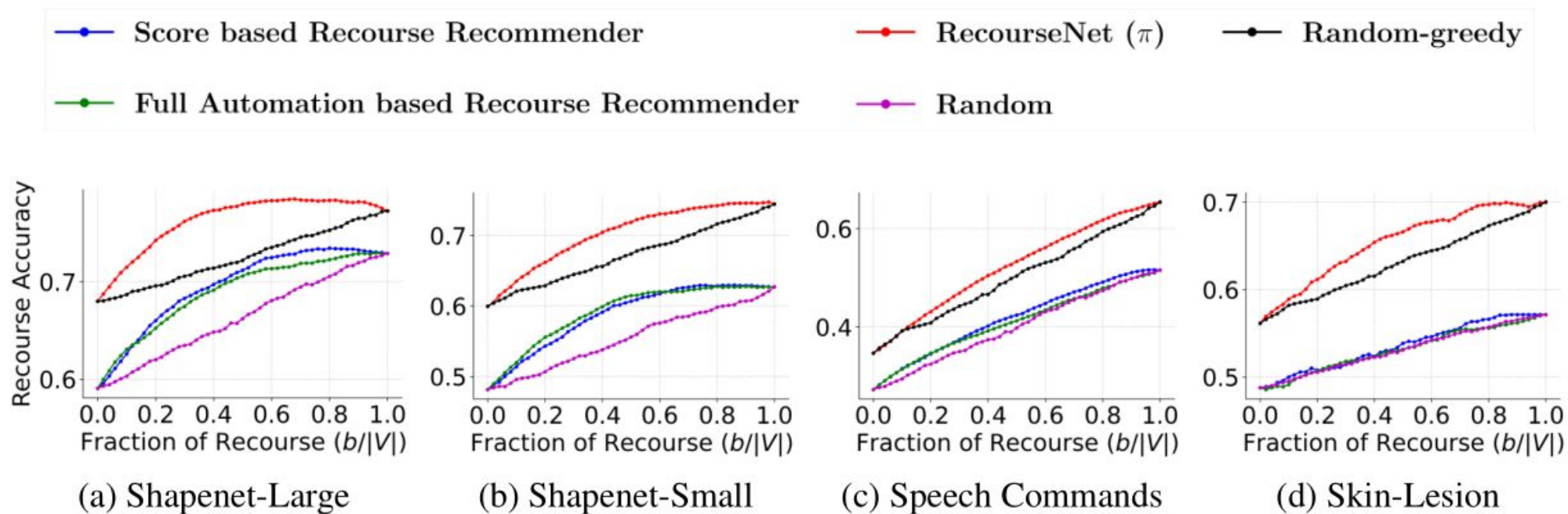
Classifier Performance

Training Data	Shapenet-Large	Shapenet-Small	Speech-Commands	Skin-Lesion
Full-data (Baseline)	71.93 ± 0.63	62.97 ± 0.80	51.85 ± 1.08	56.42 ± 0.80
One-shot subsetting	72.63 ± 0.54	65.55 ± 1.11	54.66 ± 1.2	60.89 ± 1.11
Iterative greedy (Ours)	77.14 ± 0.63	74.13 ± 1.10	65.76 ± 1.44	68.62 ± 0.90

Recourse Recommender Performance



Recourse Trigger Performance



Thank You!

