ValueScope: Unveiling Implicit Norms and Values via Return Potential Model of Social Interactions



Chan Young Park*, Shuyue Stella Li*, Hayoung Jung*, Svitlana Volkova, Tanushree Mitra, David Jurgens, Yulia Tsvetkov {chanpark, stelli, hjung10}@cs.washington.edu}

<u>We measure implicit norms&values in community by:</u>

- 1. Introducing ValueScope, a framework to quantitatively decompose social values into social norms and community preferences.
- 2. Studying dynamics in social interactions grounded in social science theory.
- 3. Providing quantitative evidences that even similar communities exhibit diverse norms, advancing practical applications of social norm studies.

Downstream Applications

- 1. Community moderation.
- 2. Recommender systems on norm preference instead of topic.
- 3. RLC(ommunity)F: preference generator for training community specific LLMs.

Return Potential Model Plot





Rewards or Punishment?

Finding: Online communities tend to prioritize using rewards to regulate behaviors, aligning with studies showing positive signals are more effective.



Impacts of External Events

Selecting Communities & Norms to Investigate



r/askmen, r/askwomen, r/asktransgender r/stocks, r/wallstreetbets, r/pennystock r/democrats, r/republican, r/libertarian

Quantifying Preference



DioloGPT: GPT-2 trained on 147M comments We fine-tune DioloGPT for each subreddit

 $e^{h(c_i,r_i^+)}$

Findings: Norms of online communities may be influenced by external events, such as the 2020 U.S. Elections and the creation of new spinoffs.



Predicting Changes in Norms

Finding: Norms with high preference magnitude but low consensus is likely to witness upcoming change.

	\mathbf{R}^2	NI-only	NI+CR	N m
	Politeness	0.17	0.23	(d
~		A A A	0 1 0	

lorm Intensity (NI): agnitude of lis)approval.



