conditioned on value profiles

Run each instance through the decoder conditioning on

Government ...

{Yes: .7, {Yes: .8,

{Yes: .6, | {Yes: .7,

{Yes: .6. No: .4}

No: .2}

No: .3}

{Yes: .9,

No: .1}

 $p(\mathcal{Y}|\mathcal{X},V)$

Euniconweutally

 v_{N_V}

{Yes: .2,

No: .8}

{Yes: .3,

No: .7}

{Yes: .1

No: .9}

{Yes: .2,

No: .8}

each value profile, save the probabilities to a table

{Yes: .9,

No: .3}

No: .4}

{Yes: .7,

Calculate output probabilities from decoder

Calculate cross-entropy loss for assigning

each rater to each value profile

= 1.67

 v_1

2.11

1.67

0.84

3.81

 r_2

 r_{N_r}

e.g., Get "fit" ratings for rater 2

 $-\log p(\text{No}|x_{10},v_1)$

Repeat for all rater / value profile combinations

0.92

3.04

6.32

2.43

 $= -\log \frac{0.9}{0.9} - \log \frac{0.7}{0.7} - \log \frac{0.3}{0.3}$

 $D_2^{\text{fit}} = \{(x_1, \text{Yes}), (x_2, \text{Yes}), (x_{10}, \text{No})\}$

 v_{N_V}

2.45

4.22

3.23

1.34

Calculate loss of assigning rater 2 to value profile 1

 $L(r_2, v_1) = -\log p(Yes|x_1, v_1) - \log p(Yes|x_2, v_1)$

e.g., w/ 2 clusters $\hat{V}_c = \{v_1, v_2\}$

 r_{N_m}

0.92 3.04 6.32 0.84

2.43

= .92 + 1.67 + ... + 2.43

Assign new raters to cluster

with value profile that

minimizes loss

3.81

 $i \in [1, N_{\rm R}]$ $\tilde{v} \in \{v_1, v_2\}$

 $oldsymbol{Q}$ Find C (# clusters) value

 $\hat{V}_c \in \mathcal{P}(V): |\hat{V}_c| = C \sum_{i \in [1, N_{\mathrm{R}}]} \tilde{v} \in \hat{V}_c$

minimized when assigning each

rater to one of the profile clusters

profiles such that overall loss is

 $\sum \min L(r_i, \tilde{v})$

 v_{N_V}

2.45

4.22

3.23

1.34