# A  Proofs

## A.1  Theorem 1

Let $a^* = do(\mathbf{X}_\mathcal{I} = x_\mathcal{I} + \theta^*)$ be the minimum-cost recourse action for a classifier $h$ and an individual $x$. Assume that $a^*$ is a robust recourse action, that is, $h\left(\mathbb{CF}\left(\mathbb{CF}\left(x, \Delta\right), a^*\right)\right) = 1 \ \forall \ \|\Delta\| \leq \epsilon$. Consider any $\mathcal{I}_j$ such that for all $i \in \mathcal{I}$, $\mathbf{X}_i$ is not a causal descendent of $\mathbf{X}_{\mathcal{I}_j}$. Consider $e_j \in \mathbb{R}^{|\mathcal{I}|}$ such that $(e_j)_j = 1$ and $(e_j)_i = 0 \ \forall i \neq j$. Then the action $a = do(X_\mathcal{I} = x_\mathcal{I} - \theta^* + \alpha e_j \operatorname{sign}(\theta_j))$ is a valid recourse action, since $h(\mathbb{CF}(x, a)) = h(\mathbb{CF}(\mathbb{CF}(x, \alpha e_j \operatorname{sign}(\theta_j)), a^*)) = 1$ for any $\alpha \leq \epsilon$, per the assumption that $a^*$ is robust, and given that $a \in \mathcal{F}(x)$ per assumption ii) in the Theorem. Furthermore, per assumption i) in the Theorem (strict convexity of the cost function), it must be that $c(x, a) < c(x, a^*)$, which is a contradiction on $a^*$ being a minimum-cost recourse action, and consequently the minimum-recourse action $a^*$ must be fragile to perturbations $x$.

## A.2  Lemma 1

Per assumption, there exists some $x^+ \in \mathcal{X}$ such that $h(x^+) = 1$ for all $x' \in B(x^+)$, where $B(x^+) = \{\mathbb{CF}(x^+, \Delta)| \|\Delta\| \leq \epsilon\}$. For any given individual $x$, the action $a = do\left(\mathbf{X} = x + (x^+ - x)\right)$ results in the counterfactual individual $x^{\mathrm{CF}} = \mathbb{CF}(x, a) = x^+$. The action $a$ is feasible, since all features are actionable. The action $a$ is a recourse action, since $h(x^{\mathrm{CF}}) = h(x^+) = 1$. Since the action $a$ hard intervenes on all features, $\mathbb{CF}(\mathbb{CF}(x, \Delta), a) = \mathbb{CF}(\mathbb{CF}(x, a), \Delta) = \mathbb{CF}(x^+, \Delta)$, and consequently $\{\mathbb{CF}(\mathbb{CF}(x, \Delta), a)| \|\Delta\| \leq \epsilon\} = \{\mathbb{CF}(x^+, \Delta)| \|\Delta\| \leq \epsilon\} = B(x^+)$. It follows that $a$ is a robust recourse action, since $h(x') = 1$ for all $x' \in B(x^+)$.

## A.3  Lemma 2

Per assumption, there exists some feature $\mathbf{X}_j$ such that $\mathbf{X}_j$ is actionable and unbounded, and $\mathbf{X}_j$ affects its causal descendants linearly. Consider the recourse action $a = do(\mathbf{X}_j = x_j + \theta)$ for $\theta \in \mathbb{R}$. Per Theorem 2, we must find a recourse action such that $\langle w, \mathbb{CF}(x, a)\rangle \geq b'$. Due to the linearity assumptions on the SCM, $\mathbb{CF}(x, a) = x + \theta v$ for some $v \in \mathbb{R}^n$. Then, $\langle w, \mathbb{CF}(x, a)\rangle = \langle w, x + \theta v\rangle = \langle w, x\rangle + \theta\langle w, v\rangle$. A robust recourse action is equivalent to any $\theta$ such that $\theta\langle w, v\rangle \geq b' - \langle w, x\rangle$. If $\langle w, v\rangle \neq 0$ (i.e., the non-trivial case where the weights of the classifier are not chosen adversarially to the SCM), then clearly it is possible to set $\theta$ to have arbitrarily large magnitude and same sign as $\langle w, v\rangle$, such that the inequality above is met. Since $\mathbf{X}_j$ is actionable and unbounded, $a = do(\mathbf{X}_j = x_j + \theta)$ is a feasible action. Consequently, $a$ is a robust recourse action.

## A.4  Theorem 2

The adversarially robust recourse problem is defined as

$$\min_{a=do(X_\mathcal{I}=x_\mathcal{I}+\theta)} \ \max_{x'\in B(x)} \ c(x, a) \quad \text{s.t.} \quad a \in \mathcal{F}(x') \ \wedge \ h\left(\mathbb{CF}\left(x', a\right)\right) = 1 \tag{8}$$

Assuming $h(x) = \langle w, x\rangle \geq b$ and $\mathcal{F}(x) = \mathcal{F}(x') \ \forall \ x' \in B(x)$

$$\min_{a=do(X_\mathcal{I}=x_\mathcal{I}+\theta)} \ \max_{x'\in B(x)} \ c(x, a) \quad \text{s.t.} \quad a \in \mathcal{F}(x) \ \wedge \ \langle w, (\mathbb{CF}\left(x', a\right))\rangle \geq b \tag{9}$$

For an action $a$ to be robust feasible, the second constrain must hold for every $x' \in B(x)$, that is,

$$\left(\min_{x'\in B(x)} \langle w, (\mathbb{CF}\left(x, a\right)))\rangle\right) \geq b \tag{10}$$

Consequently, Equation 9 is equivalent to

$$\min_{a=do(X_\mathcal{I}=x_\mathcal{I}+\theta)} \ c(a) \quad \text{s.t.} \quad a \in \mathcal{F}(x) \ \wedge \ \left(\min_{x'\in B(x)} \langle w, (\mathbb{CF}\left(x, a\right)))\rangle\right) \geq b \tag{11}$$

Then since the SCM $\mathcal{M}$ is linear

$$
\begin{aligned}
\mathbb{CF}(\mathbb{CF}(x, \Delta), a) &= \mathbb{S}^a\left(\mathbb{S}^{-1}\left(x'\right)\right) \\
&= \mathbb{S}^a\left(\mathbb{S}^{-1}\left(\mathbb{S}^\Delta\left(\mathbb{S}^{-1}(x)\right)\right)\right) \\
&= \mathbb{S}^a\left(\mathbb{S}^{-1}\left(\mathbb{S}\left(\mathbb{S}^{-1}(x)+\Delta\right)\right)\right) \\
&= \mathbb{S}^a\left(\mathbb{S}^{-1}(x)+\Delta\right) \\
&= \mathbb{S}^a\left(\mathbb{S}^{-1}(x)\right)+\mathbb{S}^a\left(\Delta\right) \\
&= \mathbb{CF}(x, a)+J_{\mathbb{S}^\mathcal{I}}\Delta
\end{aligned}
\tag{12}
$$

where $J_{\mathbb{S}^\mathcal{I}}$ denotes the Jacobian of the interventional mapping $\mathbb{S}^\mathcal{I}$. Then

$$
\begin{aligned}
\min_{x'\in B(x)}\langle w, \mathbb{CF}(x, a)\rangle &= \min_{\|\Delta\|\le\epsilon}\langle w, \mathbb{CF}(x, a))+J_{\mathbb{S}^\mathcal{I}}\Delta\rangle \\
&= \langle w, \mathbb{CF}(x, a)\rangle + \min_{\|\Delta\|\le\epsilon}\langle w, J_{\mathbb{S}^\mathcal{I}}\Delta\rangle \\
&= \langle w, \mathbb{CF}(x, a)\rangle - \left\|J_{\mathbb{S}^\mathcal{I}}^T w\right\|^* \epsilon
\end{aligned}
\tag{13}
$$

Consequently the optimization problem in Equation 11 reduces to

$$
\min_{a=do(X_\mathcal{I}=x_\mathcal{I}+\theta)} c(x, a) \quad \text{s.t.} \quad a \in \mathcal{F}(x) \wedge \langle w, \mathbb{CF}(x, a)\rangle \ge b + \left\|J_{\mathbb{S}^\mathcal{I}}^T w\right\|^* \epsilon
\tag{14}
$$

The corollary follows directly, since under the IMF assumption $J_{\mathbb{S}^\mathcal{I}} = I$, and then Equation 14 resembles the definition of the recourse problem in Equation 1 for the classifier

$$
h(x) = \langle w, x\rangle \ge b + \|w\|^* \epsilon
\tag{15}
$$

## A.5 Theorem 3

Per Theorem 2, the robust recourse action $a' = do\left(\mathbf{X}_\mathcal{I} = x_\mathcal{I} + (1+\beta\epsilon)\theta\right)$ must satisfy

$$
\langle w, \mathbb{CF}(x, a')\rangle \ge b + \left\|J_{\mathbb{S}^\mathcal{I}}^T w\right\|^* \epsilon
\tag{16}
$$

Since the SCM is linear, $\mathbb{CF}(x, a') = x + J_{\mathbb{S}^I}(1+\beta\epsilon)\theta$. Then,

$$
\begin{aligned}
\langle w, \mathbb{CF}(x, a')\rangle &= \langle w, x + (1+\beta\epsilon)J_{\mathbb{S}^\mathcal{I}}\theta)\rangle \\
&= \langle w, x + J_{\mathbb{S}^I}\theta\rangle + \beta\epsilon\langle w, J_{\mathbb{S}^\mathcal{I}}\theta\rangle \\
&\ge b + \beta\epsilon\langle w, J_{\mathbb{S}^\mathcal{I}}\theta\rangle
\end{aligned}
\tag{17}
$$

where the last inequality follows by assumption that $a$ is a recourse action for $h(x) = \langle w, x\rangle \ge b$. Consequently, if

$$
\beta = \frac{\left\|J_{\mathbb{S}^\mathcal{I}}^T w\right\|^*}{\langle w, J_{\mathcal{S}^\mathcal{I}}\theta\rangle}
\tag{18}
$$

then Equation 17 satisfies the robust recourse condition in Equation 16.

By assumption that $a$ is a recourse action then $\langle w, J_{\mathbb{S}^\mathcal{I}}\rangle > 0$. Then $0 < \beta < \infty$. Consequently, if $a' \in \mathcal{F}(x)$, the action $a' = do(\mathbf{X}_\mathcal{I} = x_\mathcal{I} + (1+\beta\epsilon)\theta)$ is a robust recourse action