Dynamic Reward Adjustment in Multi-Reward Reinforcement Learning for Counselor Reflection Generation Do June Min, Verónica Pérez-Rosas, Kenneth Resnicow, Rada Mihalcea

University of Michigan

Pre-recorded video for COLING 2024

Introduction

- Role of RL in NLP
 - Challenges of Multi-reward Optimization
 - Different classes of approaches have emerged
 - Combining vs Alternate Rewards
- New Multi-armed Bandit (MAB) Algorithms
 - DynaOpt and C-DynaOpt
 - Tested on counselor reflection generation.

(1) Input prompts are fed to the generator.

(2) The generator outputs responses.

(3) Multiple reward models compute reward scores.

(4) Multiple rewards are used for updating the generator.



Background

- Challenges in Automated Reflections: Accurately mimicking human emotional dynamics and language nuances in automated reflections requires a multireward optimization approach.
- Limitations of Current RL Approaches: Static reward systems may lead to suboptimal results because they fail to adapt and capture the evolving complexities of therapeutic dialogues during training.
- Multi-armed Bandits (MABs): A class of RL approaches that manages the exploration and exploitation trade-off, enabling dynamic adjustments to the training dynamics based on observed rewards.

Multi-armed Bandit Approach



Bandit is updated

Previous Work: DORB (Pasnuru et al, 2020)

- DORB operates by cycling through multiple rewards, optimizing one at a time.
- The key idea is to use multi-armed bandits (MABs) to model the relationship between each action (which reward to optimize) and the observed reward (after the reward is optimized, what rewards have been obtained?)
- Specifically, DORB employs the Exponential-weight algorithm for Exploration and Exploitation (Exp3), which tackles the adversarial bandit problem (Auer et al., 2002), to dynamically select from a pool of reward functions during training stages.

Algorithm: DynaOpt and C-DynaOpt

Algorithm 1 DYNAOPT Optimization

- **Input:** # of rewards N, # of train steps n_{train} , # of RL validation steps round_{bandit} Initial policy p_0 , Initial Distribution of reward weights, W (uniform distribution over N)
- 1: Make a copy p_{θ} of initial policy p_0 .
- **2**: Initialize Exp3 bandit B with N + 1 arms.
- 3: Initialize weights $w_{0,i}$ over $i \in [1, 2, \dots, N+1]$ as uniform distribution.
- 5: $W \leftarrow UpdateRewardWeight(W, B_w, a, 1) \triangleright Eqns 2,3$
- **6**: $i \leftarrow 0$
- 7: while $i < n_{\text{train}} \text{ do}$
- 8: train_responses \leftarrow Sample(p_{θ} , train_data)
- 9: $r_{\text{train}} \leftarrow \text{ComputeReward}(\text{train}_{\text{responses}}, W)$
- **10:** Optimize p_{θ} with R_{train}, p_0 \triangleright Eqn 5
- 11: if i % round_{bandit} == 0 then 12: dev_responses \leftarrow Sample(p_{θ} , dev_data) 13: $r_{\text{bandit}} \leftarrow \text{ComputeReward}(\text{dev}_\text{responses})$ uniform weights) $r \leftarrow \text{ComputeBanditReward}(r_{\text{bandit}}) \triangleright \text{Eqn 4}$ 14: ⊳ Eqns 2,3 UpdateBandit(B, a, r) 15: 16: $a \leftarrow chooseArm(B)$ 17: $W \leftarrow \mathsf{UpdateRewardWeight}(W, B_w, a, 1) \triangleright$ Eqns 2,3 end if 18: 19: $i \leftarrow i + 1$ 20: end while

Proposed Methods: DynaOpt and C-DynaOpt

• Choosing an Action

$$p_t(i) = (1 - \gamma) \frac{a_{t,i}}{\sum_{j=1}^N a_{t,j}} + \frac{\gamma}{N+1}$$
 (1)

Updating the Weights (Exp3)

$$\hat{r}_{t,j}^{B^W} = \begin{cases} \frac{r_t^{B^W}}{p_t(i)} & \text{if } j = i\\ 0 & \text{otherwise} \end{cases}$$
(2)

$$w_{t+1,i} = w_{t,i} \exp\left(\frac{\gamma \hat{r}_{t,i}^{B^W}}{K}\right)$$
(3)

Bandit Reward Computation

$$\hat{r}^t = \sum_i r^t{}_i$$
 (4)
 $r^t{}_i = \operatorname{Mean}(R_{t,i}) - \operatorname{Mean}(R_{t-1,i})$

- Gamma is a hyperparmeter
- N+1 indexes N actions + No Update

Datasets

• Datasets Used: The research employs two key datasets, PAIR (Min et al, and evaluating the models.

Statistics

- # of Exchange Pairs
- Avg # of Words
- # of Complex Reflection
- # of Simple Reflection
- # of Non-Reflection

2022) and CounselChat (Welivita and Pu, 2023), which contain real-world counseling dialogues. These datasets provide diverse scenarios for training

| Welivita and Pu (2023) |
|------------------------|
| 1184 |
| 36.92 |
| 768 |
| 416 |
| 0 |
| |

Models & Algorithms

- Model: t5-base
- Baseline Algorithms
 - DORB (Pasnuru et al, 2020): Selectively chooses one reward function at a time from a set of available options based on their performance.
 - Cross Entropy: Supervised learning baseline
 - Round: Cycles through each reward function sequentially, giving each an equal opportunity to be tested and applied during training.
 - Uniform Weighted: Assigns equal weight to all reward functions



Reward Metrics

- Reflection Score (Min et al, 2022): Measures how accurately the model's generated responses reflect the emotional content and intent of the client's statements, assessing the quality of the reflections.
- Fluency: Evaluates the smoothness and naturalness of the language used in the model's responses by using perplexity.
- Coherence: Evaluates the logical flow and consistency of the generated counselor reflections. Implemented by training a RoBERTa classifier trained to detect coherent and incoherent client prompt and counselor response pairs.

Evaluation Metrics

- Automated Metrics
 - counselor reflections.
 - Levenshtein Edit Distance: Quantifies the extent to which the model successfully avoids verbatim repetition of client words.
- Human Metrics
 - Reflection, Coherence, Fluency annotated by two MI experts.

• Diversity (Li et al., 2016): Gauges the linguistic diversity of the generated

Automated Evaluation Results

| Models | Reflection (↑) | Fluency (↑) | Coherence (↑) | ∣ Edit Rate (↑) | Diversity-2 (↑) |
|------------------------------|----------------|---------------|---------------|-----------------|-----------------|
| Round | -5.02% | 11.36% | 5.51% | -8.75% | -0.20% |
| Uniform Weighted | 4.48% | 8.13% | 5.36 % | -6.28% | -0.23% |
| DORB (Pasunuru et al., 2020) | -3.03% | 9.54% | 5.42% | -7.00% | -0.08% |
| DynaOpt | 7.80% | 7.03% | 5.02% | -4.90% | -0.63% |
| C-DynaOpt | 6.14% | 8.73 % | 5.02% | -5.75% | -0.46% |

- Reflection measurements, while both achieve similar Fluency and Coherence.
- stability in the training process.

• Performance Comparison: Combine methods outperform Alternate models, particularly in

• Reflection Metrics: Notable improvement in Reflection scores with Combine methods over the Cross Entropy baseline; Alternate models show decreased performance in this area.

Stability of Training: The Round class of models exhibits higher overall variance over random runs (reflection variance of 3.59 vs 1.43 & 1.29 of our methods), indicating less

Human Evaluation Results

| | Uniform Weighted | DORB (<mark>2020</mark>) | ΟύναΟρτ | C-DynaOpt |
|------------|---------------------|-------------------------------|---------|-----------|
| Reflection | 28.29 | 25.30 | 32.10 | 29.93 |
| Fluency | 60.31 | 55.85 | 59.38 | 58.91 |
| Coherence | 62.48 | 62.79 | 63.62 | 63.68 |

- in reflection quality.

• Comparison of Models: Human evaluations confirm that Combine models, specifically DynaOpt and C-DynaOpt, outperform the Uniform Weighted model

 Fluency and Coherence: Despite lower scores in automated metrics, DynaOpt and C-DynaOpt exhibit superior fluency and coherence in human evaluations, suggesting that higher reflection levels enhance the naturalness of responses.

Overall Results

- Not All Multi-reward Optimization Methods Are Effective for Counselor **Reflection Generation.**
 - the Alternate methods (Table 2).
- Comparative Advantage of Our Methods
 - reflection levels.

• Methods that combine weights exhibit superior performance compared to

 DynaOpt and C-DynaOpt outperform not only the Alternate methods but also the Uniform Weighted baseline in terms of both automated and human

Bandit Visualization





Conclusion

- Problem Addressed: Optimizing multiple linguistic rewards in reinforcement learning for counselor reflection in motivational interviewing.
- Strategies Explored: Investigated Alternate and Combine approaches, enhancing them with bandit-augmented versions.
- Novel Methods: Introduced DynaOpt and C-DynaOpt, which dynamically adjust reward weights using multi-armed bandits during training.
- Empirical Findings: Demonstrated that previous approaches failed to improve response quality, while DynaOpt and C-DynaOpt surpassed existing baselines.

Acknowledgments, Resources, and Contact

- Acknowledgments: Special thanks to the research team and students from the University of Michigan School of Public Health for their invaluable feedback and participation. This work was partially supported by the National Science Foundation under award #2306372. The views expressed are those of the authors and do not necessarily reflect the views of the NSF.
- Our code is available at <u>https://github.com/michigannlp/dynaopt</u>
- Contact for Questions and Discussion: For further discussions and any questions, please email dojmin@umich.edu.