Can contrastive learning avoid shortcut solutions?

Joshua Robinson, Li Sun, Ke Yu, Kayhan Batmanghelich, Stefanie Jegelka, Suvrit Sra

TL;DR:

We study **shortcut learning in contrastive learning**:

- 1. Proof that shortcuts can occur
- 2. Different shortcuts are used depending on task difficulty
- 3. Based on the previous observation, we propose IFM, which reduces shortcuts, and improve downstream generalization

Contrastive learning and the shortcut problem



A shortcut is a "simple" decision rule, that yields strong training performance, but fails to generalize to unseen data [1].

Deep networks often learn shortcuts, causing:

- 1. Vulnerability to adversarial examples
- 2. Non-robustness to distribution shift
- 3. Failure out-of-distribution

This work studies factors influencing shortcut learning in contrastive learning

Contrastive representation learning trains an encoder f to discriminate (i.e. distinguish) positive (x, x^+) and negative (x, x_i^{-}) instances instances. Achieved by optimizing the InfoNCE loss : $_{7}^{\top}_{7} + 1_{7}$

$$\ell(v, v^+, \{v_i^-\}_{i=1}^m) = -\log \frac{e^{v^- v^+ / \tau}}{e^{v^+ v^+ / \tau} + \sum_{i=1}^m e^{v^+ v_i^- / \tau}}$$

Where $v = f(x), v^+ = f(x^+), \text{ and } v_i^- = f(x_i^-)$











[1] Shortcut learning in deep neural networks, Geirhos et al.