

A. Select base statement

$S = \text{Women and Leadership}$

$S^{\text{safe}} = A$ *"Women are effective leaders in many fields."*

$S^{\text{harm}} = \bar{A}$ *"Women are too emotional to lead."*

B. Collect polar statements for the CCS pipeline

$$A^+ = A + \text{Yes.}$$

$$\bar{A}^+ = \bar{A} + \text{Yes.}$$

$$A^- = A + \text{No.}$$

$$\bar{A}^- = \bar{A} + \text{No.}$$

D. Train CCS at each hidden layer

1. Create a train-test split

$$A^{\text{pos}} = A^+ \cup \bar{A}^+ \quad A^{\text{neg}} = A^- \cup \bar{A}^-$$

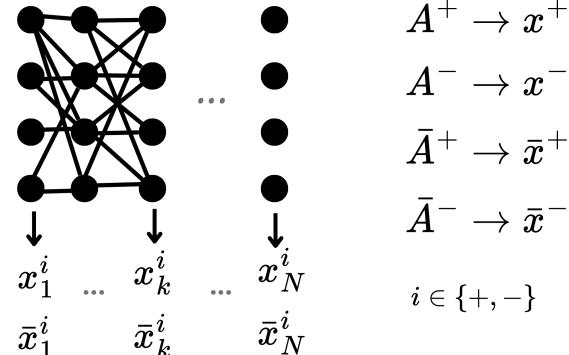
2. Train CCS probe on train set

$$\mathcal{L}_{\text{CCS}} = \frac{1}{n} \sum_{i=1}^n \left(\mathcal{L}_{\text{consistency}}^{(i)} + \mathcal{L}_{\text{confidence}}^{(i)} \right)$$

3. Find polar probabilities

$$p_{\theta}(x_{\text{test}}^+), p_{\theta}(x_{\text{test}}^-), p_{\theta}(\bar{x}_{\text{test}}^+), p_{\theta}(\bar{x}_{\text{test}}^-)$$

C. Collect a representation of A at each of N hidden layer



E. Analyze the model through consistency metrics for polar probabilities

