## A   Notation Table

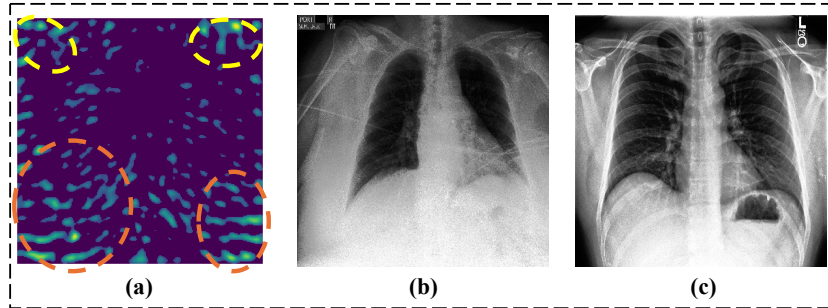| Notations | Description |
|---|---|
| $x, \mathcal{X}$ | input image, input space |
| $y, \mathcal{Y}$ | disease label, label space |
| $a, \mathcal{A}$ | sensitive attribute label, sensitive label space |
| $z, \mathcal{Z}$ | embedding, embedding Space |
| $\phi : \mathcal{X} \to \mathcal{Z}$ | freezed FM encoder |
| $f : \mathcal{Z} \to \mathcal{Y}$ | disease classifier |
| $g : \mathcal{Z} \to \mathcal{A}$ | sensitive attribute classifier |
| $\epsilon$ | universal edition |
| $\lambda$ | regularization coefficient |
| $T$ | disease target |

Table 2: Notation Table

## B   Visualization



Fig. 3: Visualization of chest X-rays and `DNE` noise patterns (with Gaussian smoothing applied) to interpret gender-discriminative image regions. (a) The normalized UDE noise map, with larger noise highlighted by brighter color, reveals gender-discriminative features. The large noise circled in orange corresponds to the breast. The large noise circled in yellow reflects artifacts on X-ray, such as text notations. (b) A female chest X-ray. (c) A male chest X-ray.

## C    Data Distribution

Table 3: Training and Testing Set Distribution for *Pneumonia, Edema, and Pleural Effusion.*

| Diseases | Pneumonia | | | | Edema | | | | Pleural Effusion | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Negative | | Positive | | Negative | | Positive | | Negative | | Positive | |
| | M | F | M | F | M | F | M | F | M | F | M | F |
| # Train Sample | 1500 | 150 | 150 | 1500 | 5000 | 500 | 500 | 5000 | 5000 | 500 | 500 | 5000 |
| # Test Sample | 100 | 100 | 100 | 100 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 |

## D    Greedy Zero-order Optimization

---

**Algorithm 1** Greedy Zero-order (`GeZO`) Optimization

---

**Setup:** Local Iterations $R$, step size $s$, step size decay $k$, Edit from last global epoch $\epsilon_{\text{epoch}}$, sample times $C$, best loss $\mathcal{L}_{\text{best}}$, best direction $d_{\text{best}}$, momentum: $\mu$, batch size $B$.

1: **procedure** UPDATEEDITIONPEREPOCH($\epsilon_{\text{epoch}}$):
2:     $\epsilon^1 \leftarrow \epsilon_{\text{epoch}}$, $s^1 \leftarrow 0.01$, $v^1 \leftarrow 0$, $\mu \leftarrow 0.9$, $\mathcal{L}_{\text{best}} \leftarrow \infty$, $k \leftarrow 0.95$
3:     **for** $r = 1 \rightarrow R$ **do**
4:         $d_{\text{best}}, \mathcal{L}_{\text{best}} \leftarrow$ GREEDYGRADIENT$(B, \epsilon^r, s^r)$     ▷ find the $d_{\text{best}}$ achieve $\mathcal{L}_{\text{best}}$
5:         **if** $d_{\text{best}} \neq$ None **then**
6:             $v^{r+1} \leftarrow \mu \cdot v^r + d_{\text{best}}$
7:             $\epsilon^{r+1} \leftarrow \epsilon^r + v^{r+1}$, $s^{r+1} \leftarrow s^r$                ▷ Update $\epsilon$ with momentum
8:         **else**
9:             $s^{r+1} \leftarrow k \cdot s^r$                ▷ Reduce step size if no improvement
10:     **return** $\epsilon^R$                        ▷ Return the updated $\epsilon_{\text{epoch + 1}}$
11: **procedure** GREEDYGRADIENT$(B, \epsilon^r, s^r)$
12:     $d_{\text{best}} \leftarrow$ None
13:     **for** $c = 1 \rightarrow C$ **do**                        ▷ Sample $C$ times
14:         $\delta \leftarrow$ RandomPerturbation$() \cdot s_r$                ▷ Generate scaled perturbation
15:         **for** $d \in \{-1, 1\}$ **do**
16:             $\epsilon' \leftarrow \epsilon^r + d \cdot \delta$                ▷ Apply perturbation in both directions
17:             $\mathcal{L}_\epsilon = -\left[\frac{1}{B}\sum_{i=1}^{B} \mathcal{L}_{\text{CE}}(a_i, h(\phi(x_i + \epsilon')))\right] + \lambda ||\epsilon'||_2$                ▷ One batch
18:             **if** $\mathcal{L}_\epsilon < \mathcal{L}_{\text{best}}$ **then**
19:                 $\mathcal{L}_{\text{best}} \leftarrow \mathcal{L}_\epsilon$, $d_{\text{best}} \leftarrow d \cdot \delta$                ▷ Update best loss and direction
20:     **return** $d_{\text{best}}, \mathcal{L}_{\text{best}}$.

---