# Contextual Bandit Applications in a Customer Support Bot

Sandra Sajeev*
ssajeev@microsoft.com
Microsoft Azure AI

Jade Huang*
jadhuang@microsoft.com
Microsoft Azure AI

Nikos Karampatziakis
nikosk@microsoft.com
Microsoft Azure AI

Matthew Hall
mathall@microsoft.com
Microsoft Azure AI

Sebastian Kochman
sebastko@microsoft.com
Microsoft Azure AI

Weizhu Chen
wzchen@microsoft.com
Microsoft Azure AI

## ABSTRACT

Virtual support agents have grown in popularity as a way for businesses to provide better and more accessible customer service. Some challenges in this domain include ambiguous user queries as well as changing support topics and user behavior (non-stationarity). We do, however, have access to partial feedback provided by the user (clicks, surveys, and other events) which can be leveraged to improve the user experience. Adaptable learning techniques, like contextual bandits, are a natural fit for this problem setting. In this paper, we discuss real-world implementations of contextual bandits (CB) for the Microsoft virtual agent. It includes intent disambiguation based on neural-linear bandits (NLB) and contextual recommendations based on a collection of multi-armed bandits (MAB). Our solutions have been deployed to production and have improved key business metrics of the Microsoft virtual agent, as confirmed by A/B experiments. Results include a relative increase of over 12% in problem resolution rate and relative decrease of over 4% in escalations to a human operator. While our current use cases focus on intent disambiguation and contextual recommendation for support bots, we believe our methods can be extended to other domains.

## CCS CONCEPTS

• **Computing methodologies** → **Reinforcement learning**; *Neural networks*; *Learning linear models*; *Natural language processing*; Batch learning; • **Applied computing** → Enterprise computing.

## KEYWORDS

multi-armed bandits, contextual bandits, chat bots

*Both authors contributed equally to this work.

## 1 INTRODUCTION

Virtual support agents have become ubiquitous in customer service for businesses. The Microsoft virtual agent is a customer support agent with which millions of users engage everyday. A typical session between the virtual agent and the user is shown in figure 1. Challenges that come with operating this system include adapting to changing user behavior and support topics. For instance, bugs associated with the release of new operating system $Y$ can result in users asking more questions about $Y$ rather than previous operating system $X$. In addition, editors will create new support articles about operating system $Y$ to help address common problems; such articles should be suggested by the system when relevant as opposed to outdated content. Another example is that the rise in employees working from home due to a pandemic may result in users asking more questions than before about dual monitor detection or how to connect audio wirelessly. Machine learning, specifically reinforcement learning, has many opportunities to help optimize virtual support agents and meet these challenges head-on.

We drew inspiration from advances in recommendation systems in news [12] and video [6, 17]. The industry standard for recommendation systems has recently been switching to CBs [11], a simplified single-step reinforcement learning (RL) paradigm which requires exploration but does not require dealing with credit assignment. In addition, CBs are data-efficient, which is helpful in our scenario where exploration of riskier actions can hurt the end-user experience and have negative monetary implications. In the context of the virtual agent, we deal with a partial information setting, meaning we only receive feedback (click, survey response, etc.) for the actions the user took. This can easily be modeled with CBs.

Our work presents applications of CBs for two scenarios of the Microsoft virtual agent: intent disambiguation and contextual recommendations. In intent disambiguation, the goal is to provide recommendations that tailor to the user query about a topic. For contextual recommendations, the objective is to provide a set of recommendations to the user as soon as they engage the virtual agent prior to any user query, based on a set of context features. These contextual features can include information such as the website from where the user engages the virtual agent.

We have deployed a MAB solution for contextual recommendations and a NLB solution for intent disambiguation, both of which showed positive impact on key performance indicators (KPIs) in online A/B experiments.

The rest of the paper is organized as follows. First, we present the problem setting and notation. Next, we describe the methods used, going into detail about the learning paradigms. Then, in the implementation section, we outline how we incorporated MABs for
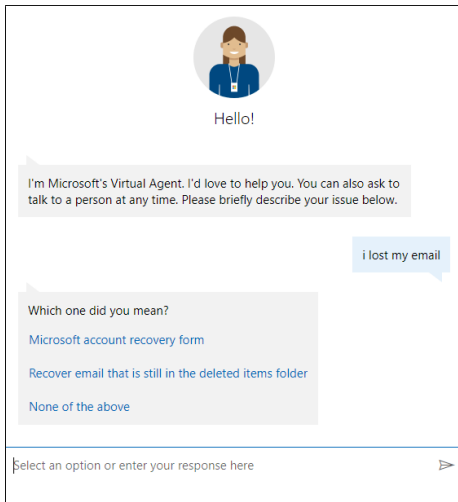
**Figure 1: The Microsoft virtual agent. In response to the user query, an intent clarification dialog presents a *slate* of topics generated by a policy trained using NLBs.**

contextual recommendations and NLBs for intent disambiguation. Finally, we present the efficacy of our solutions through results from online experiments in the evaluation section.

## 2 PROBLEM SETTING

The application motivating our work is the Microsoft virtual agent – an interactive dialogue system providing first-line customer support for many Microsoft products. The virtual agent can be accessed via multiple channels – most commonly through the Microsoft support website[1] or the Get Help app which comes with Windows 10.

In the Microsoft virtual agent's original architecture, manually configured rules dictated the behavior of the system. Most of these rules were not important from a business perspective, but rather assumptions made by application developers due to a lack of data. Moreover the original system could not adapt automatically to the diversity of user intents and a changing environment, such as software updates causing new issues or updated support content. RL addresses some of these challenges and is a good fit for this problem setting for the following reasons:

(1) Microsoft products (including Xbox, Office, Windows, Skype) have a large customer base, hence the virtual agent's incoming traffic is also significant, making RL methods feasible.

(2) We have access to multiple feedback signals including click behavior, escalation to a human agent, and responses to survey questions such as "Did this solve your problem?". Some of these signals are tracked by the product team as their KPIs.

At the same time, the application poses interesting challenges for the use of RL. Goal-oriented dialogue systems need to not only understand the user's original intent, but also be able to carry the state of the dialogue and guide the user towards their goal. While our long-term plan is to be able to have a single RL agent in charge

of the whole conversation, the requirements of such an endeavor, including the of number of samples needed to learn non-trivial policies, are substantial. Therefore, we started by applying RL to the virtual agent in isolated components first, ignoring issues of credit assignment and thus working in the setting of CBs.

In this work, we focus on two specific scenarios within the Microsoft virtual agent. In the following subsections, we describe them in more detail, and then formalize the notation that is shared between both scenarios.

### 2.1 Intent Disambiguation

The domain of customer support, especially for a large company such as Microsoft, is complex and has a significant number of intents. Our intent disambiguation policy is tasked with deciding when the user's query is clear enough to directly trigger a solution, such as a troubleshooting dialogue, or to ask a clarification question.

The inputs to this policy are the statement of the issue by the user (the query), user context features, and a list of candidate actions. The query can range from a few keywords to long and complex sentences or even paragraphs. The user context includes information like the user's operating system and its version (e.g., Office products work on Windows, Mac, iOS, and Android). The candidate actions are a collection of pre-authored dialogue intents or solutions related to the user's query pulled from the Web. Retrieval of these candidates is currently performed by several strategies. One of them includes a deep learning model similar to the one described in [8]. Another uses Bing Custom Search for customized document retrieval. These retrieval components are currently out-of-scope for RL-based optimization so we will focus on the policy that operates with a small list of retrieved candidates.

Given the input, the policy can do one of the following:

- Directly trigger a single intent or a solution: this can start a troubleshooting dialogue or display a rich text solution.
- Ask a yes/no question: "Here's what I think you are asking about: ...Is that correct?"
- Ask a multiple-choice question: "Which one did you mean?", followed by titles of two to four intents as well as option "None of the above."
- Give up: "I'm sorry, I didn't understand. It helps me when you name the product and briefly describe the issue."

Figure 1 presents an example of a multiple-choice question action taken by the disambiguation policy.

### 2.2 Contextual Recommendations

The "Settings" app is a desktop application in Windows 10, where a user can click "Get help" from any Settings page (e.g., Bluetooth, Display etc.) and interact with the Microsoft virtual agent. We will refer to the page the user is coming from as the *source page* - it is a crucial part of the user context that is available to us.

The goal is to provide contextual recommendations, i.e., to recommend a *slate* (sequence) of solution topics before the user types anything, simply based on the context sent by the app (see figure 2). The baseline experience is a fixed mapping from source pages to slates of up to six topics manually picked by human editors.

We were motivated by several reasons to leverage feedback data to learn better contextual recommendations and improve upon the
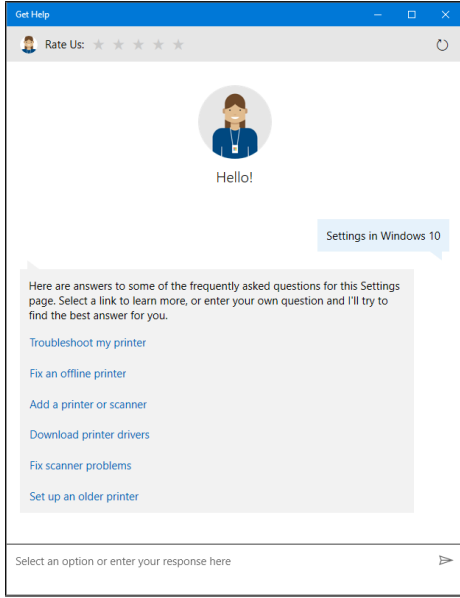
**Figure 2: Contextual recommendations in Settings in Windows 10, for the source page "Printers".**

baseline. First, maintaining manually-specified dialogues does not scale to the approximately 200 source pages present in the app (most of them served the same fixed slate of six topics). Second, a large amount of traffic flows through the app, which suggests that data-driven techniques could do well here. Third, the available user context includes more signals than just source page, which can help with suggesting relevant solutions in certain situations, such as the device network type (wired, wifi) and battery status (charging, discharging).

Given the context, the agent can take one of the following actions:

(1) Suggest up to six solutions out of about 3000 candidates.
(2) Choose to fall back to default behavior of recommending six fixed options chosen by editors for the current source page.

## 2.3 Feedback signals and goals

The feedback signals available to learn from are the same in both scenarios described above:

- **Click.** When a slate of topics is presented, the option selected by the user is recorded. This signal is censored in scenario 2.1 when the policy decides to directly trigger. While useful, it is not a metric important to the business, so we do not optimize for it alone.
- **Survey.** After providing a final answer to the user, the virtual agent asks whether the solution has resolved their problem. The user may respond "yes", "no", or decline to answer. The product team tracks this signal, aggregated per user session and averaged, as the most important KPI called Problem Resolution Rate (more details in section 5.1). Hence, this feedback signal is also our main reward signal.

- **Escalation.** The user can decide at any time to talk to a human agent. This negative reward signal is directly related to the actual cost of running a call center. In general, this metric is monitored but is not being optimized.

## 2.4 Notation

The contextual recommendation and intent disambiguation scenarios that we focus on in this paper can be modeled as combinatorial CBs. We adapt and expand notation proposed in [19] to the scenarios found in the Microsoft virtual agent.

In this setting, an agent interacts with the environment repeatedly as follows:

For time step $t = 1, 2, \ldots$:

(1) The world produces a *context* $x_t \in \mathcal{X}$ (e.g., the user sends a query $x_t$ to the system).
(2) The agent chooses a slate $s_t = (a_1, ..., a_l) \in \mathcal{S}(x_t)$, of variable length $l$ ($0 \le l \le L$), consisting of actions $a_j \in \mathcal{A}(x_t)$. We call $\mathcal{A}(x_t)$ an *item action space* and $\mathcal{S}(x_t)$ a *slate space*. Position $j$ in a slate is called a *slot*.
(3) Given the context and slate, the world produces *feedback signals*, including:
    - $click_t \in \{0, 1\}^l$ informing of the slot selected by the user. In this work, we assume it is either a one-hot vector or a zero vector (i.e., no more than a single item can be selected).
    - $survey_t \in \{\text{yes, no, skipped}\}$ indicating the user's assessment of the solution presented by the virtual agent, provided via a survey. Values "yes" and "no" mean positive and negative feedback, respectively, while "skipped" means the user has not responded to the survey.
    - $escalation_t \in \{0, 1\}$ indicates whether the user has decided to escalate the case to a human agent (1) or not (0).
(4) The reward for each slot $j$ in the slate $s_t$ is denoted by $r_{t,j}$. It is computed as a function of one or more of the *feedback signals*.

Note that the item action space $\mathcal{A}(x_t)$ is not a fixed set but can vary based on the context. The actions are parametric, with features such as title, type of content, etc. In our case, the majority of actions represent support documents and troubleshooting dialogues, but $\mathcal{A}(x_t)$ may also include the special action *null item*, denoted as $\perp$. It gives the user a choice of indicating that none of the other actions is relevant. If selected via $click_t$, the system moves on to other suggestions.

To gain intuition about this notation, let's consider some examples:

- **The intent disambiguation scenario presented in the figure 1**. The query "i lost my email" is the main property of the context $x$. Topics "Microsoft account recovery form" and "Recover email that is still in the deleted items folder" are some actions $a_1$ and $a_2$. The "None of the above" option rendered in the dialog translates to the *null item* $\perp$ - it gives the user an opportunity to explore other topics or ask another question. Hence, the slate chosen by the policy at that time step was $s = (a_1, a_2, \perp)$.

The intent disambiguation policy is also allowed, for some contexts, to directly trigger a topic. Such a slate would look like s = $(a_1)$, without the $\perp$ action. What kind of slates are allowed in what contexts depends on the function $\mathcal{S}(x)$.

- **The contextual recommendations scenario presented in the figure 2**. The source page "Printers" is the main property of the context $x$. Ignoring the presented topics and typing a query manually is equivalent to selecting the *null item* $\perp$ - hence, the slate can be represented as s = $(a_1, a_2, \ldots, a_6, \perp)$.
  In the contextual recommendations scenario, every slate in $\mathcal{S}(x)$, in every context $x$, must always include $\perp$ (i.e.: topics are never directly triggered).

The reward used by our CB algorithms is always a function of feedback signals. Decoupling them from the reward definition, while not a standard approach in the RL literature (as in [5, 11]), is useful in the learning algorithms we are presenting in later sections.

Some feedback signals may be delayed. For example, a long conversation may occur between the action $a_t$ and the point at which we ask the user if the proposed solution addressed their problem (to generate the signal $survey_t$). We currently treat all variations that can happen there as part of the environment that generates a noisy version of the reward. This simplification is a practical choice and can be removed with better modeling of the interaction between the user and the bot.

# 3 METHODS

We divide discussion about the solutions into two main categories: learning from a discrete context (corresponding to the contextual recommendations scenario described in section 2.2) and learning from a rich context (relevant in the intent disambiguation scenario from section 2.1).

## 3.1 Learning with Discrete Context

For each discrete context $x_t \in \mathcal{X}$, let there be two MABs corresponding to two feedback signals: clicks and survey responses. At time step $t$, each MAB can be described as follows:

- We have collected a history of reward successes and failures over $w_x$ total past time steps in each context $x$ separately for each action $a$. This $w_x$ is specific to each context.
- We have a function $Sample(successes, failures)$ that samples scores for all actions.
- Actions are ordered by these sampled scores.
- The slate ends at the position at which $\perp$ is ranked.

For the click bandit and context $x$, we consider

- successes: the total instances action $a$ has been clicked for the past $w_x$ time steps
- failures: the total instances an action $a$ has been observed (ranked above $\perp$) but not clicked

For the survey response bandit and context $x$, we ignore the cases that the survey is not answered, considering

- successes: total "yes" answers
- failures: total "no" answers

*3.1.1 Sampling with a Discrete Context.* To balance exploitation with exploration, we use a sampling algorithm to produce plausible estimates of the probability of a click and the probability of a "yes" survey. We compared both Thompson sampling [16, 20] and EwS [13], and qualitatively we found that results in our application looked better with Thompson sampling so we focus on it here. However, EwS is reasonable to use as well.

Let $\alpha_t^{x,a}$ represent the successes for action $a$ in context $x$, given the data until time step $t$, within window $w_x$. Likewise, let $\beta_t^{x,a}$ represent the failures for action $a$ in context $x$, given the data until $t$, within window $w_x$. We start with a uniform prior distribution $Beta(\alpha = 1, \beta = 1)$ and assume each event is iid. Then the posterior probability is $Beta(\alpha_t^{x,a} + 1, \beta_t^{x,a} + 1)$. In other words, we add one to all success and failure counts. In practice, we track successes $\alpha$ and trials $N$, and failures are then computed as $\beta = N - \alpha$.

Each bandit samples an estimated action value $Q(x_t, a)$ from the posterior distribution $Beta(\alpha_t^{x_t,a} + 1, \beta_t^{x_t,a} + 1)$ for each action $a$. To represent the final score for an action $a$ based on both bandits, we initially combined the scores from the two bandits using a naive product to model the joint probability $P(click) \cdot P(survey = yes|click)$. Later, having amassed more logged data via online experimentation, we switched to a more flexible approach of using a context-specific interpolation score $\lambda$ which weights the click score lower if there are sufficient survey response trials.

$$\log(Q_{joint}(x_t, a)) \triangleq \lambda_x \log(Q_{click}(x_t, a)) + (1 - \lambda_x) \log(Q_{survey}(x_t, a))$$
(1)

This combined score is used to rank actions in descending order. The size of the slate visible to the user is determined by the minimum of the position of $\perp$ in the slate and action space rules dictating the maximum slate length. In other words, the slate that the user actually sees is a subset of the scored actions. All actions with a higher sampled score than $\perp$ is considered observed by the user.

Having assembled the slate, the user then provides feedback signals. As mentioned in 2.4, the *click* signal is a one-hot vector: 1 for the clicked slot and 0 for the rest. The *survey* signal is considered observed only for the slot that was clicked.

Note that the downside of Thompson sampling in comparison to EwS is that it does not produce closed-form probabilities of choosing an action in a given state. Collecting such probabilities is often useful for off-policy learning and evaluation methods. This can be mitigated by logging the $\alpha$ and $\beta$ values at the time of making a decision and estimating these probabilities offline in a Monte-Carlo fashion.

## 3.2 Learning with Rich Context

A rich context is defined by large cardinality of the set $\mathcal{X}$, which makes it difficult to learn the value function for each state. Therefore value function approximation is required for rich contexts, like natural language embeddings and image data. Assume there is a representation function $\phi$ that takes in a rich context $x \in \mathcal{X}$ and action $a \in \mathcal{A}$ and outputs a set of low-dimensional features, $\phi(x, a) \in \mathbb{R}^d$. We assume that the expected reward is linear with respect to these features.

$$\mathbb{E}[r|x, a] = w^\top \phi(x, a)$$
(2)

In other words, there exists an unknown vector $w$, which models the linear relationship of the expected reward $\mathbb{E}[r|x, a]$, for action $a$ and a context $x$.

*3.2.1 Neural-Linear Bandit.* Correct quantification of uncertainty is critical in bandit algorithms because the algorithm is in charge of how data is collected. If the algorithm uses uncertainty estimates that are too wide, it over-explores and unnecessarily tries suboptimal actions. If the uncertainty estimates are too narrow, the algorithm under-explores and runs the risk of not discovering the best action.

In this work we use the Neural-Linear Bandit strategy from [15], which compares various methods for quantifying uncertainty with deep learning models. There it was demonstrated that the neural-linear approach is often the best single-model method (i.e., not relying on an ensemble) for uncertainty quantification. We have slightly adapted this method for our purposes. The NLB consists of two main parts. First is the representation function $\phi : \mathcal{X} \times \mathcal{A} \mapsto \mathbb{R}^d$ which produces a $d$-dimensional vector of neural features. The function $\phi$ is generated by fitting a neural network, $N : \mathcal{X} \times \mathcal{A} \mapsto \mathbb{R}^1$, which predicts the observed reward from query-answer pairs. The last layer of this network is removed to represent $\phi(x, a)$.

The second part is the bandit function, where we try to estimate the $w$ that satisfies (2) and the uncertainty around it. This is the simplifying assumption in NLBs: all uncertainty in the model can be approximated just by the uncertainty in the last layer of the network. Due to the assumption that the reward is linear with respect to the neural features, the focus is on creating a posterior distribution for the solution. Note that validating this assumption becomes difficult in practice due to the large dimensions of the neural features.

$$\hat{w}_t = \underset{w}{\arg\min} \sum_{i=1}^{t-1} (\phi(x_i, a_i)^\top w - r_i)^2 \qquad (3)$$

This is a simple least squares problem that we need to solve each time we update the bandit. To do this efficiently and support rolling windows of training data we maintain the sufficient statistics:

$$f_t = \sum_{i=1}^{t-1} r_i \phi(x_i, a_i) \qquad (4)$$

$$B_t = \sum_{i=1}^{t-1} \phi(x_i, a_i)\phi(x_i, a_i)^\top \qquad (5)$$

which help us rewrite eq. (3) as $B_t \hat{w}_t = f_t$. We solve the latter with the help of the Cholesky factorization $L_t L_t^\top = B_t$ which we also use for obtaining the uncertainty for new predictions $\hat{r}(x, a) = \phi(x, a)^\top \hat{w}_t$. To do this we need a notion of a count of how many times action $a$ has been tried in context $x$. For linear bandits this is given by

$$b(a) = \frac{1}{\phi(x, a)^\top B_t^{-1} \phi(x, a)} = \left\| L^{-1}\phi(x, a) \right\|^{-2} \qquad (6)$$

The bandit gives more benefit of doubt to an action $a$ with small $b(a)$. Similarly to the number of times an action has been tried in the case of a MAB, $b(a)$ is unit-less even if $\phi(x, a)$ has units: the matrix $B_t^{-1}$ cancels them. The inverse of the quadratic form is also necessary for other reasons. If $\phi$ were one-hot vectors, then $B_t$

would be diagonal with the count for action $a$ in context $x$ on the diagonal. Therefore, the inverse quadratic form is the count.

*3.2.2 Sampling with a Rich Context.* The NLB policy samples actions $a$ to generate the slate $s$ according to a sampling algorithm. The main algorithms we explored for sampling in the rich context are the linear versions of Thompson sampling [3] and EwS [1]. Both sampling algorithms result in stochastic policies.

Linear Thompson sampling [3] requires specification of a prior covariance matrix for $w$, typically chosen to be a multiple of the identity matrix. However, there is no easy way to set this hyperparameter as it can have a hard-to-predict effect on the learning dynamics. In addition, this algorithm doesn't provide the probability of sampling action $a$ in an explicit form which is useful for off-policy evaluation.

Linear EwS [1] was chosen since it has no requirement of setting a prior. In Linear EwS, the probability of selecting an action $a$ has an explicit form. Given context $x$ we first compute the action the model currently prefers

$$a^* = \underset{a \in \mathcal{A}}{\arg\max} \; \hat{w}^\top \phi(x, a)$$

and the gaps between the best action and all actions

$$g(a) = \hat{w}^\top \phi(x, a^*) - \hat{w}^\top \phi(x, a) \qquad (7)$$

Finally, as noted in [1], Linear EwS samples action $a$ proportional to $\exp(-2b(a)g(a)^2)$ where $b(a)$ is given by (6).

## 3.3 Inference Strategies

*3.3.1 Making Exploration Safer.* If a strong baseline already exists, to make exploration safer, we can compare the total value of the sampled slate to the value of the baseline slate. Whichever slate has a higher value is then presented to the user. In our case, we consider the slate value to be the sum of the sampled scores of the actions in the slate. This simple trick can reduce variance in the KPIs, normally expected from exploration, after initial deployment of an automatic learning system. Note that the baseline slate can be still used to collect feedback signals and improve future predictions.

*3.3.2 Sampling to Obtain a Slate.* One can simply sample for expected rewards from possible actions once to obtain a full slate of actions. If there are per-slot rules dictated by business requirements, they can be applied after the sampling.

## 4 IMPLEMENTATION

This section outlines key implementation details about our MAB solution for contextual recommendations and NLBs for intent disambiguation. Both training pipelines were written in Python. All training and evaluation was done asynchronously using Azure Batch. The inference code was written in C# to easily integrate with our partner team's product, the Microsoft virtual agent. The intent disambiguation models were trained using PyTorch [14] and converted to ONNX [4] for inference in the production system. ONNX Runtime is used to load the models in .NET.

## 4.1 Contextual Recommendations

*4.1.1 Problem Formulation.* The discrete context in our case is categorical: the source page from where the user is clicking "Get
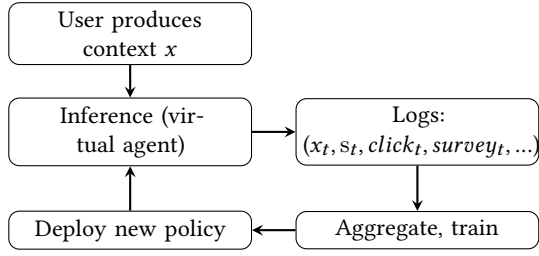
**Figure 3: Feedback loop for contextual recommendations.**

help". In offline analysis, we found that while there are other interesting attributes in the context such as battery value and network type, the source page is the most indicative of what the user may free-type when they choose not to select from the given list of "top solutions". For example, users clicking "Get help" from the Printers page tend to ask printer-related questions. The candidate content is represented simply by a unique id and no other features. We associate the null item ⊥ with the user free-typing.

User clicks and surveys are logged and aggregated every four hours using asynchronous jobs over a moving window. This moving window ensures that content that is not performing well and not interacted with will disappear from recommendations eventually. Note that this window is context-specific, as content performing poorly in one context does not have an effect on another context.

At runtime, a user clicks on "Get help" from source page $x$. First, we gather available candidates that have click and survey counts. We have limitations on the size of our logged data, so we perform an initial sampling of the candidates from over 1000 candidates down to around 25 using Thompson sampling as described in 3.1.1.

With this pared-down set of candidates, we sample once more using Thompson sampling to obtain a final slate of actions. The length of the slate is determined by the minimum of the following: the null item's position and the configured max length of the slate $L$ (in our case $L = 7$ including ⊥). We apply the trick of making exploration safer from 3.3.1 as there is a strong human-authored baseline.

*4.1.2 Automatic Action Space Expansion.* So far as we've described, this experience has a concrete set of candidates per source page from which it draws to recommend to the user. Occasionally, an editor may author new content relevant to a source page which performs well within the context of a source page but outside of the initial recommendation, e.g., the user free-types a query that triggers this content, finds it helpful, and answers 'yes' to the end-of-chat survey. Thanks to our work with intent disambiguation, these interactions are also logged.

To address this, we have another asynchronous job that runs once a day. This job merges statistics from our intent disambiguation flow that occur within the context of a source page with statistics from the initial Settings recommendations flow. For each source page, we use the null item's score as a success threshold for candidates in consideration from intent disambiguation. We first filter out candidates with insufficient trials, e.g., 10 trials. For each potential candidate $a$ we then compute a bound on the probability that the expected reward of $a$ is larger than the expected reward of ⊥.

We then filter out candidates whose bound on the probability is less than our configured allowable false positive probability $fp$. Intuitively, this allowable false positive probability is the maximum probability we allow that an individual new candidate has lower survey response success than ⊥, i.e., the probability of a candidate not being a false positive is $\geq 1 - fp$. The qualifying candidates' counts are copied over from the intent disambiguation space to the Settings space–however, we zero out the click counts as the trials between the two spaces are not comparable.

## 4.2 Intent Disambiguation

*4.2.1 Problem Formulation.* For the intent disambiguation scenario, the context is rich. The input data to the policy contains the text embedding of the user query and candidate title generated in a manner similar to [8]. It also includes a collection of categorical user context features, such as the website from where the user initiated the virtual agent. When the user engages the virtual agent by typing a question, our policy samples from a set of candidates described in 3.3.2 to generate a slate. The slate length is determined by the null item.

*4.2.2 Training.* To train the NLB policy, we first featurize the logged data, limiting to sessions that have explicit user survey responses ("yes" or "no"). We initially defined the reward to be a function of the survey feedback of only the clicked action (8).

$$r_t = \begin{cases} 1, & survey_t = \text{"yes"} \\ -1, & survey_t = \text{"no"} \end{cases} \tag{8}$$

In subsequent iterations, we experimented with combining multiple feedback signals for the reward as described in section 5.5. After featurization, we follow the learning procedure for NLBs defined in 3.2.1. The training objective of the network $N$ was set to minimize the squared error between the reward of the clicked action and the prediction. This network is used to generate $\phi$ with $d = 2048$ neural features. Note that $d$ cannot be excessively large as the bandit needs to store $B_t^{-1}$ or $L_t^{-1}$ which is of size $d \times d$ to compute eq. (6).

We train the bandit in fixed periods (e.g., one day). For each fixed period, we solve the least squares problem defined in (3). When training this model in practice, the matrix $B_t$ can be very close to low rank. To be able to have a usable inverse, we used principal component regression and retained enough principal components to capture 99% of the variation in the neural features.

The initial NLB policies were fixed: trained on a set number of days of data and deployed to production. We later explored auto-updating policies with hourly retraining.

## 5 EXPERIMENT RESULTS

### 5.1 Metrics

The virtual agent product team tracks several KPIs where some were the metrics for which we optimized. Our metrics are reported over user sessions and each user session is treated independently. We do not track an individual user's behavior over time, e.g. if a user re-engages with the Microsoft virtual agent, that is treated as a separate user session.

- **Problem Resolution Rate (PRR):** a ratio of user sessions with a positive survey response to all sessions with any

survey response. This is the main metric we track while running offline or online experiments.

- **Engaged Assisted Support (EAS):** fraction of user sessions that end up with escalation to human-assisted support.
- **Self Help Success (SHS):** fraction of user sessions that ended with assumed success (i.e., an answer is delivered and the user does not engage assisted support or indicate the answer was not useful).
- **User Engagement (UE):** fraction of user interactions where the user sends at least one message after the virtual agent's greeting (where the greeting may include a slate with contextual recommendations, in case of help for the Settings app).

We can approximate some of these metrics using feedback signals described in sections 2.3 and 2.4. E.g., for a time period $t = 1, ..., T$ we can define the following approximations:

$$P\hat{R}R_{1:T} = \frac{\sum_{t=1}^{T} \mathbb{1}\left[survey_t = yes\right]}{\sum_{t=1}^{T} \mathbb{1}\left[survey_t = yes \vee survey_t = no\right]}$$

$$E\hat{A}S_{1:T} = \frac{1}{T}\sum_{t=1}^{T} escalation_t$$

Where $\mathbb{1}$ is an indicator function.

The difference between the product team's KPIs and the approximated values is that the real metrics are aggregated across a user session rather than across an individual event within a user session. A single session lasts from the moment the user connects with the virtual agent until the last event that occurs within the conversation. Specific aggregation depends on the metric - e.g., the user may respond negatively to one survey, but then continue conversation with the virtual agent, and eventually respond to the second survey positively. The product team's KPI would indicate this as a single event with success contributing towards PRR, while our simplified definition would consider this as two events - one failure and one success.

Our CB algorithms, as currently defined, are capable of optimizing only the approximated metrics like $P\hat{R}R$ or $E\hat{A}S$. While aware of this problem, we observe that a CB policy optimizing a simplified metric will usually cause movement of the real KPI in the same direction, when deployed to production, as the following sections show.

## 5.2 Offline Evaluation

Before attempting online experimentation with the NLB policy for intent disambiguation, we first used offline evaluation to check whether a newly trained policy seemed promising. This helped limit the negative impact to actual users in production traffic. Since we made use of EwS sampling, our behavior policy was stochastic and we logged the probabilities. This allows us to run off-policy evaluations to estimate the impact of our policy.

We use the SNIPS [18] estimator to measure how much more reward our current policy would have provided compared to the logging policy. In our counterfactual evaluation, we were only able to use data where the logged policy and the new policy agreed on the clicked action, rather than the entire slate. The primary metrics we observed were how the logged policy and new policy compared in terms of rewards based on the survey and deflection feedback signals. A policy had to achieve higher reward estimates than the logging policy with low variance to be promoted to an online experiment. The requirement of low variance was enacted to reduce mismatch between offline evaluation and online experiments.

In the case of contextual recommendations, we do not perform offline evaluation beyond qualitative analysis as described in 3.1.1 and careful monitoring of the auto-updated results. The updated model is deployed automatically without any experiments or human in the loop, although with an automated monitoring system alerting about unusual behavior.

## 5.3 A/B experimentation

We conducted several online A/B experiments to verify the efficacy of both our new contextual recommendation and intent disambiguation policies in comparison to a control policy (i.e., the current deployed production policy) on user traffic from the Microsoft virtual agent. Each A/B experiment should be viewed independently, as the control policy and the amount of user traffic varies from experiment to experiment. The A/B experiments gave key insights into how well the polices improved KPIs relative to control. Some of these KPIs were only approximated by offline evaluation and others were not captured in offline evaluation, like UE or SHS (see section 5.1). Furthermore, A/B experiments reflected the actual environment in which the policy was deployed, which differed in some cases from the training data. For instance, the NLB policies were trained only with data that included the survey signals, "yes" or "no", leaving out data where the user did not answer the survey. Finally, these A/B experiments allowed us to quantify impact, which was especially key for contextual recommendations which did not have formal offline evaluation. The results from these A/B experiments were the deciding factor in whether to promote these new policies to the default experience, thus serving general production traffic in the future. When looking at results from online A/B experiments, there were several KPIs that we observed as described in section 5.1.

Optimizing for several KPIs was a balancing act. Some of these metrics like PRR were almost directly optimized for as reward signals during training, others like EAS were used as limitations where we sought to prevent harmful movement, and others such as SHS and UE were assumed to be correlated to reward signals.

## 5.4 Contextual Recommendations

We first ran an A/B experiment for an automatically updating implementation of our MAB approach based on Thompson sampling that used the naive combined score $P(click) \cdot P(survey = yes|click)$ described in section 3.1.1. The control was a human-authored fixed mapping from source pages to slate. The experiment showed gains in SHS and UE, which led us to deploy this new policy to all production traffic.

Afterwards, we set up a small experiment to monitor the auto-updating policy's performance over time. The MAB policy was set as the champion policy and the old control behavior as the challenger. The challenger only ran with a small amount of traffic as to provide minimal impact on performance. Running this monitoring experiment was a crucial decision, as over time the champion policy

began to show negative movement in PRR as well as smaller gains in SHS and UE. As PRR is our primary KPI, gains in SHS and UE was not worth a drop in PRR.

Motivated in part by the negative results of the monitoring experiment, we worked on implementing a more dynamic interpolation scoring approach described by equation 1. We ran an A/B experiment for a second automatically updating MAB based on this approach. This experiment also showed positive movement in SHS and UE for the treatment compared to the control without any drop in PRR, as shown in table 1.

**Table 1: Dynamic Thompson Sampling Policy A/B Results**

| Metric | Movement | P-Value | Sample Size |
| --- | --- | --- | --- |
| PRR | +1.13% | 0.2384 | 36K |
| EAS | -1.73% | 0.3656 | 193K |
| SHS | +2.93% | $< 10^{-10}$ | 193K |
| UE | +2.13% | $< 10^{-10}$ | 272K |

This second policy was then deployed to production, replacing the previously deployed MAB policy. Afterwards, we set up another small experiment to monitor the new auto-updating policy's performance over time. Over several weeks, the policy began to show improvements in PRR and EAS as well as greater gains in SHS and UE as seen in table 2. It is important to note that a reduction in EAS is a positive improvement, as it is reducing the number of escalations to a human agent and practically, reducing cost.

**Table 2: Dynamic Thompson Sampling Policy Monitoring A/B Results**

| Metric | Movement | P-Value | Sample Size |
| --- | --- | --- | --- |
| PRR | +2.36% | $5 \cdot 10^{-5}$ | 104k |
| EAS | -4.09% | $2 \cdot 10^{-6}$ | 549K |
| SHS | +4.70% | $< 10^{-10}$ | 549K |
| UE | +4.54% | $< 10^{-10}$ | 766k |

## 5.5 Intent Disambiguation

In the first A/B experiment we ran with our NLB policy, our initial NLB policy showed improvement over the then-current policy, a deep network trained via policy gradient, where both policies had the same input. The greatest improvement was in PRR, which showed a large increase over our previous policy, as shown in table 3. There was also a small improvement in SHS.

**Table 3: Neural-Linear Bandit Policy A/B Results**

| Metric | Movement | P-Value | Sample Size |
| --- | --- | --- | --- |
| PRR | +12.45% | $< 10^{-10}$ | 94K |
| EAS | -1.12% | 0.0599 | 473K |
| SHS | +0.49% | 0.0297 | 473K |
| UE | +0.07% | 0.5374 | 632K |

We also ran an A/B experiment that compared our NLB policy with a greedy rules-based policy. This rules-based policy was the original solution for intent disambiguation in the Microsoft virtual agent. The experiment results showed that while the NLB policy improved PRR, the greedy rules-based policy had a better EAS score. This is likely because the NLB policy can determine if there is no candidate that matches the user's intent, and thus shows no result more frequently, leading to escalations.

After promoting the new NLB policy to be the default policy for production traffic, we experimented with reducing the intent disambiguation slate size from four to three. The objective of reducing slate size was to decrease the likelihood of a user getting distracted by suboptimal candidates in the slate. This change was tested in an A/B experiment against a control of the NLB policy with a slate size of four. The experiment results showed very little movement amongst all KPIs. This gave us the confidence to reduce our slate size to three for all production traffic, also helping reduce the action space of our intent disambiguation policy.

To study the effect of optimizing different feedback signals, we ran an experiment with a model optimized for EAS in addition to PRR with the hope of lowering EAS. This A/B experiment shows promising results with improvements in both EAS and PRR as seen in table 5. Interestingly, there was also a drop in SHS. We plan to pursue this method of training in the future.

**Table 4: Neural-Linear Bandit Policy Optimized for Escalation A/B Results**

| Metric | Movement | P-Value | Sample Size |
| --- | --- | --- | --- |
| PRR | +4.66% | $< 10^{-10}$ | 42K |
| EAS | -2.29% | 0.0160 | 263K |
| SHS | -1.72% | $6 \cdot 10^{-7}$ | 263K |
| UE | -0.01% | 0.9295 | 363K |

To achieve the objective of adapting to changing user traffic, we began experimenting with auto-updating policies. For this experiment, the control was the NLB policy with a slate size of three. The treatment auto-updating NLB policy was initialized with the same parameters as the control aside from a higher exploration rate. The higher exploration rate was chosen to counteract the overfitting effects of training with less data. Then, we retrained NLB's bandit-layer hourly on new data from the treatment traffic. There were promising improvements in PRR over the two-week experiment. With this experiment, we successfully closed the loop for the intent disambiguation scenario using the auto-updating NLB policy.

**Table 5: Neural-Linear Auto-Updating A/B Results**

| Metric | Movement | P-Value | Sample Size |
| --- | --- | --- | --- |
| PRR | +1.48% | 0.0307 | 90K |
| EAS | -0.5% | 0.4356 | 590K |
| SHS | -0.0014% | 0.1450 | 590K |
| UE | +0.13% | 0.1936 | 795K |

## 6 RELATED WORK

In recent years, there have been many advances in using RL for recommendation systems. Policy gradient-based methods are one such approach to model recommendation systems, such as REINFORCE [6]. This method is more stable in regards to policy convergence compared to Q-learning techniques.

Recent work on slate recommendation like SlateQ [9] allows the decomposition of a slate's total value as a function of the items. In the SlateQ approach, a Q-Network is trained with a user choice model that approximates the click-through rate. One caveat is that SlateQ makes significant assumptions about user behavior through the use of a user choice model.

CBs are a lightweight single-step RL method that works well with interactive systems. It is the augmented form of a K-armed bandit problem with context $x$ as formalized in [11]. In addition, CBs offer key properties that simplify the problem space. For instance, the reward is only associated with the selected action. CBs enable exploration via sampling algorithms like Thompson sampling, epsilon-greedy, or EwS [11, 13, 20]. Furthermore, CBs have been successfully utilized in recommendation systems. In [12], MABs modeled personalized news article recommendation with LinUCB sampling. CBs have also been used in search engines for ranking recommendations [19].

Inspired by the application of CBs in traditional recommendation systems, we extend these ideas to be applied in intent disambiguation and contextual recommendations for a virtual support agent. For the implementation of our CB solutions, we extended Decision Service [2], a simple interface for any developer to use RL. To learn more details about the implementation of our RL system and practical lessons, refer to [10].

## 7 CONCLUSIONS AND FUTURE WORK

In this work, we describe two applications of CBs in the Microsoft virtual agent in the areas of intent disambiguation and contextual recommendations. We have demonstrated the efficacy of these solutions through online A/B experiments in the Microsoft virtual agent. The NLB solution for intent disambiguation provided 12.45% improvement in the primary KPI of the Microsoft virtual agent, PRR. Our MAB models for contextual recommendations also increased PRR by 2.36% in addition to increasing SHS and UE, while lowering EAS (escalations to a human agent). We have since deployed these solutions to production and impact millions of users each day.

For future work, we are exploring how we can apply the insights from MABs and NLBs to multi-domain support bots as well as other products. We are also exploring real-world applications of episodic RL where it is important to do proper credit assignment. For this we plan to rely on a reduction approach [7] which can operate well with the rest of our existing system.

## REFERENCES

[1] Yasin Abbasi-Yadkori. 2013. Online learning for linearly parametrized control problems. *Ph.D. thesis* (2013).

[2] Alekh Agarwal, Sarah Bird, Markus Cozowicz, Luong Hoang, John Langford, Stephen Lee, Jiaji Li, Dan Melamed, Gal Oshri, Oswaldo Ribas, et al. 2016. Making contextual decisions with low technical debt. *arXiv preprint arXiv:1606.03966* (2016).

[3] Shipra Agrawal and Navin Goyal. 2013. Thompson Sampling for Contextual Bandits with Linear Payoffs. In *Proceedings of The 30th International Conference on Machine Learning*. 127–135.

[4] Junjie Bai, Fang Lu, Ke Zhang, et al. 2019. ONNX: Open Neural Network Exchange. https://github.com/onnx/onnx.

[5] Olivier Chapelle and Lihong Li. 2011. An Empirical Evaluation of Thompson Sampling. In *Advances in Neural Information Processing Systems 24*, Vol. 24. 2249–2257.

[6] Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. 2019. Top-K Off-Policy Correction for a REINFORCE Recommender System. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. ACM, 456–464.

[7] Hal Daumé III, John Langford, and Amr Sharaf. 2018. Residual Loss Prediction: Reinforcement Learning With No Incremental Feedback. In *International Conference on Learning Representations*. https://openreview.net/forum?id=HJNMYceCW

[8] Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2020. Poly-encoders: Architectures and Pre-training Strategies for Fast and Accurate Multi-sentence Scoring. In *ICLR 2020 : Eighth International Conference on Learning Representations*.

[9] Eugene Ie, Vihan Jain, Jing Wang, Sanmit Narvekar, Ritesh Agarwal, Rui Wu, Heng-Tze Cheng, Tushar Chandra, and Craig Boutilier. 2019. SlateQ: A Tractable Decomposition for Reinforcement Learning with Recommendation Sets. In *Proceedings of the Twenty-eighth International Joint Conference on Artificial Intelligence (IJCAI-19)*. Macau, China, 2592–2599. See arXiv:1905.12767 for a related and expanded paper (with additional material and authors).

[10] Nikos Karampatziakis, Sebastian Kochman, Jade Huang, Paul Mineiro, Kathy Osborne, and Weizhu Chen. 2019. Lessons from Contextual Bandit Learning in a Customer Support Bot. arXiv:1905.02219 [cs.LG]

[11] John Langford and Tong Zhang. 2007. The epoch-greedy algorithm for contextual multi-armed bandits. In *Proceedings of the 20th International Conference on Neural Information Processing Systems*. Citeseer, 817–824.

[12] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. 2010. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*. ACM, 661–670.

[13] Odalric-Ambrym Maillard. 2011. *APPRENTISSAGE SÉQUENTIEL: Bandits, Statistique et Renforcement.* Ph.D. Dissertation. Université des Sciences et Technologie de Lille-Lille I.

[14] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 8024–8035. http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf

[15] Carlos Riquelme, George Tucker, and Jasper Snoek. 2018. Deep Bayesian Bandits Showdown: An Empirical Comparison of Bayesian Deep Networks for Thompson Sampling. arXiv:1802.09127 [stat.ML]

[16] Daniel Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. 2020. A Tutorial on Thompson Sampling. arXiv:1707.02038 [cs.LG]

[17] Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims. 2016. Recommendations as treatments: Debiasing learning and evaluation. *arXiv preprint arXiv:1602.05352* (2016).

[18] Adith Swaminathan and Thorsten Joachims. 2015. The Self-Normalized Estimator for Counterfactual Learning. In *Advances in Neural Information Processing Systems* (advances in neural information processing systems ed.), Vol. 28. Curran Associates, Inc., 3231–3239. https://www.microsoft.com/en-us/research/publication/self-normalized-estimator-counterfactual-learning/

[19] Adith Swaminathan, Akshay Krishnamurthy, Alekh Agarwal, Miroslav Dudík, John Langford, Damien Jose, and Imed Zitouni. 2017. Off-policy evaluation for slate recommendation. arXiv:1605.04812 [cs.LG]

[20] WILLIAM R THOMPSON. 1933. ON THE LIKELIHOOD THAT ONE UNKNOWN PROBABILITY EXCEEDS ANOTHER IN VIEW OF THE EVIDENCE OF TWO SAMPLES. *Biometrika* 25, 3-4 (12 1933), 285–294. https://doi.org/10.1093/biomet/25.3-4.285 arXiv:https://academic.oup.com/biomet/article-pdf/25/3-4/285/513725/25-3-4-285.pdf