# PERSONALIZED REWARD LEARNING WITH INTERACTION-GROUNDED LEARNING (IGL)

#### Anonymous authors

Paper under double-blind review

### Abstract

In an era of countless content offerings, recommender systems alleviate information overload by providing users with personalized content suggestions. Due to the scarcity of explicit user feedback, modern recommender systems typically optimize for a fixed combination of implicit feedback signals across all users. However, this approach disregards a growing body of work that (i) implicit signals can be used by users in diverse ways, signaling anything from satisfaction to active dislike, and (ii) different users communicate preferences in different ways. We propose applying the recent Interaction Grounded Learning (IGL) paradigm to address the challenge of learning representations of diverse user communication modalities. Rather than taking a fixed, human-designed reward function, IGL is able to learn personalized reward functions for different users and then optimize directly for the latent user satisfaction. We demonstrate the success of IGL with experiments using simulations as well as with real-world production traces.

# **1** INTRODUCTION

From shopping to reading the news, modern Internet users have access to an overwhelming amount of content and choices from online services. Recommender systems offer a way to improve user experience and decrease information overload by providing a customized selection of content. A key challenge for recommender systems is the rarity of explicit user feedback, such as ratings or likes/dislikes (Grčar et al., 2005). Rather than explicit feedback, practitioners typically use more readily available implicit signals, such as clicks (Hu et al., 2008), webpage dwell time (Yi et al., 2014), or inter-arrival times (Wu et al., 2017) as a proxy signal for user satisfaction. These implicit signals are used as the reward objective in recommender systems, with the popular Click-Through Rate (CTR) metric as the gold standard for the field (Silveira et al., 2019). However, directly using implicit signals as the reward function presents several issues.

*Implicit signals do not directly map to user satisfaction.* Although clicks are routinely equated with user satisfaction, there are examples of unsatisfied users interacting with content via clicks. Clickbait exploits cognitive biases such as caption bias (Hofmann et al., 2012) or the curiosity gap (Scott, 2021) so that low quality content attracts more clicks. Direct optimization of the CTR degrades user experience by promoting clickbait items (Wang et al., 2021). Recent work shows that users will even click on content that they know a priori they will dislike. In a study of online news reading, Lu et al. (2018a) discovered that 15% of the time, users would click on articles that they strongly disliked. Similarly, although longer webpage dwell times are associated with satisfied users, a study by Kim et al. (2014) found that dwell time is also significantly impacted by page topic, readability and content length.

*Different users communicate in different ways.* Demographic background is known to have an impact on the ways in which users engage with recommender systems. A study by Beel et al. (2013) shows that older users have CTR more than 3x higher than their younger counterparts. Gender also has an impact on interactions, e.g. men are more likely to leave dislikes on YouTube videos than women (Khan, 2017). At the same time, a growing body of work shows that recommender systems do not provide consistent performance across demographic subgroups. For example, multiple studies on ML fairness in recommender systems show that women on average receive less accurate recommendations compared to men (Ekstrand et al., 2018; Mansoury et al., 2020). Current systems are also unfair across different age brackets, with statistically significant recommendation utility

degradation as the age of the user increases (Neophytou et al., 2022). The work of Neophytou et al. identifies usage features as the most predictive of mean recommender utility, hinting that the inconsistent performance in recommendation algorithms across subgroups arises from the differences in how users interact with the recommender system.

These challenges motivate the need for personalized reward functions. However, extensively modeling the ways in which implicit signals are used or how demographics impact interaction style is costly and inefficient. Furthermore, as recommender systems and their users evolve, so do the ways in which users implicitly communicate preferences. Any extensive models developed now could easily become obsolete within a few years time.

To this end, we propose Interaction Grounded Learning (IGL) Xie et al. (2021) for personalized reward learning (IGL-P). IGL is a learning paradigm where a learner optimizes for unobservable rewards by interacting with the environment and associating observable feedback with the true latent reward. Prior IGL approaches assume the feedback either depends on the reward alone Xie et al. (2021), or on the reward and action Xie et al. (2022). These methods are unable to disambiguate personalized feedback which depends on the context. Other approaches such as reinforcement learning and traditional contextual bandits suffer from the choice of reward function. However our proposed personalized IGL, IGL-P, resolves the 2 above challenges while making minimal assumptions about the value of observed user feedback. Our new approach is able to incorporate both explicit and implicit signals, leverage ambiguous user feedback and adapt to the different ways in which users interact with the system.

**Our Contributions:** We present the first IGL strategy for context-dependent feedback, the first use of inverse kinematics as an IGL objective, and the first IGL strategy for more than two latent states. Using simulations and real production data, we demonstrate that recommender systems require at least 3 reward states, and that IGL is able to address two of the biggest challenges for modern online recommender systems.

## 2 PROBLEM SETTING

#### 2.1 CONTEXTUAL BANDITS

The contextual bandit (Auer et al., 2002; Langford & Zhang, 2007) is a statistical model of myopic decision making which is pervasively applied in recommendation systems (Bouneffouf et al., 2020). IGL operates via reduction to contextual bandits, hence, we briefly review contextual bandits here.

The contextual bandit problem proceeds over T rounds. At each round  $t \in [T]$ , the learner receives a context  $x_t \in \mathcal{X}$  (the context space), selects an action  $a_t \in \mathcal{A}$  (the action space), and then observes a reward  $r_t(a_t)$ , where  $r_t : \mathcal{A} \to [0, 1]$  is the underlying reward function. We assume that for each round t, conditioned on  $x_t$ ,  $r_t$  is sampled from a distribution  $\mathbb{P}_{r_t}(\cdot \mid x_t)$ . A contextual bandit algorithm attempts to minimize the regret

$$\mathbf{Reg}_{\mathsf{CB}}(T) := \sum_{t=1}^{T} r_t(\pi^*(x_t)) - r_t(a_t)$$
(1)

relative to an optimal policy  $\pi^*$  over a policy class  $\Pi$ .

In general, both the contexts  $x_1, \ldots, x_T$  and the distributions  $\mathbb{P}_{r_1}, \ldots, \mathbb{P}_{r_T}$  can be selected in an arbitrary, potentially adaptive fashion based on the history. In the sequel we will describe IGL in a stochastic environment, but the reduction induces a nonstationary contextual bandit problem, and therefore the existence of adversarial contextual bandit algorithms is relevant.

#### 2.2 INTERACTION GROUNDED LEARNING

IGL extends the contextual bandit framework by eliding the reward from the learning algorithm and providing feedback instead (Xie et al., 2021). We describe the stochastic setting where  $(x_t, r_t, y_t) \sim D$  triples are sampled iid from an unknown distribution; the learner receives the context  $x_t \in \mathcal{X}$ , selects an action  $a_t \in \mathcal{A}$ , and then observes the feedback  $y_t(a_t)$ , where  $y_t : \mathcal{A} \to [0, 1]$  is the underlying feedback function. Note  $r_t(a_t)$  is never revealed to the algorithm: nonetheless, the



Figure 1: IGL in the recommender system setting. The learner observes the context x, plays an action a, and then observes a feedback y (that is dependent on the latent reward r), but not r itself.

regret notion remains the same as Eq. (1). An information-theoretic argument proves assumptions relating the feedback to the underlying reward are necessary to succeed (Xie et al., 2022).

## 2.2.1 Specialization to Recommendation

For specific application in the recommendation domain, we depart from prior art in IGL (Xie et al., 2021; 2022) in two ways: first, in the assumed relationship between feedback and underlying reward; and second, in the number of latent reward states.

**Feedback Dependence Assumption** Xie et al. (2021) assumed full contextual independence of the feedback on the context and chosen action, i.e.  $y \perp x, a | r$ . For recommender systems, this implies that all users communicate preferences identically for all content. In a subsequent paper, Xie et al. (2022) loosen full conditional independence by considering context conditional independence, i.e.  $y \perp x | a, r$ . For our setting, this corresponds to the user feedback varying for combinations of preference and content, but remaining consistent across all users. Neither of these two assumptions are natural in the recommendation setting because different users interact with recommender systems in different ways. (Beel et al., 2013; Shin, 2020). In this work, we assume  $y \perp a | x, r$ , i.e., the feedback y is independent of the displayed content a given the user x and their disposition toward the displayed content r. Thus, we assume that users may communicate in different ways, but a given user expresses satisfaction, dissatisfaction and indifference to all content in the same way.

**Number of Latent Reward States** Prior work demonstrates a binary latent reward assumption, along with an assumption that rewards are rare under a known reference policy, is sufficient for IGL to succeed. Specifically, optimizing the contrast between a learned policy and the oblivious uniform policy is able to succeed when feedback is both context and action independent Xie et al. (2021); and optimizing the contrast between the learned policy and all constant-action policies succeeds when the feedback is context independent Xie et al. (2022).

Although the binary latent reward assumption (e.g., satisfied or dissatisfied) appears reasonable for recommendation scenarios, it fails to account for user indifference versus user dissatisfaction. This observation was first motivated by our production data, where a 2 state IGL policy would sometimes maximize feedback signals with obviously negative semantics. Assuming users ignore most content most of the time (Nguyen et al., 2014), negative feedback can be as difficult to elicit as positive feedback, and a 2 state IGL model is unable to distinguish between these extremes. Hence, we posit a minimal latent state model for recommender systems involves 3 states: (i) r = 1, when users are satisfied with the recommended content, (ii) r = 0, when users are indifferent or inattentive, and (iii) r = -1, when users are dissatisfied.

# 3 DERIVATIONS

Prior approaches to IGL use contrastive learning objectives (Xie et al., 2021; 2022), but the novel feedback dependence assumption in the prior section impedes this line of attack. Essentially, given arbitrary dependence upon x, learning must operate on each example in isolation without requiring comparison across examples. This motivates attempting to predict the current action from the current context and the currently observed feedback, i.e., inverse kinematics.

**Inverse Kinematics** We motivate our inverse kinematics strategy using exact expectations. When acting according to any policy P(a|x), we can imagine trying to predict the action taken given the context and feedback; the posterior distribution is

$$P(a|y,x) = \frac{P(a|x)P(y|a,x)}{P(y|x)}$$

$$= P(a|x)\sum_{x} \frac{P(y|r,a,x)}{P(y|x)}P(r|a,x)$$
(Bayes rule)
(Total Probability)

$$=P(a|x)\sum_{r}\frac{P(y|r,x)}{P(y|x)}P(r|a,x) \hspace{1.5cm} (y\perp a|x,r)$$

$$= P(a|x) \sum_{r} \frac{P(r|y,x)}{P(r|x)} P(r|a,x)$$
 (Bayes rule)

$$=\sum_{r} P(r|y,x) \left( \frac{P(r|a,x)P(a,x)}{\sum_{a} P(r|a,x)P(a|x)} \right). \quad \text{(Total Probability)} \tag{2}$$

We arrive at the inner product between a reward decoder term (P(r|y, x)) and a reward predictor term (P(r|a, x)).

**Extreme Event Detection** Direct extraction of a reward predictor using maximum likelihood on the action prediction problem with Eq. (2) is frustrated by two identifiability issues: first, this expression is invariant to a permutation of the rewards on a context dependent basis; and second, the relative scale of two terms being multiplied is not uniquely determined by their product. To mitigate the first issue, we assume  $\sum_a P(r = 0|a, x)P(a|x) > \frac{1}{2}$ , i.e., nonzero rewards are rare under P(a|x); and to mitigate the second issue, we assume the feedback can be perfectly decoded, i.e.,  $P(r|y, x) \in \{0, 1\}$ . Under these assumptions we have

$$r = 0 \implies P(a|y,x) = \frac{P(r=0|a,x)P(a|x)}{\sum_{a} P(r=0|a,x)P(a|x)}$$
$$\leq 2P(r=0|a,x)P(a|x) \leq 2P(a|x). \tag{3}$$

Eq. (3) forms the basis for our extreme event detector: anytime the posterior probability of an action is predicted to be more than twice the prior probability, we deduce  $r \neq 0$ .

Note a feedback merely being apriori rare or frequent (i.e., the magnitude of P(y|x) under the policy P(a|x)) does not imply that observing such feedback will induce an extreme event detection; rather the feedback must have a probability that strongly depends upon which action is taken. Because feedback is assumed conditionally independent of action given the reward, the only way for feedback to help predict which action is played is via the (action dependence of the) latent reward.

**Extreme Event Disambiguation** With 2 latent states,  $r \neq 0 \implies r = 1$ , and we can reduce to a standard contextual bandit with inferred rewards  $\mathbb{1}(P(a|y,x) > 2P(a|x))$ . With 3 latent states,  $r \neq 0 \implies r = \pm 1$ , and additional information is necessary to disambiguate the extreme events. We assume partial reward information is available via a "definitely negative" function<sup>1</sup> DN :  $\mathcal{X} \times \mathcal{Y} \rightarrow \{-1,0\}$  where P(DN(x,y) = 0|r = 1) = 1 and P(DN(x,y) = -1|r = -1) > 0. This reduces extreme event disambiguation to one-sided learning (Bekker & Davis, 2020) applied only to extreme events, where we try to predict the underlying latent state given (x, a). We assume partial labelling is selected completely at random Elkan & Noto (2008) and treat the (constant) negative labelling propensity  $\alpha$  as a hyperparameter. We arrive at our 3-state reward extractor

$$\rho(x, a, y) = \begin{cases}
0 & P(a|y, x) \le 2P(a|x) \\
-\alpha^{-1} & P(a|y, x) > 2P(a|x) \text{ and } DN(x, y) = -1 , \\
1 & \text{otherwise}
\end{cases}$$
(4)

equivalent to Bekker & Davis (2020, Equation 11). Setting  $\alpha = 1$  embeds 2-state IGL.

<sup>&</sup>lt;sup>1</sup>"Definitely positive" information can be incorporated analogously.

Algorithm 1 IGL; Inverse Kinematics; 2 or 3 Latent States; On or Off-Policy. Input: Contextual bandit algorithm CB-Alg. Input: Calibrated weighted multiclass classification algorithm MC-Alg.  $\# DN(\ldots) = 0$  for 2 state IGL **Input:** Definitely negative oracle DN. **Input:** Negative labelling propensity  $\alpha$ . #  $\alpha = 1$  for 2 state IGL **Input:** Action set size K. 1:  $\pi \leftarrow \text{new CB-Alg.}$ 2: IK  $\leftarrow$  new MC-Alg. 3: for t = 1, 2, ...; do Observe context  $x_t$  and action set  $A_t$  with  $|A_t| = K$ . 4: 5: if On-Policy IGL then 6:  $P(\cdot|x_t) \leftarrow \pi.\operatorname{predict}(x_t, A_t).$ 7: Play  $a_t \sim P(\cdot|x_t)$  and observe feedback  $y_t$ . 8: else 9: Observe  $(x_t, a_t, y_t, P(\cdot|x_t))$ .  $w_t \leftarrow 1/(KP(a_t|x_t)).$ 10: # Synthetic uniform distribution  $\hat{P}(a_t|y_t, x_t) \leftarrow \texttt{IK.predict}((x_t, y_t), A_t, a_t).$ # Predict action probability 11: if  $K\hat{P}(a_t|y_t, x_t) \leq 2$  then 12:  $\# \hat{r}_t = 0$ 13:  $\pi$ .learn $(x_t, a_t, A_t, r_t = 0)$ #  $\hat{r}_t \neq 0$ 14: else if  $DN(\ldots) = 0$  then 15:  $\pi$ .learn $(x_t, a_t, A_t, r_t = 1, P(\cdot | x_t))$ 16: 17: # Definitely negative else  $\pi$ .learn $(x_t, a_t, A_t, r_t = -\alpha^{-1}, P(\cdot | x_t))$ 18: 19: IK.learn $((x_t, y_t), A_t, a_t, w_t)$ .

**Implementation Notes** In practice, P(a|x) is known but the other probabilities are estimated.  $\hat{P}(a|y,x)$  is estimated online using maximum likelihood on the problem predicting *a* from (x,y), i.e., on a data stream of tuples ((x,y),a). The current estimates induce  $\hat{\rho}(x,a,y)$  based upon the plug-in version of Eq. (4). In this manner, the original data stream of (x, a, y) tuples is transformed into stream of  $(x, a, \hat{r} = \hat{\rho}(x, a, y))$  tuples and reduced to a standard online contextual bandit problem.

As an additional complication, although P(a|x) is known, it is typically a good policy under which rewards are not rare (e.g., offline learning with a good historical policy; or acting online according to the policy being learned by the IGL procedure). Therefore we use importance weighting to synthesize a uniform action distribution P(a|x) from the true action distribution.<sup>2</sup> Ultimately we arrive at the procedure of Algorithm 1.

## 4 EMPIRICAL EVALUATIONS

**Evaluation Settings:** Evaluation settings include simulation using a supervised classification dataset, online news recommendation on Facebook, and a production image recommendation scenario.

Abbreviations: Algorithms are denoted by the following abbreviations: Personalized IGL for 2 latent states (IGL-P(2)); Personalized IGL for 3 latent states (IGL-P(3)); Contextual Bandits for the Facebook news setting that maximizes for emoji, non-like click-based reactions (CB-emoji); Contextual Bandits for the Facebook news setting that maximizes for comment interactions (CB-comment).

**General Evaluation Setup:** At each time step t, the context  $x_t$  is provided from either the simulator (Section 4.1, Section 4.2) or the logged production data (Section 4.3). The learner then selects an action  $a_t$  and receives feedback  $y_t$ . In these evaluations, each user provides feedback in exactly one interaction and different user feedback signals are mutually exclusive, so that  $y_t$  is a one-hot

<sup>&</sup>lt;sup>2</sup>When the number of actions is changing from round to round, we use importance weighting to synthesize a non-uniform action distribution with low rewards, but we elide this detail for ease of exposition.



Figure 2: The proposed personalized IGL algorithm successfully disambiguates both different user communication styles and different event semantics. See Section 4.1 for details.

vector. In simulated environments, the ground truth reward is sometimes used for evaluation but never revealed to the algorithm.

**Code:** Our code will be made available at {url redacted} for all publicly replicable experiments (i.e., except for the production data).

#### 4.1 COVERTYPE IGL SIMULATION

To highlight that personalized IGL can distinguish between different user communication styles, we create a simulated 2-state IGL scenario from a supervised classification dataset. First, we apply a supervised-to-bandit transform to convert the dataset into a contextual bandit simulation (Bietti et al., 2021), i.e., the algorithm is presented the example features as context, chooses one of the classes as an action, and experiences a binary reward which indicates whether or not it matches the example label. In the IGL simulation this reward is experienced but not revealed to the algorithm. Instead, the latent reward is converted into a feedback signal as follows: each example is assigned one of N different user ids and the user id is revealed to the algorithm as part of the example features. The simulated user will generate feedback in the form of one of M different word ids. Unknown to the algorithm, the words are divided equally into "good" and "bad" words, and the users are divided equally into "normal" and "bizarro" users. "Normal" users indicate positive and zero reward via "good" and "bad" words respectively, while "bizarro" users employ the exact opposite communication convention.

We simulated using the Covertype (Blackard & Dean, 1999) dataset with M = N = 100, and an (inverse kinematics) model class which embedded both user and word ids into a 2 dimensional space. Fig. 2 demonstrates both the user population and the words are cleanly separated into two latent groups. Additional results showcasing the learning curves for inverse kinematics, reward and policy learning are shown in **??**.

#### 4.2 FAIRNESS IN FACEBOOK NEWS RECOMMENDATION

Personalized reward learning is the key to more fair recommender systems. Previous work (Neophytou et al., 2022) suggests that inconsistent performance in recommender systems across user subgroups arises due to differences in user communication modalities. We now test this hypothesis in the setting of Facebook news recommendation. Our simulations are built on a dataset (Martinchek, 2016) of all posts by the official Facebook pages of 3 popular news outlets (Fox News, The Huffington Post and TIME Magazine) that span the political spectrum. Posts range from May to November 2016 and contain text content information, as well as logged interaction counts, which include comments and shares, as well as diverse click-based reactions (see Fig. 3).

Constructing a hand-engineered reward signal using these feedbacks is difficult, and Facebook itself came under fire for utilizing reward weights that disproportionately promote toxic, low quality



Figure 3: Facebook click-based reactions: like, love, haha, wow, sad and angry (image source: Meta). The reactions allow users to engage with content using diverse communication signals.

news. One highly criticized iteration of the reward ranking algorithm treated emoji reactions as five times more valuable than likes (Merrill & Oremus, 2021). Future iterations of the ranking algorithm promoted comments, in an attempt to bolster "meaningful social interactions" (Hagey & Jeff Horwitz, 2021). Our experiments evaluate the performance of CB algorithms using these two reward functions, referring to them as CB-emoji and CB-comment.

We model the news recommendation problem as a 3 latent state problem, with readers of different news outlets as different contexts. Given a curated selection of posts, the goal of the learner is to select the best article to show to the reader. The learner can leverage action features including the post type (link, video, photo, status or event) as well as embeddings of the text content that were generated using pre-trained transformers (Reimers & Gurevych, 2019). User feedback is drawn from a fixed probability distribution (unknown to the algorithms) that depends on the user type and latent reward of the chosen action. As an approximation of the true latent reward signal, we use low dimensional embeddings of the different news outlets combined with aggregate statistics from the post feedback to categorize whether the users had a positive (r = 1), neutral (r = 0) or negative experience (r = -1) with the post. This categorization is not available to the evaluated algorithms. Finally, we implement IGL-P (3) with the angry reaction as a negative oracle to disambiguate the positive and negative reward states.





(a) Average fraction of rewards that are positive

(b) Average fraction of rewards that are negative

Figure 4: IGL uses personalized reward learning to achieve fair news recommendations across diverse reader bases, while CB policies based off of rewards used in practice by Facebook perform inconsistently, with subsets of users receiving both fewer high quality recommendations and more low quality recommendations. Standard error on all averages shown is < 0.005.

Fig. 4 shows the results of our online news recommendation experiments. While the performance of both CB-emoji and CB-comment varies significantly across the different reader groups, IGL-P(3) maintains relatively stable performance for both positive and negative rewards. The CB algorithm that maximizes emoji reactions achieves the best performance for TIME readers at the cost of very bad performance for the Fox News and Huffington Post Readers. On the other hand, the CB algorithm that maximizes comment engagement achieves best performance with Fox News readers, however it still performs worse than IGL-P(3). Finally, our simulations show that the CB-comment objective that was introduced to decrease low quality news actually significantly *increased* it for the Fox News reader population.

#### 4.3 **PRODUCTION RESULTS**

Our production setting is a real world image recommendation system that serves hundreds of millions of users. In our recommendation system interface, users provide feedback in the form of clicks,

Algorithm	Clicks	Likes	Dislikes
IGL-P(3)	[0.999, 1.067, 1.152]	[0.985, 1.029, 1.054]	[0.751, 1.072, 1.274]
IGL-P(2)	[0.926, 1.005, 1.091]	[0.914, 0.949, 0.988]	[1.141, 1.337, 1.557]

Table 1: Relative metrics lift over a production baseline. The production baseline uses a handengineered reward function which is not available to IGL algorithms. Shown are point estimates and associated bootstrap 95% confidence regions. IGL-P(2) erroneously increases dislikes to the detriment of other metrics. IGL-P(3) is equivalent to the hand-engineered baseline.

likes, dislikes or no feedback. All four signals are mutually exclusive and the user only provides one feedback after each interaction. For these experiments, we use data that spans millions of interactions. The production baseline is a contextual bandit algorithm with a hand-engineered multi-task reward function which dominates approaches that only use click feedback<sup>3</sup>. Consequently, any improvements over the production policy imply improvement over any bandit algorithm optimizing for click feedback.

We implement IGL-P(2) and IGL-P(3) and report the performance as relative lift metrics over the production baseline. Unlike the simulation setting, we no longer have access to the user's latent reward after each interaction. As a result, we evaluate IGL by comparing all feedback signals. An increase in both clicks and likes, and a decrease in dislikes, are considered desirable outcomes. Table 1 shows the results of our empirical study.

IGL-P(2) exhibits an inability to avoid extreme *negative* events. Although the true latent state is unknown, IGL-P(2) is Pareto-dominated due to an increase in dislikes. IGL-P(3) does not exhibit this pathology. These results indicate the utility of a 3 latent state model in real world recommendation systems.

## 5 RELATED WORK

Recommender systems are a well-studied field due to their direct link to product revenues (Naumov et al., 2019; Steck et al., 2021). The rapid growth of online content has generated interest in ML-based solutions (Li et al., 2011) that are able to offer more diverse personalized recommendations to the internet users. Traditional vanilla recommendation approaches can be divided into three types. Content-based approaches (Balabanović & Shoham, 1997; IJntema et al., 2010; Kompan & Bieliková, 2010; Lops et al., 2019; Argyriou et al., 2020; Javed et al., 2021) maintain representations for users based on their content and recommend new content with good similarity metrics for particular users. In contrast, collaborative filtering approaches (Balabanović & Shoham, 1997; Schafer et al., 2007; Hu et al., 2008; Argyriou et al., 2020; Steck & Liang, 2021) employ user rating predictions based on historical consumed content and underlying user similarities. Finally, there are hybrid approaches (Balabanović & Shoham, 1997; Funakoshi & Ohguro, 2000; Burke, 2007; Argyriou et al., 2020; Javed et al., 2021) that combine the previous two contrasting approaches to better represent user profiles for improved recommendations. Our work is a significant departure from these approaches, in that we learn representations for users via their content interaction history for improved diverse personalized recommendations.

Recommendation as a contextual bandits problem has a rich history (Li et al., 2010; 2011; Bouneffouf et al., 2020). Typically, implicit signals such as the CTR metric are meticulously incorporated into manually-engineered reward functions (Li et al., 2010; 2011; Bouneffouf et al., 2020). Variants include using ensembles of contextual bandits Tang et al. (2014), applying collaborative filtering approaches Gentile et al. (2014); Wu et al. (2016), and using content-based hybrid methods (Li et al., 2010; Ding et al., 2021). Other works formulate recommendation systems as a reinforcement learning problem (Zou et al., 2019; Lin et al., 2021; Afsar et al., 2021). All of these crucially depend on the implicit signals for driving the recommendation systems, but suffer from the disconnect that these implicit signals have with respect to the true user satisfaction (Section 1), which hinders the practicality of these approaches. Alternatively, inverse reinforcement learning based recommendation systems (Chen et al., 2021b; Hu et al., 2022) learn the complex reward structure through

<sup>&</sup>lt;sup>3</sup>The utility of multi-task learning for recommendation systems is well-established, e.g., Chen et al. (2021a); Lu et al. (2018b); Chen et al. (2019).

demonstrations of expert recommendations. However inverse reinforcement learning is incapable of generalizing Chen et al. (2022) to different feedback signals, often posed as an optimization problem to imitate the expert demonstrations that may itself be bad recommendation agents, and generally need expensive compute and enormous sampling capabilities for learning.

Our work is significantly different from these in the following way: (*i*) we formulate and propose a new recommendation system based on the IGL paradigm, (*ii*) we leverage the capability of IGL learning rewards based off of implicit and explicit feedback signals and avoid the costly, inefficient, status quo process of reward engineering, and (*iii*) we propose a novel personalized IGL algorithm based on the inverse kinematics strategy as described in Section 3. The works that are closest to ours are Xie et al. (2021; 2022) which introduce and solve for the IGL paradigm under different assumptions. However, we propose a personalized IGL algorithm for recommendation systems with improved reward predictor models, more practical assumptions on feedback signals, and more intricacies described in Section 2.2.1.

# 6 **DISCUSSION**

We evaluated the proposed personalized IGL approach (IGL-P) in three different settings: (1) A simulation using a supervised classification dataset shows that IGL-P can learn to successfully distinguish between different communication modalities; (2) A simulation for online news recommendation based on real data from Facebook users shows that IGL-P leverages insights about different communication modalities to learn better policies and achieve fairness with consistent performance among diverse user groups; (3) A real-world experiment deployed in an image recommendation product showcases that the proposed method outperforms the hand-engineered reward baseline, and succeeds in a practical application.

This work assumes that user may communicate in different ways, but a given user expresses (dis)satisfaction or indifference to all content in the same way. This assumption was critical to deriving the inverse kinematics approach, but in practice user feedback can also depend upon content (Freeman et al., 2020). IGL with arbitrary joint content-action dependence of feedback is intractable, but plausibly there exists a tractable IGL setting with a constrained joint content-action dependence which is a better fit for recommendation scenarios. Furthermore, although we established that utility of a three state model over a two state model in our experiments, perhaps more than three states is necessary for more complex recommendation scenarios.

We also explore the possibilities of achieving fairness through personalized reward learning. Our findings build off observations in the literature and are empirical results showing the potential for consistent performance using IGL-P. By utilizing personalized rewards, unlike alternative approaches, IGL is not incentivized to sacrifice performance for sub-populations in order to achieve better average performance across all users. Theoretical guarantees of IGL-P fairness remain an interesting future direction. Potential obstacles to guaranteed fairness might arise due to insufficient data from underrepresented populations, as well as scarcity of a negative oracle signal in the 3 latent state setting. Although IGL-P provably learns even with very rare presence of a negative oracle signal, inconsistent performance across user subsets can arise due to faster learning and convergence for populations that utilize the negative oracle signal more frequently.

# ETHICS STATEMENT

In our paper, we use two real-world interaction datasets. The first is a publicly available dataset of interactions with public Facebook news pages. All interactions are anonymous and the identities of users are excluded from the dataset, preserving their privacy. Similarly, our production data does not include user information and was used with the consent and permission of all relevent parties.

# **Reproducibility Statement**

We have taken considerable measures to ensure the results are as reproducible as possible. We have provided the code to replicate our experiment results (except for the production results in Section 4.3) as part of the supplementary material. The code will be made publicly available at

{url redacted}. We used publicly available datasets for our simulated experiments in Section 4.1 (Blackard & Dean, 1999) and Section 4.2 Martinchek (2016). The experiment code for Section 4.1, when executed, will automatically download the dataset. The dataset for Section 4.2 is included as part of our supplementary material. Finally, our supplementary material also includes a conda environment file to help future researchers recreate our development environment on their machines when running the experiments.

## REFERENCES

- M Mehdi Afsar, Trafford Crump, and Behrouz Far. Reinforcement learning based recommender systems: A survey. ACM Computing Surveys (CSUR), 2021.
- Andreas Argyriou, Miguel González-Fierro, and Le Zhang. Microsoft recommenders: Best practices for production-ready recommendation systems. In *Companion Proceedings of the Web Conference* 2020, pp. 50–51. Association for Computing Machinery, 2020. ISBN 9781450370240.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- Marko Balabanović and Yoav Shoham. Fab: content-based, collaborative recommendation. Communications of the ACM, 40(3):66–72, 1997.
- Joeran Beel, Stefan Langer, Andreas Nürnberger, and Marcel Genzmehr. The impact of demographics (age and gender) and other user-characteristics on evaluating recommender systems. In *International Conference on Theory and Practice of Digital Libraries*, pp. 396–400. Springer, 2013.
- Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: A survey. *Machine Learning*, 109(4):719–760, 2020.
- Alberto Bietti, Alekh Agarwal, and John Langford. A contextual bandit bake-off. J. Mach. Learn. Res., 22:133–1, 2021.
- Jock A Blackard and Denis J Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and electronics in agriculture*, 24(3):131–151, 1999.
- Djallel Bouneffouf, Irina Rish, and Charu Aggarwal. Survey on applications of multi-armed and contextual bandits. In 2020 IEEE Congress on Evolutionary Computation (CEC), pp. 1–8. IEEE, 2020.
- Robin Burke. Hybrid web recommender systems. The adaptive web, pp. 377-408, 2007.
- Lei Chen, Jie Cao, Guixiang Zhu, Youquan Wang, and Weichao Liang. A multi-task learning approach for improving travel recommendation with keywords generation. *Knowledge-Based Systems*, 233:107521, 2021a.
- Xiaocong Chen, Lina Yao, Aixin Sun, Xianzhi Wang, Xiwei Xu, and Liming Zhu. Generative inverse deep reinforcement learning for online recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 201–210, 2021b.
- Xiaocong Chen, Lina Yao, Xianzhi Wang, Aixin Sun, and Quan Z Sheng. Generative adversarial reward learning for generalized behavior tendency inference. *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- Zhongxia Chen, Xiting Wang, Xing Xie, Tong Wu, Guoqing Bu, Yining Wang, and Enhong Chen. Co-attentive multi-task learning for explainable recommendation. In *IJCAI*, pp. 2137–2143, 2019.
- Qinxu Ding, Yong Liu, Chunyan Miao, Fei Cheng, and Haihong Tang. A hybrid bandit framework for diversified recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 4036–4044, 2021.

- Michael D Ekstrand, Mucun Tian, Ion Madrazo Azpiazu, Jennifer D Ekstrand, Oghenemaro Anuyah, David McNeill, and Maria Soledad Pera. All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. In *Conference on fairness, accountability and transparency*, pp. 172–186. PMLR, 2018.
- Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 213–220, 2008.
- Cole Freeman, Hamed Alhoori, and Murtuza Shahzad. Measuring the diversity of facebook reactions to research. *Proceedings of the ACM on Human-Computer Interaction*, 4(GROUP):1–17, 2020.
- Kaname Funakoshi and Takeshi Ohguro. A content-based collaborative recommender system with detailed use of evaluations. In KES'2000. Fourth International Conference on Knowledge-Based Intelligent Engineering Systems and Allied Technologies. Proceedings (Cat. No. 00TH8516), volume 1, pp. 253–256. IEEE, 2000.
- Claudio Gentile, Shuai Li, and Giovanni Zappella. Online clustering of bandits. In *International Conference on Machine Learning*, pp. 757–765. PMLR, 2014.
- Miha Grčar, Dunja Mladenič, Blaž Fortuna, and Marko Grobelnik. Data sparsity issues in the collaborative filtering framework. In *International workshop on knowledge discovery on the web*, pp. 58–76. Springer, 2005.
- Keach Hagey and Jeff Jeff Horwitz. Facebook tried to make its platform a healthier place. it got angrier instead. *The Wall Street Journal*, 2021. URL https://www.wsj.com/articles/facebook-algorithm-change-zuckerberg-11631654215.
- Katja Hofmann, Fritz Behr, and Filip Radlinski. On caption bias in interleaving experiments. In Proceedings of the 21st ACM international conference on Information and knowledge management, pp. 115–124, 2012.
- Chenhao Hu, Shuhua Huang, Yansen Zhang, and Yubao Liu. Learning to infer user implicit preference in conversational recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 256–266, 2022.
- Yifan Hu, Yehuda Koren, and Chris Volinsky. Collaborative filtering for implicit feedback datasets. In 2008 Eighth IEEE international conference on data mining, pp. 263–272. Ieee, 2008.
- Wouter IJntema, Frank Goossen, Flavius Frasincar, and Frederik Hogenboom. Ontology-based news recommendation. In Proceedings of the 2010 EDBT/ICDT Workshops, pp. 1–6, 2010.
- Umair Javed, Kamran Shaukat, Ibrahim A Hameed, Farhat Iqbal, Talha Mahboob Alam, and Suhuai Luo. A review of content-based and context-based recommendation systems. *International Jour*nal of Emerging Technologies in Learning (iJET), 16(3):274–306, 2021.
- M Laeeq Khan. Social media engagement: What motivates user participation and consumption on YouTube? *Computers in human behavior*, 66:236–247, 2017.
- Youngho Kim, Ahmed Hassan, Ryen W White, and Imed Zitouni. Modeling dwell time to predict click-level satisfaction. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pp. 193–202, 2014.
- Michal Kompan and Mária Bieliková. Content-based news recommendation. In International conference on electronic commerce and web technologies, pp. 61–72. Springer, 2010.
- John Langford and Tong Zhang. The epoch-greedy algorithm for contextual multi-armed bandits. *Advances in neural information processing systems*, 20(1):96–1, 2007.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pp. 661–670, 2010.

- Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextualbandit-based news article recommendation algorithms. In *Proceedings of the Forth International Conference on Web Search and Web Data Mining, WSDM 2011, Hong Kong, China*, pp. 297–306. ACM, 2011.
- Yuanguo Lin, Yong Liu, Fan Lin, Pengcheng Wu, Wenhua Zeng, and Chunyan Miao. A survey on reinforcement learning for recommender systems. *arXiv preprint arXiv:2109.10665*, 2021.
- Pasquale Lops, Dietmar Jannach, Cataldo Musto, Toine Bogers, and Marijn Koolen. Trends in content-based recommendation. User Modeling and User-Adapted Interaction, 29(2):239–249, 2019.
- Hongyu Lu, Min Zhang, and Shaoping Ma. Between clicks and satisfaction: Study on multi-phase user preferences and satisfaction for online news reading. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pp. 435–444, 2018a.
- Yichao Lu, Ruihai Dong, and Barry Smyth. Why i like it: multi-task learning for recommendation and explanation. In *Proceedings of the 12th ACM Conference on Recommender Systems*, pp. 4–12, 2018b.
- Masoud Mansoury, Himan Abdollahpouri, Jessie Smith, Arman Dehpanah, Mykola Pechenizkiy, and Bamshad Mobasher. Investigating potential factors associated with gender discrimination in collaborative recommender systems. In *The Thirty-Third International Flairs Conference*, 2020.
- Patrick Martinchek. 2012-2016 Facebook Posts. https://data.world/martinchek/ 2012-2016-facebook-posts, 2016.
- Jeremy Merrill and Will Oremus. Five points for anger, one for a 'like': How facebook's formula fostered rage and misinformation. *The Washington Post*, 2021. URL https://www.washingtonpost.com/technology/2021/10/26/ facebook-angry-emoji-algorithm/.
- Meta. Reactions now available globally. *Facebook News*. URL https://about.fb.com/ news/2016/02/reactions-now-available-globally/.
- Maxim Naumov, Dheevatsa Mudigere, Hao-Jun Michael Shi, Jianyu Huang, Narayanan Sundaraman, Jongsoo Park, Xiaodong Wang, Udit Gupta, Carole-Jean Wu, Alisson G Azzolini, et al. Deep learning recommendation model for personalization and recommendation systems. *arXiv* preprint arXiv:1906.00091, 2019.
- Nicola Neophytou, Bhaskar Mitra, and Catherine Stinson. Revisiting popularity and demographic biases in recommender evaluation and effectiveness. In *European Conference on Information Retrieval*, pp. 641–654. Springer, 2022.
- Tien T Nguyen, Pik-Mai Hui, F Maxwell Harper, Loren Terveen, and Joseph A Konstan. Exploring the filter bubble: the effect of using recommender systems on content diversity. In *Proceedings* of the 23rd international conference on World wide web, pp. 677–686, 2014.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bertnetworks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019.
- J Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen. Collaborative filtering recommender systems. In *The adaptive web*, pp. 291–324. Springer, 2007.
- Kate Scott. You won't believe what's in this paper! Clickbait, relevance and the curiosity gap. *Journal of pragmatics*, 175:53–66, 2021.
- Donghee Shin. How do users interact with algorithm recommender systems? The interaction of users, algorithms, and performance. *Computers in Human Behavior*, 109:106344, 2020.
- Thiago Silveira, Min Zhang, Xiao Lin, Yiqun Liu, and Shaoping Ma. How good your recommender system is? A survey on evaluations in recommendation. *International Journal of Machine Learning and Cybernetics*, 10(5):813–831, 2019.

- Harald Steck and Dawen Liang. Negative interactions for improved collaborative filtering: Don't go deeper, go higher. In *Fifteenth ACM Conference on Recommender Systems*, pp. 34–43, 2021.
- Harald Steck, Linas Baltrunas, Ehtsham Elahi, Dawen Liang, Yves Raimond, and Justin Basilico. Deep learning for recommender systems: A netflix case study. AI Magazine, 42:7–18, Nov. 2021. doi: 10.1609/aimag.v42i3.18140. URL https://ojs.aaai.org/index.php/ aimagazine/article/view/18140.
- Liang Tang, Yexi Jiang, Lei Li, and Tao Li. Ensemble contextual bandits for personalized recommendation. In Proceedings of the 8th ACM Conference on Recommender Systems, pp. 73–80, 2014.
- Wenjie Wang, Fuli Feng, Xiangnan He, Hanwang Zhang, and Tat-Seng Chua. Clicks can be cheating: Counterfactual recommendation for mitigating clickbait issue. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1288–1297, 2021.
- Qingyun Wu, Huazheng Wang, Quanquan Gu, and Hongning Wang. Contextual bandits in a collaborative environment. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pp. 529–538, 2016.
- Qingyun Wu, Hongning Wang, Liangjie Hong, and Yue Shi. Returning is believing: Optimizing long-term user engagement in recommender systems. In Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, pp. 1927–1936, 2017.
- Tengyang Xie, John Langford, Paul Mineiro, and Ida Momennejad. Interaction-Grounded Learning. In *International Conference on Machine Learning*, pp. 11414–11423. PMLR, 2021.
- Tengyang Xie, Akanksha Saran, Dylan J Foster, Lekan Molu, Ida Momennejad, Nan Jiang, Paul Mineiro, and John Langford. Interaction-Grounded Learning with Action-inclusive Feedback. *arXiv preprint arXiv:2206.08364*, 2022.
- Xing Yi, Liangjie Hong, Erheng Zhong, Nanthan Nan Liu, and Suju Rajan. Beyond clicks: dwell time for personalization. In *Proceedings of the 8th ACM Conference on Recommender systems*, pp. 113–120, 2014.
- Lixin Zou, Long Xia, Zhuoye Ding, Jiaxing Song, Weidong Liu, and Dawei Yin. Reinforcement learning to optimize long-term user engagement in recommender systems. In *Proceedings of* the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 2810–2818, 2019.