

Making Contextual Decisions with Low Technical Debt

Alekh Agarwal, Sarah Bird, Markus Cozowicz, Luong Hoang, John Langford, Stephen Lee*, Jiaji Li*
Dan Melamed, Gal Oshri*, Oswaldo Ribas*, Siddhartha Sen, Alex Slivkins

Microsoft Research, *Microsoft

Abstract

Applications and systems are constantly faced with decisions that require picking from a set of actions based on contextual information. Reinforcement-based learning algorithms such as contextual bandits can be very effective in these settings, but applying them in practice is fraught with technical debt, and no general system exists that supports them completely. We address this and create the first general system for contextual learning, called the Decision Service.

Existing systems often suffer from technical debt that arises from issues like incorrect data collection and weak debuggability, issues we systematically address through our ML methodology and system abstractions. The Decision Service enables all aspects of contextual bandit learning using four system abstractions which connect together in a loop: explore (the decision space), log, learn, and deploy. Notably, our new explore and log abstractions ensure the system produces correct, unbiased data, which our learner uses for online learning and to enable real-time safeguards, all in a fully reproducible manner.

The Decision Service has a simple user interface and works with a variety of applications: we present two live production deployments for content recommendation that achieved click-through improvements of 25-30%, another with 18% revenue lift in the landing page, and ongoing applications in tech support and machine failure handling. The service makes real-time decisions and learns continuously and scalably, while significantly lowering technical debt.

1 Introduction

Machine learning has well-known high-value applications, yet effective deployment is fraught with difficulties in practice [12, 43]. Machine-learned models for recommendation, ranking, and spam detection all become obsolete unless they are regularly infused with new data. Using data is actually quite tricky—thus the growing demand for data scientists. How do we create a machine learning system that collects and uses data in a safe and correct manner?

Recently, Sculley et al. [43] used the framework of *technical debt*—the long-term costs that accumulate when expedient (but suboptimal) decisions are made in the short run—to argue that ML systems incur massive hidden costs at the system level, beyond the basic code complexity issues of traditional software systems. Thus, traditional mechanisms for coping with technical debt, such as refactoring code or improving unit tests, are insufficient for ML systems.

In this work, we create a *system for contextual decision making* that addresses an important subset of these failure

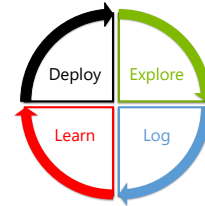


Figure 1: Complete loop for effective contextual learning, showing the four system abstractions we define.

modes and technical debts. Prior experiences of deploying one-off solutions [4, 33, 34, 12], and failures we have experienced ourselves, call for a systematic approach that addresses these debts by design. Specifically, we target failure modes and technical debts incurred by: feedback loops and bias, distributed data collection, changes in the environment, and weak monitoring and debugging (§2). Existing systems neither solve the class of problems we handle, nor address these debts.

As it turns out, the problem domain we are interested in—contextual decisions in interactive settings—leads to higher incidence of failure modes and technical debt than traditional supervised prediction problems. In a traditional prediction problem, a context (*e.g.*, image features) is given, a prediction (*e.g.*, dog or cat) is made and the true label is then revealed. Crucially, the label can be used to measure the quality of *any potential prediction* on this context, which is key in the supervised learning techniques used for this setting. In our setting, an application repeatedly takes *actions* (*e.g.*, which type of news article to show) when faced with a particular *context* (*e.g.*, user demographics) to obtain a desirable outcome quantified as a *reward* (*e.g.*, a click). The goal is to find a good *policy* mapping contexts to actions, *e.g.*, show politics articles to adults and sports articles to teenagers. The feedback mode puts the algorithm in a reinforcement learning setting, which is significantly more challenging [2] than supervised learning [20]. This paradigm covers a huge range of applications including virtually every online service; Table 2 has examples.

Two properties make this setting particularly challenging:

1. The reward is observed only for the chosen action; nothing is learned about unexplored actions, leading to *partial feedback*. In supervised learning, the true label determines the quality of every prediction leading to *full feedback*.
2. Rewards arrive after a delay, often via a separate data path.

• **Partial feedback.** Partial feedback settings are prone to bias, because the decisions made by a system affect the data

collected to train the system [12, 43]. For example, it is not possible to learn that a news article about sports generates more clicks than one about politics when shown to a particular demographic, without showing them both types of articles at least some of the time. In addition, the demographic’s preferences may change over time. Hence, there is a need to *explore* so as to acquire the right data; a *biased* dataset, no matter how large, does not enable good learning.

- **Delayed rewards.** In most applications, reward information (*e.g.*, did the user click the article?) arrives after a considerable delay, ranging from seconds to days. Often, this data is collected by a separate subsystem and follows a separate data path. Since an event is incomplete until the reward is known, this requires distributed data collection, which is a common source of data errors. In addition, since rewards for some actions take longer to arrive than others, there is a real possibility of delay-related bias entering the dataset.

Addressing these issues requires a synthesis of ML techniques as well as careful system design. On the ML side, exploration addresses some of the issues caused by partial feedback. A common methodology for exploration is *A/B testing* [28, 29]: policy A is tested against policy B by running both live on a percentage of user traffic drawn at random. This requires data scaling linearly with the number of policies to be tested. On the other hand, *contextual bandits* [3, 31] is a type of machine learning that allows testing and optimization over exponentially more policies for the same amount of data¹. Moreover, the policies being evaluated need not be deployed or even known ahead of time, saving vast business and engineering effort. We refer to this dramatic improvement in capability as *multiworld testing* (MWT), and realize it with contextual bandits and policy evaluation techniques [33, 34].

Our ML methodology has two important properties: it addresses the biased feedback loop and enables advanced monitoring/debugging through its policy evaluation capabilities (§3). It does not, however, fully address a nonstationary environment nor the data collection issues caused by delayed rewards. In fact, the ML theory simply assumes correct data is provided as input, which is hard to ensure in practice.

We argue that well-defined system abstractions are required for correct data collection. We propose the 4-step loop shown in Fig. 1: *explore* to collect the data, *log* this data correctly, *learn* a good model, and *deploy* it in the application. Existing systems typically focus on the learn and deploy steps, while ignoring or mishandling the data collection steps (explore and log); further, the steps are managed by independent processes, resulting in the complete process being inefficient, expensive, and error-prone. By connecting these abstractions in the right way, we address the remaining sources of failures and technical debt (§4.2). The explorer invokes the logger at the time of decision to record the decision event, and later to join the reward. To avoid delay bias,

the logger imposes a uniform delay before releasing joined data to the learner. Finally, models are trained and deployed online to cope with dynamic environments. The abstractions are modular, and can be implemented, maintained and tested separately so as to not add new sources of debt resulting from a monolithic system.

Our system abstractions and ML methodology enable other techniques for reducing technical debt, such as full reproducibility and real-time safeguards (§4.3). The resulting system is called the Decision Service: it is fully functional, publicly available, and open-sourced [16]. It exposes a simple interface and supports several deployment options (§5), each of which has been exercised by a real deployment. We describe three live production deployments for content recommendation in MSN, Complex, and TrackRevenue, and two ongoing applications in assisting tech support staff and cloud machine failures (§6). Each deployment has achieved a double-digit lift between 14–25% in their target metric (*e.g.*, clicks, revenue) relative to strong baselines.

Our evaluation (§7) shows that the Decision Service makes decisions with low latency, incorporates new data into deployed models quickly, and scales to the traffic of our customers and beyond. More importantly, we show through experiments on production data and anecdotal experiences, that the service reduces the technical debts we sought to address.

In summary, we make the following contributions:

- We implement an ML methodology that achieves MWT capability and eliminates important failures by design (§3).
- We define four system abstractions for realizing this methodology in a robust manner (§4.2). We describe advanced techniques that further reduce technical debt, such as full reproducibility and real-time safeguards (§4.3).
- We present the first general-purpose service for contextual reinforcement learning, with a simple API, several deployment options, and source code (§5). We describe three live production deployments and two ongoing applications (§6).
- We evaluate the performance of the service (in production when possible) using systems as well as learning criteria. We also evaluate its effectiveness at reducing technical debt (§7).

2 Motivation

Our motivation to create the Decision Service came from failures and hidden costs both we and others have encountered² while deploying contextual learning solutions for production applications [4, 33, 34, 12]. Each application required substantial code and infrastructure, and was subtle enough that correct implementation eluded most developers.

Below are the sources of failure that we target, at a high level. If not addressed by design, most of these issues can cause difficult-to-debug performance degradation later in the deployment and lead to technical debts.

¹*E.g.*, 1 billion policies for the data collection cost of 21 A/B tests [38].

²Unfortunately, such failures are typically not published.

(F1) Partial feedback and bias. Due to partial feedback, a system that does not explore will reinforce its own biased decisions. Rewards that arrive after variable delays can introduce bias in favor of actions with faster positive feedback. These biases create a discrepancy between (predicted) performance at learning time and (actual) performance at decision time. *Examples:* The baseline systems in our content recommendation deployments do not explore, and are unable to evaluate actions they do not pick.³ Fast clicks on “click-bait” articles can temporarily bias the system in their favor.

(F2) Incorrect data collection. Distributed and delayed data collection increases inconsistencies between the data seen at decision time and the data used for learning. With complex pipelines, it is common to log not the action chosen by the ML algorithm, but the outcome at the end of the pipeline. The benefits of MWT are lost if the data is incorrect. *Examples:* In MSN and TrackRevenue, the context features (*e.g.*, user browsing history and click statistics) can be updated by independent processes that pre-date our deployment and are not necessarily designed with ML in mind. MSN editors periodically lock content on the site (*e.g.*, breaking news) which may override our decisions and typical systems record just the editorial override.

(F3) Changes in the environment. Real environments are *non-stationary*: the distribution of inputs to the system—*e.g.*, user requests or input features generated by upstream components—as well as the processing of outputs by downstream components, changes over time. Adapting to these changes is necessary to maintain performance. *Examples:* In our deployments, we have seen breaking news events sway MSN users’ interests en masse. Downstream business rules at MSN and Complex also change over time.

(F4) Weak monitoring and debugging. The ability to reproduce an online run offline is key to debugging a learning system, but rarely supported in full. Most systems lack the ability to ask “what if”, or *counterfactual* questions, making it hard to implement safeguards or automated responses even when real-time monitoring is available. *Examples:* Early in our MSN deployment, the cause for a buggy ML model was correctly traced only after full reproducibility was achieved. All product teams we engaged with used expensive A/B tests to monitor the performance of alternative solutions.

There is a common theme above of ensuring consistency between learning-time and decision-time performance. This is not surprising: being able to accurately predict online performance is essential to a learning system.

We discuss related work in §8, but note here that the existing ML systems fail to address most of these issues. In particular, the supervised learning systems listed in Table 1 do not support exploration, and hence do not address contextual

Category	Example ML systems
Supervised learning	Caffe, CNTK, GraphLab, LUIS, Minerva, mldb.ai, MXNet, MLlib, NEXT, Param.Server, ScikitLearn, TensorFlow, Torch
Bandit learning	Clipper, Google Analytics, LASER, VW, Yelp MOE
Cloud service	Amazon ML, AzureML, Google Cloud ML
A/B testing	Google Analytics, MixPanel, Optimizely

Table 1: A simplified categorization of ML systems (overlaps exist, *e.g.*, cloud services support supervised learning). None of the systems address our main technical debts (§2).

learning settings or the issues in (F1). Several systems support bandit learning but not *contextual* bandit learning, which is significantly more powerful. None of the systems support data collection to properly address (F2), except LUIS (for a different setting of active learning) and NEXT, which additionally handles multi-armed bandits but not general contextual bandits and neither system supports offline evaluation/monitoring highlighted in (F4). The cloud services support retraining and deploying models, but not in an online fashion. Clipper [15] incorporates feedback in real-time and thus better addresses (F3), but only explores over the predictions of existing (batch-trained) models. Several systems provide some debugging functionality, and some (*e.g.*, NEXT) support reproducible runs, addressing (F4). While A/B testing platforms can answer counterfactual questions by running a live experiment for each question, the systems in Table 1 cannot answer new counterfactual questions from already-collected data, let alone match the sample efficiency of MWT. The Decision Service fills all these gaps.

The paper addresses the motivating issues as follows. The ML methodology in §3 partially addresses (F1), (F3) and (F4), and the system abstractions in §4 aim to address all issues at once. Our evaluation (§7.2) investigates the specific examples listed above (among others).

3 Machine learning methodology

From the machine learning perspective, we implement a capability we call *multiworld testing* (MWT), which can test and optimize over K policies using data and computation that scales as $\log K$, without any prior knowledge of the policies. We show that MWT addresses (F1) and helps address (F4).

Our methodology synthesizes ideas from contextual bandits (*e.g.*, [3, 31]) and policy evaluation (*e.g.*, [33, 34]). The methodology is modular: exploration and logging support offline evaluation and learning over arbitrary policy sets. This modularity maps to modularity in our system design.

Contextual decisions. Consider an application APP that interacts with its environment, such as a news website with users or a cloud controller with machines. Each interaction follows the same broadly applicable protocol:

1. A *context* x arrives and is observed by APP.
2. APP chooses *action* $a \in A$ to take (A may depend on x).
3. A *reward* r for a is observed by APP.

Table 2 shows examples from our deployments. Contexts and actions are usually represented as feature vectors. APP

³An illustration of such problems can be found in the blog post: <http://hunch.net/?s=randomized+experimentation>

	News website (News)	Tech support assistance (TechSupp)	Cloud controller (Cloud)
Decision to optimize	article to display on top	response to query	wait time before reboot unresponsive machine
Context	location, browsing history,...	previous dialog elements	machine hardware/OS, failure history,...
Feasible actions	available news articles	pointers to potential solutions	minutes in {1,2,...,9}
Reward	click/no-click	(negative) human intervention request	(negative) total downtime

Table 2: Example applications of the Decision Service. Each is representative of a real deployment discussed in §6.

chooses actions by applying a *policy* π that takes a context as input and returns an action. The goal is to find a policy that maximizes average reward over a sequence of interactions. We assume for now that APP faces a stationary environment.⁴ **Exploration and logging.** An *exploration policy* is used to randomize each choice of action. The randomization need not be uniform and for best performance should not be [3, 6]. A simple policy is *EpsilonGreedy*: with probability ϵ_0 it chooses an action uniformly, and uses a *default policy* π_0 otherwise. π_0 might be the baseline deployed in production or the current best guess for an optimal policy. ϵ_0 controls an *explore-exploit tradeoff*: π_0 guarantees a minimum performance while randomization explores for better alternatives. The parameters ϵ_0 and π_0 can be changed over time. The contextual bandit literature provides several exploration policies including near-optimal schemes.

Each interaction is logged as a tuple (x, a, r, p) , where p is the exploration policy’s probability of choosing a given x . These datapoints are called *exploration data*. Recording p enables unbiased policy learning, which addresses (F1).

Policy learning. Given N exploration datapoints, we can *evaluate* any policy π (i.e., estimate its average reward) regardless of how the data was collected. The simplest approach is to use *inverse propensity scoring* (i_{ps}):

$$\text{i_{ps}}(\pi) = \frac{1}{N} \sum_{t=1}^N \mathbb{1}\{\pi(x_t) = a_t\} r_t/p_t, \quad (1)$$

where the indicator is 1 when π ’s action matches the exploration data and 0 otherwise. This estimator has three important properties. First, it is *data-efficient*: each interaction can be used to evaluate any π that has a matching action, *regardless of the policy collecting the data*. In contrast, A/B testing only uses data collected using π to evaluate π . Second, the importance weighting by p_t makes it statistically *unbiased*: it converges to the true reward as $N \rightarrow \infty$. Third, the estimator can be *computed incrementally* as new data arrives.

Thus, using a fixed exploration dataset, we can compute accurate counterfactual estimates of how arbitrary policies *would have performed* without actually running them, in real-time. We use this in our system design to enable advanced monitoring and safeguards, addressing (F4). In contrast, A/B testing would have to run a live experiment to test each policy.

The ability to reuse data is what makes this approach exponentially more efficient than A/B testing, in a manner we can quantify. Suppose we wish to evaluate K different policies. Let ϵ be the minimum probability given to each action

⁴Stationary here means the context and the reward given the context-action pair are drawn independently from fixed distributions.

in the exploration data (for *EpsilonGreedy*, $\epsilon = \epsilon_0/|A|$), and assume all rewards lie in $[0, 1]$. Then, with probability $1 - \delta$ the i_{ps} estimator yields a confidence interval of size $\sqrt{\frac{C}{\epsilon N} \log \frac{K}{\delta}}$ for the values of all K policies simultaneously, where C is a small constant. Crucially, the error scales $O(\log K)$. In contrast, with A/B testing the error could be as large as $C\sqrt{\frac{K}{N} \log \frac{K}{\delta}}$, which is *exponentially worse*. This also underscores the need to explore over all relevant actions: if $\epsilon = 0$, we cannot correctly evaluate arbitrary policies.

Policy evaluation allows us to search a policy class Π to find the best policy, with accuracy similar to the above bounds (replace K with $|\Pi|$). This is called *policy training*. Typically Π is defined by a tunable template, such as linear vectors, decision trees, or neural nets. Note that we need not test every policy in Π ; instead, we can use a reduction to *cost-sensitive classification* [19], for which many practical algorithms exist.

Non-stationarity. The estimates of policy performance are only predictive of future performance if the environment is stationary, but in practice applications exhibit only periods of (near-)stationarity. To cope with a changing environment, our system design implements a continuous loop in the vein of Fig. 1. This has two implications for policy training. First, an online learning algorithm is (almost) necessary, so that new datapoints can be incorporated quickly without restarting the training. Second, since most online learning algorithms gradually become less sensitive to new data (via a shrinking learning rate in the optimization), we periodically reset the learning rate (e.g., for MSN we reset each day) or use a constant rate throughout (e.g., for Complex), thereby partially addressing (F3).

Problem framing. “Framing” the problem—defining the context/features, action set, and reward metric—is often non-trivial and is a common difficulty in many applications of machine learning. The Decision Service does not directly address problem framing, but eases the task in two ways: auto-generating features for content recommendation applications (§4.3.3), and allowing testing of different framings (including changes to the reward metric, with proper logging) without collecting new data (see “Offline Learner”, §4.2).

4 System Design

While our ML methodology breaks the biased feedback loop in (F1) and enables arbitrary policy evaluation, it takes systems support to fully address changing environment, data collection issues, and delay-related bias. For example, the ML methodology simply assumes (x, a, r, p) data is cor-

rectly provided as input, but this is nontrivial in practice.

Our system design fills these gaps to operationalize our ML methodology. We list our design goals (§4.1), define a set of abstractions and an architecture that implements them (§4.2), and highlight key techniques to meet our goals (§4.3).

4.1 Design goals

Our goals are addressing the twin concerns of keeping a low technical debt and having a performant system.

Failures / technical debt. So far we have discussed these issues at a high-level. We now create concrete design goals motivated by specific examples.

Logging at the point of decision: MWT relies crucially on accurate logging of the (x, a, r, p) tuples. The system must record (x, a, p) at the time of decision and match the appropriate reward r to it. There are many ways we have seen this go wrong in practice. For example, features may be stored as references to database entries that are updated by a separate process. Consequently, the feature values available at learning time might differ from those at decision time, whether due to updates/additions/removals or access failures. When optimizing an intermediate decision in a complex system, the action chosen initially might be overridden by downstream business logic (like editorial locking), and this is the action that gets logged. In this case, the probabilities p correspond to the choice of a , not the recorded action, and are thus incorrect. Sometimes, probabilities are stashed as part of the context and included as an action feature by accident.

Experimental unit for joining: Different rate of reward arrivals on different actions can lead to biases as discussed in §2. This can be avoided by waiting a pre-set duration for reward arrival before the joining is done, and assigning a default reward (say 0) if no feedback arrives in that duration.

Continuous learning: A system that does not continuously learn (or fails to reset the learning rate as specified by our methodology) does not adapt well in non-stationary environments where a low-latency learning loop is required.

Reproducibility: Interactive systems which span many components are challenging to debug. Events may be delayed, reordered, or dropped and affect the system in complex ways, making it difficult to reproduce a bug. The difficulty is magnified when the system is continuously learning, because the system is no longer stationary, and it is difficult to disentangle issues in the learning algorithms from systems issues. The ability to fully reproduce an online run offline is essential to effectively debug learning.

Systems goals. In addition to supporting the technical debt minimization goals above, the system needs to meet the performance and functionality demands of our customers.

For interactive high-value applications, serving latency tends to be directly linked to user experience and revenue [42]. To optimize these reward metrics, the system must provide decisions in $\sim 10\text{ms}$ or less to keep application response times under $\sim 100\text{ms}$. Some applications need to

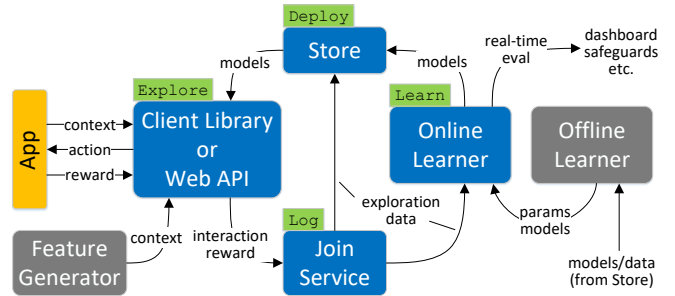


Figure 2: Decision Service architecture

quickly incorporate new data into the machine-learned policy (e.g., for breaking news stories), so the system needs a learning loop that can update the policy every few minutes.

As with any service, the system should be scalable and fault-tolerant. Two issues are unique in our setting. First, delayed rewards increase the active state of the system where the volume of state depends on the size of the outcome observation, the interaction arrival rate, and the typical delay until the reward is observed. Second, ensuring reproducibility requires care when handling reordered events from scale-out components or failures. The system should not lose data that was trained on and should recover the previously learned policies and other valuable state of the learning algorithm.

The system should provide flexible deployment options and be easy to use and customize. A modular design with well-defined interfaces admits multiple implementations, allowing us to adapt and evolve the system to applications’ needs. The system should expose the power of our ML methodology to support offline experimentation, i.e., using exploration data (instead of live experiments) to tune parameters, try other exploration/learning algorithms, etc.. To reduce setup complexity, sensible defaults should be provided for all components.

4.2 Abstractions and architecture

Our system is designed to match the modularity of the ML methodology including the exploration and logging components and the policy evaluation and training steps. Doing so, we define an abstraction for each step of the loop in Fig. 1:

- **Explore:** This component interfaces with the APP. It takes as input context features x and an event key k from APP, and outputs an action a . Separately, a keyed tuple $\langle k, (x, a, p) \rangle$ is sent to the Log component, where p is the probability of the chosen action according to the exploration policy. Later, a reward r and key k are input from APP, triggering a transmission of $\langle k, (r) \rangle$ to the Log component.
- **Log:** This component generates exploration data by joining each (x, a, p) tuple with its reward r . It defines a configurable parameter, *experimental unit*, that specifies how long to wait for a reward to arrive. Thus, the component takes a stream of keyed observations $(k, o)^*$ and emits a stream of joined observations $(o_1, o_2, \dots)^*$, where observations are

joined if they share the same key k and appear within time $[t, t + exp_unit]$, where t is the time k was first observed.

- **Learn:** This component performs policy training *online*. It takes a stream $(x, a, r, p)^*$ of exploration datapoints from the Log component and outputs a unique model $(id, model)$. A new model can be output after each datapoint.
- **Deploy:** This component stores models via a get/put interface. The Learn component puts $(id, model)$ into the store and the Explore component gets a *model* using its *id* (or the most recent model if no *id* is specified).

The abstractions by themselves address several of our technical debt reduction goals. Fig. 2 shows the architecture of the Decision Service and the component that implements each abstraction. Each component is a scalable service.

The *Client Library* implements the Explore abstraction. It implements various exploration policies from the contextual bandit literature, addressing (F1), and logs the correct data to the Log component, addressing (F2). In particular, this enforces logging at the decision point (§4.1). The Client Library downloads models from the Deploy component at a configurable rate; before the first model is available, a user-provided default policy/action may be used. The Client Library can be linked to directly for low latency decisions, or accessed via a portable, scalable web API (see Fig. 3).

The *Join Service* implements the Log abstraction by joining rewards to decisions, producing correct exploration data to address (F2). It also enforces a uniform delay (*exp_unit*) before releasing data to the Learn component to avoid delay-related biases, addressing (F1) (as discussed in §4.1 above). The exploration data is also copied to the Store for offline experimentation (discussed below).

The *Online Learner* implements the Learn abstraction using any ML software that provides online learning from contextual bandit exploration data. It incorporates data continuously and checkpoints models to the Deploy component at a configurable rate, addressing our goal of continuous learning and handling (F3). The Client Library must use compatible ML software to invoke these models for prediction. The Online Learner uses Eq. 1 to evaluate arbitrary policies in real-time, realizing MWT capability. These statistics enable advanced monitoring and safeguards (§4.3.2), addressing (F4).

The *Store* implements the Deploy abstraction to provide model and data storage. This data is also used by the *Offline Learner* for offline experimentation, such as tuning hyperparameters, evaluating other learning algorithms or policy classes, changing the reward metric, etc. Our ML methodology guarantees these experiments are counterfactually accurate [17, 19, 33, 34]. Improvements from offline experimentation can be integrated into the online loop by simply restarting the Online Learner with the new policy or settings. This addresses some usability and customizability goals.

The *Feature Generator* eases usability further by auto-generating features for certain types of content (§4.3.3).

4.3 Key techniques

We highlight several techniques that build on our abstractions to further address our design goals.

4.3.1 Full reproducibility

The Decision Service supports full offline reproducibility of online runs, despite the presence of randomized exploration and continuously updated policies. This addresses (F4).

The Client Library implements randomization using a two-layer design: various exploration policies sitting in the lower layer take as input a context and output a distribution over actions; then the top layer samples randomly from this distribution. Thus all randomization logic exists in one place. Randomization occurs by seeding a pseudorandom number generator (PRG) with the key k and an application ID; each PRG is used for exactly one interaction. Including the application ID in the seed ensures that the randomization from multiple uses of the Decision Service are not correlated.

Policies trained by the Online Learner are stored with a unique ID in the Store. The Client Library records the ID of the policy used for each decision. This, combined with our randomization scheme, ensures decisions are reproducible.

The Online Learner may encounter reordered events from the Join Service due to scale out, so it records the order in which interactions are processed with the stored policies. Other learning events, such as periodic resets of the learning rate to favor recent events (see §3), are also recorded. Together, this ensures each trained policy is reproducible.

We have not yet evaluated full reproducibility when scaling the Online Learner as not even our largest deployments have required it. In principle this can be achieved by recording the parallel learning configuration.

4.3.2 Real-time safeguards

Each of our deployments has stressed the need for real-time monitoring and safeguards against the unexpected. The Decision Service supports this, addressing (F4).

The incremental nature of ips (1) ensures we can do policy evaluation in real-time. The system supports this ability in the Online Learner, which has a mechanism to specify candidate evaluation policies, and the resulting estimates are displayed through a dashboard (Fig. 2). This provides an accurate estimate of a policy’s performance before it is deployed, and also allows real-time comparisons with default or “safe” policies. This monitoring is novel and powerful because it is counterfactually unbiased, enabling automated responses with confidence. For example, one could roll back a policy if it deviates too much from how the previous best policy would have performed, or trigger alerts when the policy starts losing out to some reasonable baseline. More advanced controls are possible, depending on the application.

4.3.3 Auto-generated features

Through working with different customers, we observed that a large class of applications—online content

recommendation—benefited from the same kinds of featurization. The Feature Generator makes it very easy for such applications to use the Decision Service (see Fig. 3). Any publicly accessible content (*e.g.*, articles, videos, etc.) can be passed to the Feature Generator ahead of time to be automatically featurized and cached—*e.g.*, text sentiment analysis, video content indexing, etc.. These features are inserted on-the-fly by the Client Library before making a prediction involving this content.

Our implementation of the Feature Generator (§5) does not directly reduce the technical debt associated with input dependencies [43], because it relies on independent cloud services to featurize the content. However, it shifts and consolidates the technical debt from multiple customers to a single service (ours), which reduces costs in the long term.

4.3.4 Low-latency learning

The Decision Service enables efficient, low-latency learning across its components, addressing a key systems goal.

The Client Library supports multi-threaded predictions over a locally-stored model, which is critical for front-end applications that serve many concurrent requests (*e.g.*, MSN). It also avoids sending repeated fragments of a context (*e.g.*, features of an article that appears in many decisions) using a simple encoding scheme: the fragments are sent explicitly only periodically, and are replaced with references otherwise. Our implementation has further optimizations (§5).

The Join Service adds no fundamental delay on top of the experimental unit, but due to the above encoding scheme, it may reorder events in a way that forces the Online Learner to buffer references to features that have not yet been seen. The Online Learner and Client Library use ML software supporting millisecond training and sub-millisecond prediction.

All components support batching to improve throughput. Without batching, the end-to-end latency of submitting an event and receiving a policy that trained on it is ~ 7 sec (§7).

5 Implementation and Deployment Options

Our implementation of the Decision Service is on Microsoft Azure [36], but our design is cloud-agnostic and could be migrated to another provider. We describe noteworthy aspects of our implementation and the deployment options. Our self-hosted service is open source [16].

The Decision Service presents a simple API, shown in Fig. 3, implementing the definition of a contextual decision from §3. The API is implemented by the Client Library or a web service that proxies requests to it; the latter is built on Azure Web Apps [10]. The Client Library links to Vowpal Wabbit (VW) [52] to predict locally on models trained by the Online Learner, which also uses VW. It is implemented in C#, C++, and Java. The C# library is about 5K lines of code: 1.5K for exploration algorithms ranging from `EpsilonGreedy` to the theoretically-optimal `Cover` [3], and 1.5K to handle batched uploads to the Join Service. We have

since migrated the code for exploration algorithms into VW (~ 1 K lines), ensuring consistent parameters and logic between training and prediction.

The Join Service is implemented using Azure Stream Analytics (ASA) [9], which allows us to specify a delayed streaming join using 45 lines of query language.

The Online Learner is implemented in 3.2K lines of C# as an Azure worker role. It uses VW for online policy training and evaluation; we modified VW to support real-time evaluation of arbitrary policies (*i.e.*, MWT capability). VW reduces policy training to cost-sensitive classification, for which many base algorithms are supported. We currently train a linear model (a vector of weights), though decision trees, neural networks, and other representations are available (*e.g.*, we are integrating CNTK [14] into VW as a base learning algorithm).

The Store is implemented using Azure Blob storage. The Offline Learner schedules jobs for policy training/evaluation on stored data using a distributed analytics platform (we are currently migrating to Azure Data Lake/U-SQL [7]).

The Feature Generator may provide additional context for predictions. It uses Microsoft Cognitive Services to featurize content from a URL and caches the results, such as Text Analytics, Computer Vision, and Video Breakdown.

Communication between the Client Library, Join Service, and Online Learner (shown as arrows in Fig. 2) happens via JSON messages sent on queues called Event Hubs.

Scalability and faults. Except for the Online Learner, all components can scale out and tolerate faults using the Azure services they are built on. The Event Hubs connecting the components are also fault tolerant: they can replay data up to 7 days in the past, which eases recovery. The Online Learner recovers from failures more coarsely: it loads the last checkpointed model from the Store and replays the Event Hub data from that point onwards.

The Online Learner can scale out using VW’s `AllReduce` communication primitive [1], which enables parallel learning across VW instances.

Optimizations for speed. To reduce the memory overhead of parallel predictions in the Client Library, we modified VW slightly to support sharing model state across VW instances. Also, since the same context class (*e.g.*, `UserContext` in Fig. 3) is processed repeatedly, we construct and reuse an abstract syntax tree to speed up serialization of the class during uploads to the Join Service. On the other side, the Online Learner supports multi-threaded deserialization of joined data to keep up with VW’s training (which is very fast).

When using the web API, we leverage Azure FrontDoor’s (AFD) global edge network, which reduces client latency by maintaining persistent connections. Additionally, the API allows contextual information to be embedded—either manually or using the Feature Generator—in a single background HTTP call (Fig. 3). When called from a browser, this allows the decision to manifest before the page finishes loading.

Sample interaction between APP and web API to choose content:

1. APP sends request to Decision Service: `https://ds.microsoft.com/api/decision/func/APP/a1/a2/.../aN.js`
User/content features can be embedded manually (e.g., `"/APP!location=NY.../a1;trending=3.2"`) or using Feature Generator.
2. Decision Service responds with action: `Action:a2,EventId:X`
This is passed via JSONP to the `func` callback to render `a2`.
3. APP reports a reward (e.g., `click=1`) using provided event ID: `https://ds.microsoft.com/api/reward?reward=1&eventId=X` Or, relies on a tracking pixel embedded into each content's page.

Sample interaction between APP and Client Library choose content:

```
var serviceConfig = new
  DecisionServiceConfiguration(settingsBlobUri: "<from
  deployment page>");
var service =
  DecisionService.Create<MyContext>(serviceConfig);
int action = service.ChooseAction(uniqueKey, myContext);
service.ReportReward(reward, uniqueKey);
```

The `ChooseAction` and `ReportReward` calls correspond to steps 1 and 3 in the web API. `MyContext` is any class that has been annotated with JSON properties to identify features, non-features, etc. [16]).

Figure 3: The Decision Service API, accessible via a cross-platform web API (left) or a low-latency Client Library (right).

Deployment options. The Decision Service supports push-button deployment (~6 min.) after a simple registration [39]. Each option below has been used for a real deployment.

- **Self-hosted:** This deploys a Decision Service loop in *your* Azure account. This is ideal for customers like MSN who wish to control the deployment or limit exposure to data.
- **Hosted:** This deploys a Decision Service loop in *our* Azure account, allowing us to share resources across tenants. We omit discussion of our multi-tenant design, but the key elements are: distributing model state across the web API front-end servers and using AFD rules to direct each application's requests to the right server, creating a multi-tenant Join Service, and consolidating Online Learners onto the same physical machines. The hosted service is being used by Complex.
- **Local mode:** This deploys a Decision Service loop locally in your machine using extended functionality in the Client Library. The Join Service is replaced by an in-memory buffered join, and models are trained and invoked for prediction by local VW instances. Local mode is particularly useful for testing in simulated environments (e.g., network/cloud simulators). It is being used by Azure Compute.

6 Deployments

We describe five applications of the Decision Service across diverse domains, as represented in Table 2. We focus on problem framing and deployment characteristics here; evaluations and lessons learned are discussed in §7.

MSN. The Decision Service personalizes news stories displayed on the MSN homepage shown in Fig. 4. The deployment is now the production default, serving 10s of millions of users who issue thousands of requests per second.

MSN's problem is essentially the News problem from Table 2: a user requests the homepage and MSN's front-end servers must decide how to order the articles on the page. If a user is logged in, there is context: demographics (e.g., age, location) and the topics of news stories they have clicked on in the past; otherwise only location is available. The action choices are the current set of news articles selected and ranked by editors (tens of articles, typically). Each article has features that describe its topic. The reward signal is clicks.

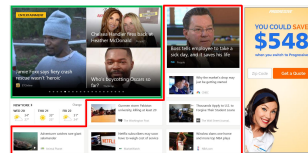


Figure 4: MSN homepage. The Slate (boxed red) and Panel (boxed green) use the Decision Service in production as of early 2016. Ads to the right (not optimized) also saw a lift.

The Decision Service optimizes the click-through rate (CTR) on the article in the most prominent slot of a segment, unless it is locked by the editors, and uses the resulting policy to pick articles for all slots.⁵ The default exploration policy `EpsilonGreedy` is used with $\epsilon = 33\%$. The experimental unit is set to 10 minutes, and a new model is deployed to the Client Library every 5 minutes. These settings are driven by the data rate and how non-stationary the news application is.

Before using the Decision Service in production, MSN used standard A/B testing to compare its performance with the editorial ordering, the default policy which beat previous attempts to use machine learning. Experiments on the Slate segment of the page showed a >25% CTR improvement over a two-week period; the Panel showed a 5.3% improvement, which is even more significant as it receives 10x more traffic than the Slate. These gains were achieved while maintaining or improving long-term engagement metrics such as sessions per unique user and average session length, showing that the CTR metric is aligned with longer-term goals in this case.

The success of the experiments led MSN to make it the default in production for all logged-in users in early 2016. The MSN team has since deployed the Decision Service a half dozen ways. MSN's deployments are self-hosted and they have customized some of our components. For example, they implemented a multi-tenant Join Service using Redis Cluster [41] that is shared across applications (e.g., different page segments). Our modular architecture allowed this customization, while ensuring correct semantics were preserved. MSN uses the C# Client Library in their front-end servers for low-latency, but clicks are reported directly from user web browsers via the web API. A single Online Learner is used for each application to isolate the trained models.

⁵There are approaches for optimizing all slots simultaneously [30].

Complex. Complex wanted to use the Decision Service in two ways: to recommend videos to associate with news articles and to recommend top news articles. These applications are similar to MSN but with several significant variations:

1. The amount of traffic is much lower, $\sim 35\text{K}$ events/day.
2. Features from Complex are much sparser, but high-dimensional, derived from keywords or other metadata.
3. Complex uses a content distribution network implying that all significant computation must be offloaded.

The first two variations make application of the Decision Service much more statistically sensitive than with MSN. The third variation forces the use of a cloud API (Fig. 3, left). We leveraged our feature generators to enhance Complex’s features with more informative ones. The Decision Service provide a lift of $>30\%$ over editorial.

TrackRevenue. TrackRevenue optimizes landing pages for in-app advertisements so as to maximize revenue with event rates of around 200k/day. This varies from the previous two content recommendation applications because the revenue signal is scalar rather than binary and much sparser. Request rates also vary enormously depending on the campaign. The baseline here is an ϵ -greedy bandit algorithm, that does not incorporate context. This baseline is structurally similar to our approach in this deployment besides the use of context, so the 18% lift provides a measure of the value of context.

Toronto. Toronto wants to reduce technical support load by automatically providing answers to technical support questions. There is a large negative reward associated with customers requesting human intervention. This problem significantly differs from the previous in the context (technical support queries), the nature of the actions (pointers to solutions), a low event rate, and privacy concerns. Fortunately, Toronto has existing systems and expertise winnowing down the set of actions to a reasonable set and good actions are stable over time allowing us to handle the low data rate. Thee privacy concerns are addressed by simply deploying the service into their own account. Results are too preliminary to report.

Azure Compute. Azure Compute wants to contextually optimize dealing with nonresponsive virtual machines in a cloud service. In particular, if you migrate too soon the cost of migration may exceed the cost of waiting for the VM to become responsive. Alternatively, waiting for a VM that never becomes responsive is a pure waste of time. A very large number of sparse features are available related to the machine, the VM, and the process on the VM.

Unlike previous applications a large effectively supervised dataset exists because the current baseline policy in their system is to simply wait for 10 minutes, making it possible to simulate the effect of any action involving migration before 10 minutes. Taking such a dataset collected with this baseline, we estimated a 19% reduction in wasted time, significant enough to proceed to production. Note that once the system is deployed, we can no longer evaluate all possible

options for free, as the system may choose to immediately migrate, and the MWT capability of the Decision Service is critical to computing real-world performance estimates.

7 Evaluation and Lessons

We evaluate the Decision Service based on how well it meets our systems and technical debt minimization goals (§2, §4.1). We use both live deployments and the exploration data collected from them.

7.1 System evaluation

From a systems perspective, we answer the following:

1. What is the latency and overhead of making decisions? How quickly is data incorporated into trained policies?
2. Can the system scale to high event rates?
3. Is there adequate support for offline experimentation?

Experimental methodology. We deployed a self-hosted Decision Service loop and drove traffic to it using the sample code distributed with the Client Library [16], running on several large A4 instances (8 cores, 14GB memory, 1Gbps NIC). Except the Online Learner, all components including the Event Hubs connecting them can be scaled by configuring the underlying Azure service. The Online Learner is a stand-alone worker running on a D5 instance (16 cores, 56GB memory).

Some of our experiments use exploration data collected from the MSN and Complex production deployments during April 2016 and April 2017, respectively. This data contains the real (x, a, p, r) tuples generated on those days with x consisting of $\sim 1\text{K}$ features per action.

7.1.1 Latency of decisions and learning

We are interested in both *decision latency* and *learning latency*. Decision latency is the time to make a decision in the Client Library (*i.e.*, the `ChooseAction` call). Learning latency is the time from when an interaction is complete (*i.e.*, `ReportReward` has been called) to when it affects a deployed policy in the Client Library.

We measured the decision latency by training a policy on one hour of MSN data, deploying this policy in the Client Library, and then repeatedly calling `ChooseAction`. The average latency is 0.2ms, well within the needs of our customers. Latency is not the only metric that matters, however. For example, MSN’s front-end servers are CPU limited, so CPU/req is monitored very carefully. The Decision Service increased CPU/req by 4.9%, which was deemed acceptable.

In our implementation, we described an optimization that leverages Azure’s edge network and JSONP calls to our web API to make fast decisions from a browser. This is enabled in the hosted Decision Service used by Complex. Using page load metrics across 20 days, we measured the difference between page load time and decision time (steps 1 and 2 in Fig. 3) for Complex. Fig. 5 shows the results: except for 5.9% of requests, all decisions complete before the page fin-

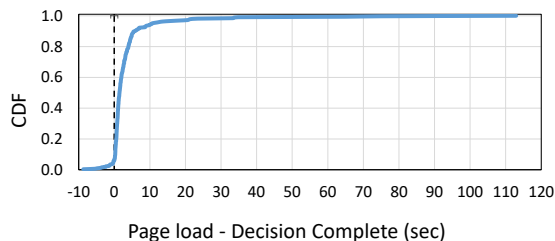


Figure 5: Decisions complete before page loads in Complex.

ishes loading, and less than 1% take >4 more seconds (the max is 8.7sec).

To measure learning latency, we first removed any configurable sources of delay, such as batching and caching (which could cause downstream reordering). We configured the Client Library to poll for new models every 100ms, we set the Join Service experimental unit to 1 second, and we configured the Online Learner to publish an updated policy after each event. We then replayed one event from the MSN data and waited for a policy to appear. The average learning latency is 7.3 sec.

7.1.2 High event rates

Most of our components are built on scalable services managed by Azure. The Feature Generator relies on Microsoft Cognitive Services APIs, which are also scalable (though outside our control). Thus we focus on the Online Learner.

The data rates seen in all our production deployments have been adequately handled by learning on a single core. To saturate the Online Learner, we preloaded MSN data into the Join Service’s output Event Hub and processed it at full speed. MSN uses the Client Library’s encoding scheme (§4.3.4) to reduce data size. The throughput achieved by the learner was stable at 2000 events/sec, implying 100 million events/day applications are viable. Buffering for events whose encoded features had not arrived was minimal, remaining at 0 most of the time with occasional spikes up to 2250 events.

We observed two different bottlenecks in the Online Learner. When using encoded contexts (*e.g.*, MSN), the smaller data makes policy training (VW) the bottleneck. For example, increasing the number of articles (actions) per decision in MSN from 10, 20, 30—which slows down training without increasing data size too much—led to a throughput decline of 2778, 1666, 1149 events/sec. When contexts are not encoded (*e.g.*, Complex), the bottleneck is the deserialization of data received from the Join Service. We verified this by changing the VW parameters to slow down training without seeing any decline in throughput, at event/sec.

Although we have been able to scale the deserialization overhead near-linearly using multiple Online Learners, we have not yet incorporated VW’s parallel learning so cannot comment on model performance. We plan to address this when the need arises and do not expect new insights over prior work [1].

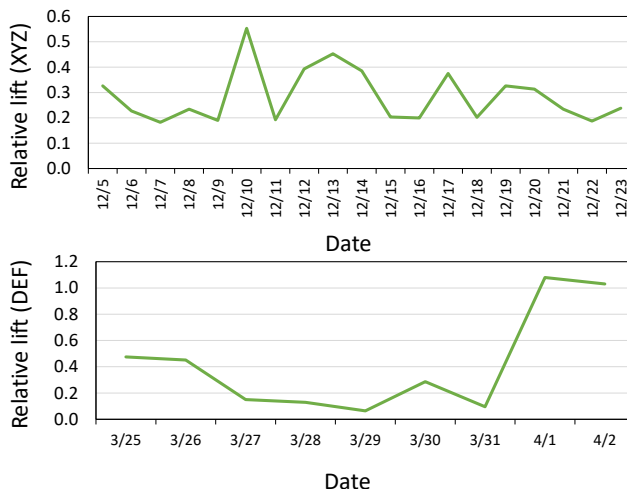


Figure 6: Daily CTR lift for MSN flight (Dec. 2015) and Complex flight (Mar. 2016).

7.1.3 Offline experimentation

The Offline Learner has enabled us to run hundreds of experiments on exploration data, across our deployments, to tune parameters, try new learning algorithms, etc.. For example, Complex began with a simple EpsilonGreedy policy to collect initial data. Using this data, we ran parallel offline experiments and discovered higher CTRs when using an ensemble exploration algorithm called Bag [3], a different learning rate, and certain feature interactions and omitted features. This was determined without any additional live experiment, and our ML methodology guaranteed that the results were counterfactually accurate.

7.2 Technical debt evaluation

From the perspective of technical debt, we ask how well the Decision Service addresses problems F1-F4. As mentioned in §2, a common theme across the problems is ensuring that performance predicted at learning time matches observed performance at decision time. Thus, most of our experiments measure the discrepancy between the two caused by failures.

Experimental methodology. To simulate failures on exploration data, we take a day’s data and randomly allocate 80% of it for training and 20% for testing while preserving time order, approximating the normal train/test split methodology of supervised learning. We modify the training data with the failure and train a policy on both the faulty and correct data; this yields an estimate of performance via the policy evaluation capability of the Online Learner (based on Eq. 1). We then evaluate both policies on the test data and obtain another estimate of performance. We report the discrepancy between the training and testing estimates, which for the correct data is always within 5%.

7.2.1 (F1) Partial feedback and bias

By using contextual bandit exploration, the Decision Service avoids the bias inherent in approaches like supervised learning. This is evidenced by the superior performance it

Failure	Train/test perf. discrepancy
Reward delay bias	1.3x
Incorrect probability	3.0x
Decision as feature	8.7x
Modified feature	1.2x
Deleted feature	2.4x

Table 3: Various data collection errors (§7.2.1,§7.2.2).

achieves compared to the editorial ranking in both MSN and Complex, shown in Fig. 6. Not only are the per-day improvements high (between 18% and 55% for MSN, between 7% and 112% for Complex), but the improvement is maintained over time (see F3 below).

Another form of bias occurs if rewards for some actions are delayed more than others. We simulated this effect in the MSN data by moving the click events of certain popular articles before those of the rest, as if the Join Service had immediately emitted them instead of enforcing a uniform experimental unit. As the Online Learner incorporates new data, the trained policy may forget the importance of these articles. Indeed, shifting the clicks of the top 2 articles yields a training performance that is 1.3x higher than test performance (Table 3). If a constant learning rate is used (as in Complex), this increases to 1.7x.

7.2.2 (F2) Incorrect data collection

We evaluated four different data collection failures using MSN data and our described methodology, all of which we have experienced in practice. Table 3 summarizes the results.

- *Incorrect probability.* A common error is when editors or business logic override the chosen action and record the override making the recorded probability incorrect. We simulated this by overriding 10% of the training data actions.
- *Decision as feature.* Another common error is when the identity or probability of the chosen action is used as a feature for downstream learning. We simulated this by adding a new feature to the training data that is 1 for the chosen action and 0 for all other actions. Note that this feature would not appear in test data because it is only available after decision.
- *Modified feature.* Features such as the user browsing history in MSN are often modified by separate processes, and hence may change between decision and learning time. We simulated this by replacing the browsing history in 20% of the training data with default values, *e.g.*, as if they were still being computed for a new user.
- *Deleted feature.* A deleted feature (*e.g.*, due to a database failure) may be present at learning but not at decision time. We simulated this by removing user demographic features.

In all of the above cases, the estimate of the policy’s performance during training deviated from its performance during testing, by a factor ranging from 1.2x to 8.7x. This undermines the entire value of a counterfactually accurate system. By correctly logging at the point of decision, the Decision Service avoids these failures.

7.2.3 (F3) Changes in the environment

We observed significant non-stationarity in the MSN and Complex data, which is likely due to the continuous arrival of new content and user interests swaying to breaking news or events. The Decision Service is able to sustain its improvement in Fig. 6 by continuously learning *and* periodically adjusting its learning rate to favor recent data, as prescribed by our ML methodology. Without this, performance would degrade over time.

Periodic resets of the learning rate—*e.g.*, in all MSN deployments the learning rate is reset each day—may affect the Online Learner’s ability to converge to a good policy quickly. To investigate this, we played a full day of MSN data and tracked the performance of the trained policy relative to the editorial policy (Editorial 1 in Fig. 7). Both the trained and editorial policies exhibit high variance initially, but the estimates become statistically significant with more data. The trained policy starts outperforming editorial after just 65K interactions—for a request rate of 1000/sec, this is about 1 minute—and eventually achieves a 42% improvement by end of day.

To demonstrate the cost of not continuously training, we took three days of MSN data, trained a policy on day 1 and tested it on days 2 and 3 (without updates). The performance relative to a policy trained on the corresponding day is low:

Policy from:	Day 1	Day 2	Day 3
Same day	1.0	1.0	1.0
Day 1	1.0	0.73	0.46

In other words, day 1’s policy achieves 73% of the CTR of day 2’s policy when tested on day 2, and 46% of day 3’s policy on day 3. This suggests that the environment and articles have changed, and day 1’s policy is stale.

Some of the data collection errors in the previous section are caused by changes in how features are generated. For example, the baseline policy in TrackRevenue uses click statistics that are aggregated with a certain decay. When we tried to use VW’s built-in marginal statistics to match this feature, it took several weeks to obtain consistent results. This was partly because TrackRevenue’s method of collecting the statistics changed over time, and different details were conveyed to us at different points. Continuous learning can help cope with such uncertainties.

7.2.4 (F4) Weak monitoring and debugging

The Decision Service supports real-time evaluation of arbitrary policies, an extremely powerful capability. To demonstrate this, we took 3 simple policies which always choose the first, second and third article respectively from the editorial ranking in MSN. The first policy corresponds to the editorial baseline, while others are reasonable alternatives. Fig. 7 shows a real-time comparison of policy performance generated by a running instance of the Decision Service that we replayed MSN data through. Each datapoint was captured within an average of 10 sec from when the corresponding interaction completed. Although not in production yet,

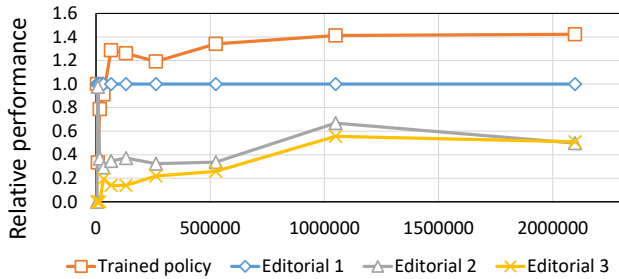


Figure 7: Real-time comparison of trained policy against editorial policies, normalized to Editorial 1’s performance.

several of our deployments have plans to install alerts based on these estimates.

The Decision Service supports full offline reproducibility of online runs by design. We cannot over-emphasize the value of this property, as it has helped us diagnose numerous failures and performance issues in the course of our deployments. Several of these failures were diagnosed by simply detecting the point at which the offline run deviated from the online run, *e.g.*, revealing that data was dropped. However, the real value came in diagnosing problems in the ML software (VW). Typically such bugs are difficult to diagnose because the ML software is treated as black-box, but with full reproducibility we can quickly rule out any of our components as being the culprit. For example, this helped us debug a poorly performing MSN model that was initialized from another model trained with an insufficient range of data. This triggered a bug in VW that inadvertently reset the min/max prediction range in a way that was incompatible with the actual rewards of $\{-1, 0\}$.

8 Related work

Here we discuss other machine learning approaches related to MWT, including other experimentation and ML systems.

8.1 Machine learning with exploration

There are hundreds of papers related to exploration which broadly fall into 3 categories. In *active learning* [11, 26, 27, 44], the algorithm helps select examples to label for a supervised learner. A maximally general setting is *reinforcement learning* [46, 47] where an algorithm repeatedly chooses among actions and optimizes *long-term reward*. A simpler setting is multi-armed bandits (MAB) where actions are chosen without contextual information [13, 22]. We build on contextual bandits with policy sets [2, 3, 6, 18, 31], as well as offline policy evaluation [17, 19, 33, 34]. The MWT capability enabled in this approach is typically absent from alternatives such as Thompson Sampling [50].

8.2 Systems for ML and experimentation

We previously discussed these systems with regards to technical debt (§2). A more general discussion follows.

A/B testing. A/B testing refers to randomized experiments with subjects randomly partitioned amongst treatments. It is routinely used in medicine and social science, and has be-

come standard in many Internet services [29, 28], as well supported by statistical theory [21]. A more advanced version, “multi-variate testing”, runs many A/B tests in parallel. Several commercialized systems provide A/B testing in web services (Google Analytics [23], Optimizely [40], Mix-Panel [37], etc.). The Decision Service instead builds on MWT, a paradigm exponentially more efficient in data usage than A/B testing.

Bandit learning systems. Several platforms support bandit learning for web services. Google Analytics [23] supports Thompson Sampling [50]. Yelp MOE [54] is an open-source software package which implements optimization over a large parameter space via sequential A/B tests. Bayesian optimization is used to compute parameters for the “next” A/B test.⁶ Clipper [15] uses bandit algorithms to adapt over supervised ML systems. However, these systems do not support *contextual* bandit learning, and they do not instrument automatic deployment of learned policies (and hence do not “close the loop” in Figure 1).

Contextual bandits deployments. There have been several applications of contextual bandit learning in web services (*e.g.*, news recommendation [4, 33, 34] and advertising [12]). However, they have all been one-offs rather than a general-purpose system like the Decision Service.

Systems for supervised learning. There are many systems designed for supervised machine learning such as CNTK [14], GraphLab [25], Parameter Server [35], MLlib [45], TensorFlow [48, 49], Torch [51], Minerva [53] and Vowpal Wabbit [52] to name a few. These principally support Machine Learning model development. A few more, such as Google Cloud ML [24], Amazon ML [5], and AzureML [8] are designed to support development and *deployment*. However, these systems do not support data gathering or exploration.

We know of two other systems that fully support data collection with exploration, model development, and deployment: LUIS [32] (based on ICE [44]), and NEXT [27]. These systems support *active learning*,⁷ and hence make exploration decisions for labeling in the back-end (unlike the Decision Service which makes decisions for customer-facing applications), and do not provide MWT capability.

9 Conclusion

We have presented the Decision Service, the first general-purpose service for contextual learning. It supports the complete data lifecycle and combines an ML methodology with careful system design to address many common failure modes. Going forward, our goal is to make the service completely parameter free. We also plan to use MWT capability to provide more sophisticated safeguards for production deployments.

⁶SigOpt.com is a commercial platform which builds on Yelp MOE.

⁷NEXT does have some support for bandit algorithms, but does not provide the MWT capability or fully general contextual bandits.

References

- [1] A. Agarwal, O. Chapelle, M. Dudík, and J. Langford. A reliable effective terascale linear learning system. *The Journal of Machine Learning Research*, 15(1):1111–1133, 2014.
- [2] A. Agarwal, M. Dudík, S. Kale, J. Langford, and R. E. Schapire. Contextual bandit learning with predictable rewards. In *15th Intl. Conf. on Artificial Intelligence and Statistics (AISTATS)*, pages 19–26, 2012.
- [3] A. Agarwal, D. Hsu, S. Kale, J. Langford, L. Li, and R. Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *31st Intl. Conf. on Machine Learning (ICML)*, 2014.
- [4] D. Agarwal, B.-C. Chen, Q. He, Z. Hua, G. Lebanon, Y. Ma, P. Shivaswamy, H.-P. Tseng, J. Yang, and L. Zhang. Personalizing linkedin feed. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1651–1660. ACM, 2015.
- [5] Amazon Machine Learning - Predictive Analytics with AWS. <https://aws.amazon.com/machine-learning/>.
- [6] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002. Preliminary version in *36th IEEE FOCS*, 1995.
- [7] Azure Data Lake Analytics. <https://azure.microsoft.com/en-us/services/data-lake-analytics/>.
- [8] Azure Machine Learning. <https://azure.microsoft.com/en-us/services/machine-learning>.
- [9] Azure Stream Analytics. <https://azure.microsoft.com/en-us/services/stream-analytics/>.
- [10] Azure Web Apps. <https://azure.microsoft.com/en-us/services/app-service/web/>.
- [11] A. Beygelzimer, J. Langford, Z. Tong, and D. J. Hsu. Agnostic active learning without constraints. In *Advances in Neural Information Processing Systems*, pages 199–207, 2010.
- [12] L. Bottou, J. Peters, J. Quinero-Candela, D. X. Charles, D. M. Chickering, E. Portugaly, D. Ray, P. Simard, and E. Snelson. Counterfactual reasoning and learning systems: The example of computational advertising. *J. of Machine Learning Research (JMLR)*, 14(1):3207–3260, 2013.
- [13] S. Bubeck and N. Cesa-Bianchi. Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Foundations and Trends in Machine Learning*, 5(1), 2012.
- [14] Computational Netwrk Toolkit. <http://www.cntk.ai/>.
- [15] D. Crankshaw, X. Wang, G. Zhou, M. J. Franklin, J. E. Gonzalez, and I. Stoica. Clipper: A low-latency online prediction serving system. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 613–627, 2017.
- [16] Decision Service source code and documentation. [http://\[anonymized\]/](http://[anonymized]/).
- [17] M. Dudík, D. Erhan, J. Langford, and L. Li. Sample-efficient nonstationary policy evaluation for contextual bandits. In *28th Conf. on Uncertainty in Artificial Intelligence (UAI)*, pages 247–254, 2012.
- [18] M. Dudik, D. Hsu, S. Kale, N. Karampatziakis, J. Langford, L. Reyzin, and T. Zhang. Efficient optimal learning for contextual bandits. In *27th Conf. on Uncertainty in Artificial Intelligence (UAI)*, 2011.
- [19] M. Dudík, J. Langford, and L. Li. Doubly robust policy evaluation and learning. In *28th Intl. Conf. on Machine Learning (ICML)*, pages 1097–1104, 2011.
- [20] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [21] A. S. Gerber and D. P. Green. *Field Experiments: Design, Analysis, and Interpretation*. W.W. Norton&Co, Inc., 2012.
- [22] J. Gittins, K. Glazebrook, and R. Weber. *Multi-Armed Bandit Allocation Indices*. John Wiley & Sons, 2011.
- [23] Google Analytics. <http://www.google.com/analytics>. See <http://services.google.com/webstioptimizer> for documentation on bandits.
- [24] Google Cloud Machine Learning. <https://cloud.google.com/ml/>.
- [25] GraphLab. <http://graphlab.org> and <https://dato.com/products/create>.
- [26] S. Hanneke. Theory of disagreement-based active learning. *Foundations and Trends® in Machine Learning*, 7(2-3):131–309, 2014.
- [27] K. G. Jamieson, L. Jain, C. Fernandez, N. J. Glattard, and R. Nowak. NEXT: A system for real-world development, evaluation, and application of active learning. In *Advances in Neural Information Processing Systems*, pages 2638–2646, 2015.
- [28] R. Kohavi and R. Longbotham. Online controlled experiments and a/b tests. In Claude Sammut and Geoff Webb, editor, *Encyclopedia of Machine Learning and Data Mining*. Springer, 2015. To appear.
- [29] R. Kohavi, R. Longbotham, D. Sommerfield, and R. M. Henne. Controlled experiments on the web: survey and practical guide. *Data Min. Knowl. Discov.*, 18(1):140–181, 2009.
- [30] A. Krishnamurthy, A. Agarwal, and M. Dudík. Efficient contextual semi-bandit learning. *arxiv.org*, abs/1502.05890, 2015.
- [31] J. Langford and T. Zhang. The Epoch-Greedy Algorithm for Contextual Multi-armed Bandits. In *21st Advances in Neural Information Processing Systems (NIPS)*, 2007.
- [32] Language Understanding Intelligent Service (LUIS). <https://www.luis.ai>.
- [33] L. Li, W. Chu, J. Langford, and R. E. Schapire. A contextual-bandit approach to personalized news article recommendation. In *19th Intl. World Wide Web Conf. (WWW)*, 2010.
- [34] L. Li, W. Chu, J. Langford, and X. Wang. Unbiased of-line evaluation of contextual-bandit-based news article recommendation algorithms. In *4th ACM Intl. Conf. on Web Search and Data Mining (WSDM)*, 2011.

- [35] M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B.-Y. Su. Scaling distributed machine learning with the parameter server. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 583–598, 2014.
- [36] Microsoft Azure. <https://azure.microsoft.com/en-us/>.
- [37] Mixpanel: Mobile Analytics. <https://mixpanel.com/>.
- [38] Multi-world testing: A system for experimentation, learning, and decision-making. White paper, [http://\[anonymized\]/](http://[anonymized]/), 2016.
- [39] Multiworld testing Decision Service. [http://\[anonymized\]/](http://[anonymized]/).
- [40] Optimizely: A/B Testing & Personalization Platform. <https://www.optimizely.com/>.
- [41] Redis Cluster. <http://redis.io/topics/cluster-spec>.
- [42] E. Schurman and J. Brutlag. The user and business impact of server delays, additional bytes, and http chunking in web search. In *Velocity*, 2009.
- [43] D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, and M. Young. Machine learning: The high-interest credit card of technical debt. In *SE4ML: Software Engineering 4 Machine Learning*, 2014.
- [44] P. Simard, D. Chickering, A. Lakshmiratan, D. Charles, L. Bottou, C. G. J. Suarez, D. Grangier, S. Amershi, J. Verwey, and J. Suh. Ice: enabling non-experts to build models interactively for large-scale lopsided problems. *arXiv preprint arXiv:1409.4814*, 2014.
- [45] SPARK MLlib. <http://spark.apache.org/mllib>.
- [46] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [47] C. Szepesvári. *Algorithms for Reinforcement Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2010.
- [48] TensorFlow. <https://www.tensorflow.org/>.
- [49] Tensorflow serving. <https://tensorflow.github.io/serving>.
- [50] W. R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- [51] Torch. <http://torch.ch/>.
- [52] Vowpal Wabbit (Fast Learning). <http://hunch.net/~vw/>.
- [53] M. Wang, T. Xiao, J. Li, J. Zhang, C. Hong, and Z. Zhang. Minerva: A scalable and highly efficient training platform for deep learning. In *NIPS Workshop, Distributed Machine Learning and Matrix Computations*, 2014.
- [54] Yelp MOE (Metrics Optimization Engine). <http://yelp.github.io/MOE>.