



Investigating UD Treebanks via Dataset Difficulty Measures

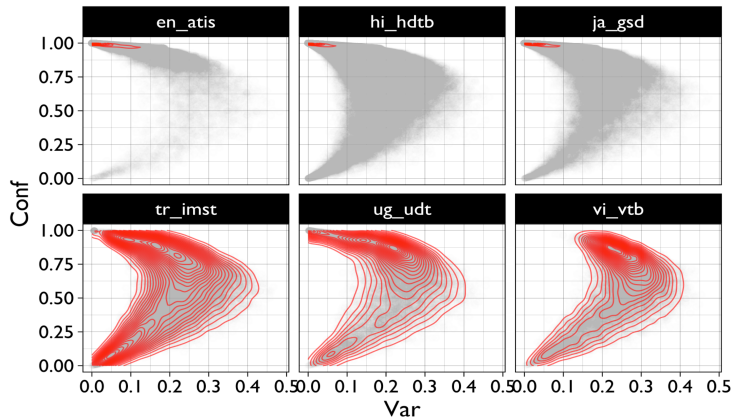
Artur Kulmizev Joakim Nivre

Dataset Cartography

Which treebanks appear hard or easy to parse, given a model's confidence throughout training, and variability therein?

$$\text{CONF}_i = \frac{1}{E} \sum_{e=1}^E p_{\theta_e}(y^*|\mathbf{x})$$

$$\text{VAR}_i = \sqrt{\frac{\sum_{e=1}^E (p_{\theta_e}(y^*|\mathbf{x}) - \text{CONF}_i)^2}{E}}$$



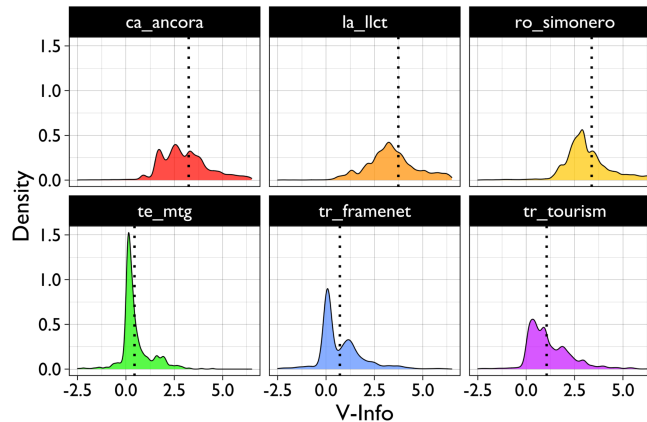
V-Information

Which treebanks contain the most (or least) information that is actually usable by a parser, with respect to a naive baseline?

$$H_V(Y|X) = \inf_{f \in \mathcal{V}} [-\log f_{\theta}[X](Y)]$$

$$H_V(Y) = \inf_{f \in \mathcal{V}} [-\log f_{\theta}[\emptyset](Y)]$$

$$I_V(X \rightarrow Y) = H_V(Y) - H_V(Y|X)$$



Minimum Description Length

Which treebanks are the most (or least) sample efficient, i.e. most easily fit by a parser, irrespective of training set size?

$$L^{\text{online}}(y_{1:n}|x_{1:n}) = \sum_{s=0}^{S-1} \sum_{n=t_s}^{t_{s+1}} -\log_2 p_{f_s}(y_n|x_n)$$

