
Clustering Context in Off-Policy Evaluation

Daniel Guzmán-Olivares
Bulil Technologies, UAM*
daniel.guzman@buliltec.com

Philipp Schmidt
Amazon
phschmid@amazon.com

Jacek Golebiowski
distil labs*
golebiowski.j@gmail.com

Artur Bekasov
Amazon
abksv@amazon.com

Abstract

Off-policy evaluation can leverage logged data to estimate the effectiveness of new policies in e-commerce, search engines, media streaming services, or automatic diagnostic tools in healthcare. However, the performance of baseline off-policy estimators like IPS deteriorates when the logging policy significantly differs from the evaluation policy. Recent work proposes sharing information across similar actions to mitigate this problem. In this work, we propose an alternative estimator that shares information across similar contexts using clustering. We study the theoretical properties of the proposed estimator, characterizing its bias and variance under different conditions. We also compare the performance of the proposed estimator and existing approaches in various synthetic problems, as well as a real-world recommendation dataset. Our experimental results confirm that clustering contexts improves estimation accuracy, especially in deficient information settings.¹

1 INTRODUCTION

The contextual bandit process models many real-world problems across industry and research, including healthcare, finance, and recommendation systems (Bouneffouf, Rish, and Aggarwal, 2020). In this setting, an agent observes a *context*, chooses an action according to a *policy*, and observes a *reward*. *Off-policy evaluation* (OPE)

methods aim to estimate the effectiveness of a policy without empirically testing it, which can be particularly useful when A/B tests are costly, or if there is an inherent risk associated with poor policy performance, as is often the case in healthcare (Bastani and Bayati, 2019). Existing OPE methods can be broadly divided into parametric methods based on the *direct method* (DM), non-parametric methods based on *inverse propensity score* weighting (IPS, Horvitz and Thompson, 1952), and a combination of the two, such as the *doubly robust* method (DR, Dudík, Langford, and Li, 2011). When every action with non-zero probability under the evaluation policy also has a non-zero probability under the logging policy, IPS is unbiased. This condition is rarely satisfied in real-world problems, however, so IPS is typically biased in practice, especially for actions that violate the condition, or have close-to-zero probabilities in the logging policy (Sachdeva, Su, and Joachims, 2020; Dudík, Langford, and Li, 2011; Saito and Joachims, 2022).

Recently proposed *Marginalized Inverse Propensity Score* estimator (MIPS, Saito and Joachims, 2022) improves upon IPS in large action spaces by pooling information across *action embeddings*. At the same time, MIPS suffers from the same problem as IPS for contexts in which a significant proportion of actions have low probability under the logging policy. In this case, MIPS lacks information about the actions to accurately estimate the importance weights, resulting in additional bias. In our work, we hypothesize that closeness at the context level should translate into similar behaviour for actions and rewards (for example, two movies of the same franchise in a recommendation system). Based on this hypothesis, we propose an estimator that *clusters* the context space, and pools information across all the contexts within a cluster. Informally, the proposed method solves the problem of deficient action information for a particular context by leveraging the information from all other contexts within the same cluster.

We define and analyze the theoretical bandit setup with context clusters in Section 3, which leads to the formal derivation of the CHIPS estimator, for which we analyze bias and variance. In section 4, we compare the estimator’s performance to the baselines on several synthetic and real-world datasets, verifying the theoretical findings, and

*Work done while at Amazon.

¹The code for reproducing our experimental implementation is available at <https://github.com/amazon-science/opec-cluster-context>

demonstrating its effectiveness. Finally section 5 explores future lines of work and CHIPS' limitations.

2 BACKGROUND ON OFF-POLICY EVALUATION AND RELATED WORK

The off-policy evaluation problem (OPE) is usually framed inside the general contextual bandit setup. Given an agent, determined by the policy $\pi : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, the bandit's data generation process is defined as iterative logging of the agent's behavior when presented with different contexts. In each iteration, a context $x \in \mathcal{X} \subseteq \mathbb{R}^{d_x}$ is drawn i.i.d. from an unknown probability distribution $p(x)$ over the context space, an action $a \sim \pi(a|x)$ is selected from a finite action space \mathcal{A} , and a bounded reward $r \in [0, R_{\max}]$ is observed as a sample from an unknown conditional distribution $p(r|a, x)$. The off-policy evaluation problem has been extensively studied from both a theoretical (McNelis et al., 2017; Saito et al., 2021; Dumitrascu, Feng, and Engelhardt, 2018; Irpan et al., 2019; Wang, Agarwal, and Dudík, 2017) and a practical point of view given its applications in fields such as recommendation systems (Li et al., 2011; Bendada, Salha, and Bontempelli, 2020; Saito et al., 2020) or healthcare (Varatharajah and Berry, 2022).

We measure the performance of a policy π through its *value*, that we define as:

$$V(\pi) := \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r] = \mathbb{E}_{p(x)\pi(a|x)}[q(a, x)] \quad (1)$$

Here $q(a, x) = \mathbb{E}_{p(r|a,x)}[r]$ denotes the conditional expected reward given an action a and a context x .

In practice, we are interested in finding a policy maximizing the expected reward observed in the bandit process. A vital part of this process is the off-policy evaluation problem, in which we estimate the value of a policy π given a dataset $\mathcal{D} := \{(x_i, a_i, r_i)\}_{i=1}^N$ collected under a logging policy π_0 (i.e. $\mathcal{D} \sim \prod_{i=1}^N p(x)\pi_0(a|x)p(r|a, x)$). We use the mean squared error (MSE) to quantify how well the estimate $\hat{V}(\pi)$ approximates the real policy value $V(\pi)$:

$$\begin{aligned} \text{MSE}(\hat{V}) &= \mathbb{E}_{\mathcal{D}}[(V(\pi) - \hat{V}(\pi; \mathcal{D}))^2] \\ &= \text{Bias}(\hat{V}(\pi; \mathcal{D}))^2 + \mathbb{V}_{\mathcal{D}}[\hat{V}(\pi; \mathcal{D})] \end{aligned}$$

A wide variety of approaches have been proposed in the literature to estimate $V(\pi)$. From them, three can be distinguished for being commonly used as starting points for developing new estimators. The first one is the Direct Method (DM), which tries to estimate $q(a, x)$ directly from Equation (1):

$$\hat{V}_{\text{DM}}(\pi; \mathcal{D}, \hat{q}) = \frac{1}{N} \sum_{i=1}^N \sum_{a \in \mathcal{A}} \hat{q}(a, x_i)$$

The bias of DM depends on the accuracy of the $\hat{q}(a, x) \approx q(a, x)$ approximation, but the variance is usually lower

than in other approaches. Supervised learning in the DM's approach can be particularly useful when generalization of an agent's behaviour is needed due to limited information in the logging data (Sachdeva, Su, and Joachims, 2020). However, when the reward function has a high variance, or the representation capacity is limited for the context-action pairs in the evaluation policy domain, $\hat{q}(a, x)$ could fail to accurately approximate $q(a, x)$ (Farajtabar, Chow, and Ghavamzadeh, 2018; Beygelzimer and Langford, 2009; Kallus and Uehara, 2019). This problem, known as *reward misspecification*, can be quite difficult to detect in real-world examples (Farajtabar, Chow, and Ghavamzadeh, 2018; Voloshin et al., 2021), and is the reason why DM is generally regarded as a highly biased estimator.

The second base approach is Inverse Propensity Scoring (IPS, Horvitz and Thompson, 1952), which approximates the policy value by reweighting the rewards to correct the shift in action probabilities between the logging and evaluation policies:

$$\hat{V}_{\text{IPS}}(\pi; \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \frac{\pi(a_i|x_i)}{\pi_0(a_i|x_i)} r_i = \frac{1}{N} \sum_{i=1}^N w(a_i, x_i) r_i$$

As per this definition, the context-action pairs selected by π in which $\pi_0(a|x) = 0$ could be problematic, which motivates the following assumption:

Assumption 2.1. (*Common Support*) Given an evaluation policy π and a logging policy π_0 , the latest has common support for π if

$$\pi_0(a|x) > 0 \quad \forall a \in \mathcal{A}, x \in \mathcal{X} : \pi(a|x) > 0$$

The IPS estimator is unbiased under Assumption 2.1. However, even when assumption 2.1 holds, IPS can present excessive variance due to the weights $w(a_i, x_i)$ taking larger values (Dudík, Langford, and Li, 2011; Saito and Joachims, 2022). This case is especially notable when π_0 and π are significantly different or when trying to achieve universal support ($\pi_0(a|x) > 0 \forall a \in \mathcal{A}, x \in \mathcal{X}$) in large action spaces (Saito and Joachims, 2022; Peng et al., 2023; Saito et al., 2021). Controlling the scaling of the propensity scores has motivated many approaches based on IPS, using techniques such as weight clipping (Su et al., 2020; Su, Srinath, and Krishnamurthy, 2020; Swaminathan and Joachims, 2015a) and self normalization (Swaminathan and Joachims, 2015b; Kuzborskij et al., 2020). The Doubly Robust (DR) estimator combines DM and IPS, aiming to obtain a low-bias, low-variance estimate:

$$\begin{aligned} V_{\text{DR}}(\pi; \mathcal{D}, \hat{q}) &:= V_{\text{DM}}(\pi; \hat{q}) \\ &\quad + \frac{1}{N} \sum_{i=1}^N w(a_i, x_i) (r_i - \hat{q}(a_i, x_i)) \end{aligned}$$

The DR estimator has been the cornerstone of multiple approaches that modify the base estimator to address problems such as low overlap between π and π_0 (Wang, Agarwal, and Dudík, 2017; Metelli, Russo, and Restelli, 2021;

Zhan et al., 2021; Guo et al., 2024), reward misspecification (Farajtabar, Chow, and Ghavamzadeh, 2018), and limited samples in logging data (Su et al., 2020; Felicioni et al., 2022). Unfortunately, the DR estimator can still inherit the large variance problem from IPS, for example, when dealing with large action spaces (Saito, Ren, and Joachims, 2023; Saito and Joachims, 2022; Shimizu and Forastiere, 2023; Sachdeva et al., 2023; Taufiq et al., 2023). The problem of dealing with large action spaces was recently studied, resulting in the *Marginalized Inverse Propensity Scoring* (MIPS) (Saito and Joachims, 2022) estimator, in which the authors pool information between similar actions given some embedding representation $e \in \mathcal{E} \subset \mathbb{R}_e^d$ of them to address deficient actions in the logging policy. For this purpose, they introduce an IPS-based estimator marginalizing the probability over the action space:

$$\begin{aligned}\hat{V}_{\text{MIPS}}(\pi; \mathcal{D}) &:= \frac{1}{n} \sum_{i=1}^n \frac{p(e_i | x_i, \pi)}{p(e_i | x_i, \pi_0)} r_i \\ &= \frac{1}{n} \sum_{i=1}^n w(x_i, e_i) r_i.\end{aligned}\quad (2)$$

Where $p(e|x, \pi) := \sum_{a \in \mathcal{A}} p(e|x, a) \pi(a|x)$. The idea of estimating deficient items' behaviour by *closely* observed ones inspired new approaches, like partitioning the action space in clusters (Peng et al., 2023; Saito, Ren, and Joachims, 2023), or an adaptive method for ranking policies by optimizing user classification into given behavioural models and estimating independently for each group (Kiyohara et al., 2023). The MR estimator (Taufiq et al., 2023) diverged from the action space transformations and proposed marginalization over the rewards density through a regression estimate of the importance weights:

$$\hat{V}_{\text{MR}}(\pi; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n w(r_i) r_i \quad (3)$$

Where $w(r)$ is defined as:

$$\begin{aligned}w(r) &:= f_{\phi^*}(r) := \operatorname{argmin} \mathbb{E}_{\phi} \left[(w(a, x) - f_{\phi}(r))^2 \right] \\ f_{\phi} &\in \{f_{\phi} : \mathbb{R} \rightarrow \mathbb{R} \mid \phi \in \Phi\}\end{aligned}\quad (4)$$

Motivated by these approaches, as well as the fact that estimating from *similar* actions or make a regression over rewards could prove challenging if a significant proportion of these actions are missing for a given context, we propose the *Context-Huddling Inverse Propensity Score* (CHIPS) estimator that we introduce in the next section.

3 THE CHIPS ESTIMATOR

The CHIPS estimator is based on the idea of partitioning the context space into clusters to extrapolate the behaviour of an agent when presented with a previously unseen or underrepresented context x . The assumption needed for this

approximation to the OPE problem is that, given a policy, all contexts belonging to a cluster c should have a similar probability of observing an action a and will observe similar rewards when that action is chosen. Formally, we will consider a finite partition of the context space as the cluster space $\mathcal{C} := \{\mathcal{C}_i\}_{i=1}^K$ with $\mathcal{C}_i \subset \mathcal{X}$ and $\mathcal{C}_i \cap \mathcal{C}_j = \emptyset$. We assume that we are given a $c \in \mathcal{C}$ for each context $x \in \mathcal{X}$, where we assume that c is drawn i.i.d from an unknown distribution $p(c|x)$. Thus, given a policy π , we can compute its value by refining Equation (1):

$$\begin{aligned}V(\pi) &:= \mathbb{E}_{p(x)p(c|x)\pi(a|x)p(r|a,c,x)}[r] \\ &= \mathbb{E}_{p(x)p(c|x)\pi(a|x)}[q(a, c, x)].\end{aligned}\quad (5)$$

Where we denote $q(a, c, x) := \mathbb{E}_{p(r|a,c,x)}[r]$ and it is important to note that $\mathbb{E}_{p(c|x)\pi(a|x)}[q(a, c, x)] = \mathbb{E}_{\pi(a|x)}[q(a, x)]$, and therefore the refinement is consistent with Equation (1). Similar to the common support condition in IPS, we formulate the following property as the equivalent for the CHIPS estimator of Assumption 2.1.

Assumption 3.1. (Common Cluster Support) Given an evaluation policy π and a logging policy π_0 , the latest has common cluster support for π if

$$p(a|c, \pi_0) > 0 \quad \forall a \in \mathcal{A}, c \in \mathcal{C} : p(a|c, \pi) > 0$$

Where we denote

$$p(a|c, \pi) = \int_{\mathcal{X}} \pi(a|x) p(x|c) dx$$

Assumption 3.1 is weaker than Assumption 2.1 since for a given triplet $(x, c, a) \in \mathcal{X} \times \mathcal{C} \times \mathcal{A}$, the fact that $\pi_0(a|x) = 0, \pi(a|x) > 0$ does not ensure the same holds for every context within c . The idea of a homogeneous behaviour for every context inside a given cluster would make the CHIPS estimator circumvent the bias increase when Assumption 2.1 is not met for the IPS estimator (if Assumption 3.1 holds). Regarding the reward, this concept is formalized in the following assumption.

Assumption 3.2. (Reward Homogeneity) We say that we observe reward homogeneity if the context x does not affect on the reward r given some action a and some context c (i.e., $r \perp x \mid c, a$).

The reward homogeneity assumption eliminates the dependency of the context on the reward when provided with the cluster and the action. Note that complying with Assumption 3.2 implies $q(a, c, x) = q(a, c, y) = q(a, c)$, where $x, y \in \mathcal{X}$, which together with Assumption 3.1 gives an alternative expression for the policy value in the following proposition:

Proposition 3.3. *Given a policy π , if Assumptions 3.1 and 3.2 hold, then we have that*

$$V(\pi) := \mathbb{E}_{p(c)p(a|c,\pi)}[q(a, c)] \quad (6)$$

Please refer to Appendix A.1 for a complete proof.

Considering the similarity of Equation (6) with the original policy value definition (Equation (1)), Proposition 3.3 naturally motivates the analytical expression of the CHIPS estimator:

$$\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D}) := \frac{1}{N} \sum_{i=1}^N \frac{p(a_i|c_i, \pi)}{p(a_i|c_i, \pi_0)} r_i = \frac{1}{N} \sum_{i=1}^N w(a_i, c_i) r_i$$

3.1 Theoretical Analysis

First, we characterize the bias of the CHIPS estimator depending on the compliance with Assumptions 3.1 and 3.2.

Proposition 3.4. *Under the Common Cluster Assumption (3.1) and the Cluster Homogeneity Assumption (3.2), the CHIPS estimator is unbiased for any given policy π :*

$$\mathbb{E}_{\mathcal{D}} [\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D})] = V(\pi)$$

Please refer to Appendix A.2 for a complete proof.

We note here that Proposition 3.4 implies that even when the Common Support Assumption (2.1) fails to ensure the unbiasedness of the IPS estimator, the CHIPS estimator can still use the more permissive Common Cluster Support (3.1), and the Reward Homogeneity (3.2) Assumption to ensure an unbiased estimate. Although Assumption 3.2 guarantees homogeneity at the reward level, a *completely* homogeneous behaviour would also eliminate the context dependency at the action level, implying a deterministic policy given cluster, i.e. $p(a|c, \pi) = \pi(a|x) \forall x \in c$. Both homogeneity conditions present a desirable scenario for the CHIPS estimator; however, they rarely occur when working in real-world data environments, which motivate the following assumption as a relaxation of the action-context independence:

Assumption 3.5. (δ -Homogeneity) Given a policy π , we say that the policy presents δ -homogeneity if for any given action $a \in \mathcal{A}$, and any given cluster $c \in \mathcal{C}$, there exist $\delta_{\pi, c, a}^- \leq 1$ and $\delta_{\pi, c, a}^+ \geq 1$ such that:

$$\delta_{\pi, c, a}^- \leq \frac{\pi(a|x)}{p(a|c, \pi)} \leq \delta_{\pi, c, a}^+ \quad \forall x \in \mathcal{X}$$

It is worth noting that if $p(a|c, \pi) \neq 0 \forall (x, c, a) \in \mathcal{D}$ then it is always possible to find $\delta_{\pi, c, a}^-$, $\delta_{\pi, c, a}^+$ satisfying δ -Homogeneity. The following proposition gives an upper bound for the bias of the CHIPS estimator when Assumption 3.2 cannot be ensured:

Proposition 3.6. *Given the logging data $\{(x_i, a_i, r_i)\}_{i=1}^N$ observed under some logging policy π_0 , and an evaluation policy π if the latest has common cluster support over the earliest, then we have that*

$$|\text{Bias}(\hat{V}_{\text{CHIPS}}(\pi))| \leq |\mathbb{E}_{p(c)p(x|c)p(a|c, \pi)}[q(a, c, x) \cdot \Delta_{c, a}]|$$

Where by Assumption 3.5 we have bounds $(\delta_{\pi, c, a}^-, \delta_{\pi, c, a}^+)$ for π , $(\delta_{\pi_0, c, a}^-, \delta_{\pi_0, c, a}^+)$ for π_0 , and we denote $\Delta_{c, a} =$

$\max\{\delta_{\pi, c, a}^+, \delta_{\pi_0, c, a}^+\} - \min\{\delta_{\pi, c, a}^-, \delta_{\pi_0, c, a}^-\}$. Please refer to Appendix A.3 for a complete proof.

Proposition 3.6 formalizes the intuition on how the bias of the estimator under Assumption 3.1 depends on the extent to which the contexts inside a cluster behave homogeneously under a given policy. Formally, the gap $\delta_{\pi}^+ - \delta_{\pi}^-$ determines how close the CHIPS is to being unbiased, being the case $\delta_{\pi, c, a}^- = \delta_{\pi, c, a}^+ = 1$ the perfect scenario. In this case, we have that $\pi(a|x) = p(a|c, \pi)$, which means that the weights in IPS $w(a, x) = w(a, c)$, and we could in theory substitute any context for any other within the same cluster for calculations, mitigating the problems that arise when Assumption 2.1 does not hold. Additionally, we can also provide an expression for the difference in mean squared error with respect to IPS in the same conditions as Proposition 3.6:

Proposition 3.7. *Under the same conditions as in Proposition 3.6, the difference in mean squared error between CHIPS and MIPS can be expressed as*

$$\begin{aligned} \text{MSE}(\hat{V}_{\text{IPS}}(\pi)) - \text{MSE}(\hat{V}_{\text{CHIPS}}(\pi)) \\ = \mathbb{V}_D[\hat{V}_{\text{IPS}}(\pi)] - \mathbb{V}_D[V_{\text{CHIPS}}(\pi; D)] \\ - \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi))^2 \end{aligned}$$

Please refer to Appendix A.4 for a complete proof.

It is also worth studying the bias of the CHIPS estimator when the Common Cluster Support assumption does not hold, while the Assumption 3.2 holds. For this purpose, we acknowledge that the bias of the IPS estimator when Assumption 2.1 is not met can be given in terms of the actions violating such assumption (Sachdeva, Su, and Joachims, 2020):

$$|\text{Bias}(\hat{V}_{\text{IPS}}(\pi; \mathcal{D}))| = \mathbb{E}_{p(x)} \left[\sum_{\mathcal{U}(x, \pi_0)} \pi(a|x) q(a, c, x) \right]$$

Where $\mathcal{U}(x, \pi_0) := \{a \in \mathcal{A} \mid \pi_0(a, x) = 0\}$ are known as the *deficient* actions. Following a similar approach we introduce the following proposition:

Proposition 3.8. *Given the logging policy π_0 and some evaluation policy π , the absolute bias of the CHIPS estimator when Assumption 3.2 holds can be expressed as*

$$|\text{Bias}(\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D}))| = \mathbb{E}_{p(c)} \left[\sum_{\mathcal{U}(c, \pi_0)} p(a|\pi, c) q(a, c) \right]$$

Where $\mathcal{U}(c, \pi_0) := \{a \in \mathcal{A} \mid p(a|\pi_0, c) = 0\}$. Please refer to Appendix A.5 for a complete proof.

Corollary 3.9. *Under the conditions of Proposition 3.8, we have that*

$$|\text{Bias}(\hat{V}_{\text{IPS}}(\pi; \mathcal{D}))| - |\text{Bias}(\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D}))|$$

$$= \mathbb{E}_{p(c)} \left[\sum_{\mathcal{U}(x, \pi_0) \setminus \mathcal{U}(c, \pi_0)} p(a | \pi, c) q(a, c) \right].$$

Please refer to Appendix A.5 for a complete proof.

Note that in this case, the CHIPS' reduction in absolute bias depends directly on the number of actions that violate Assumption 2.1, but still comply with Assumption 3.2. Thus, the greater the number of deficient actions by Common Support condition covered by the Common Cluster Support, the more significant the bias reduction with respect to IPS. In this conditions, its also interesting to study the difference in bias with respect to the other two transformation-based methods (MR and MIPS), a result given by the next proposition:

Proposition 3.10. *Let f_{ϕ^*} be defined as in Equation (4) with $f_{\phi^*} = w(a, x) + \epsilon$ for some $\epsilon \in \mathbb{R}$ and $e \in \mathcal{E}$ give action embeddings. Under the conditions of the Proposition 3.8, we have that:*

$$\begin{aligned} & |\text{Bias}(\hat{V}_{MR}; \mathcal{D})| - |\text{Bias}(\hat{V}_{CHIPS}; \mathcal{D})| \\ &= -\mathbb{E}_{p(c)} \left[\sum_{a \in (\mathcal{U}(c, \pi_0) \setminus \mathcal{U}(c, \pi_0))} q(a, c) p(a | \pi, c) \right] \\ & \quad + \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a | \pi_0, c) \right] \\ & |\text{Bias}(\hat{V}_{MIPS}; \mathcal{D})| - |\text{Bias}(\hat{V}_{CHIPS}; \mathcal{D})| \\ &= \mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{U}(e, \pi_0)} p(e | x, \pi) q(x, e) \right] \\ & \quad - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi_0)} p(a | c, \pi) q(a, c) \right]. \end{aligned}$$

Please refer to Appendix A.6 for a complete proof.

When studying homogeneity at an action level, we have focused on the probability of observing an action for a particular context x within a cluster c (i.e., $\pi(a|x)$). Conversely, we can also study the *predictability* of a context given an action and a cluster under a policy π , which we denote as $p(x|a, c) = \pi(x|a, c)$. Ideally, we would have that the conditional probability distribution of the context given the action and the cluster is uniform (i.e., $\pi(x_i|a, c) = \pi(x_j|a, c) \forall x_i, x_j \in c$). Predictability is used in the following proposition, that characterizes the relation between the reduction in variance of the CHIPS estimator with respect to IPS:

Proposition 3.11. *Given a logging policy π_0 , under the Common Support Assumption (2.1) and the Reward Homogeneity Assumption (3.2) we have that*

$$N(\mathbb{V}_{\mathcal{D}}[\hat{V}_{IPS}(\pi; \mathcal{D})] - \mathbb{V}_{\mathcal{D}}[\hat{V}_{CHIPS}(\pi; \mathcal{D})])$$

$$= \mathbb{E}_{p(c)p(a|c, \pi_0)} [\mathbb{V}_{\pi_0(x|a, c)} [w^2(a, x)] \mathbb{E}_{p(r|a, c)} [r^2]].$$

Note that this quantity is always positive, implying that CHIPS always reduces the variance of IPS. Please refer to Appendix A.7 for a complete proof.

Proposition 3.11 indicates that when Assumptions 2.1 and 3.2 hold, the variance reduction of CHIPS compared to IPS corresponds to the total decrease in mean squared error when approximating the actual policy value $V(\pi)$, as both estimators are unbiased under these conditions. This mean squared error gap is influenced by two factors: First, $\mathbb{E}_{p(r|a, c)} [r^2]$, reflecting the noise in rewards for actions within the same cluster (related to Assumption 3.2). Second, the variance of IPS weights conditioned on the predictability $p(x|a, c)$, which increases when $w(a, x)$ varies widely (e.g., when logging and evaluation policies differ) or when $\pi(x|a, c)$ is uninformative (contexts behave uniformly given the cluster and action). Thus, the variance reduction in CHIPS is particularly pronounced when IPS exhibits high variance and contexts within a cluster are similar. Furthermore, if MIPS and CHIPS are in the same space (considering contexts $c \in \mathcal{C}$ as described and action embeddings $e \in \mathcal{E}$), Proposition 3.11 can be extended to show that CHIPS has less variance than MIPS:

Proposition 3.12. *In context-action-embedding joint space $(\mathcal{X} \rightarrow \mathcal{C} \rightarrow \mathcal{A} \rightarrow \mathcal{E} \rightarrow [0, R_{max}])$, if Assumptions 3.1 and 3.2 hold, as well as their MIPS counterparts (Common Embedding Support and No Direct Effect), then we have that*

$$\mathbb{V}_{\mathcal{D}}(\hat{V}_{IPS}(\pi)) \geq \mathbb{V}_{\mathcal{D}}(\hat{V}_{MIPS}(\pi)) \geq \mathbb{V}_{\mathcal{D}}(\hat{V}_{CHIPS}(\pi)) \geq 0$$

Please refer to Appendix A.8 for a complete proof.

3.2 Empirical Calculations

The alternative analytical expression for the policy value given in Equation 6 eliminates the dependency on the original definition of policy value and motivates the CHIPS estimator under assumptions 3.1 and 3.2. However, in practice, assessing if such conditions hold is complicated, particularly if we have limited logging data. To mitigate this problem and justify using CHIPS in real-world settings, we need to make an approximation to context-homogeneous behavior on both action and reward levels within a cluster. In practice, we have a clustering method $\xi: \mathcal{X} \rightarrow \mathcal{C}$, and we use the transformation:

$$\begin{aligned} \tau: (\mathcal{X}, \mathcal{A}, [0, R_{max}]) &\rightarrow (\mathcal{X}, \mathcal{C}, \mathcal{A}, [0, R_{max}]) \\ (x, a, r) &\mapsto (x, \xi(x), a, r). \end{aligned}$$

Given a policy π and a cluster c , we use the definition to estimate $p(a|c, \pi)$:

$$p(a|c, \pi) = \int_{\mathcal{X}} \pi(a|x) p(x|c) dx$$

$$\begin{aligned}
 &= \int_{x \in c} \pi(a|x) p(x|c) dx \\
 &\approx \frac{1}{|\mathcal{D}_c|} \sum_{x \in \mathcal{D}_c} \pi(a|x), \quad (7)
 \end{aligned}$$

Here, we denote $\mathcal{D}_c = \{(x, \tilde{c}, a, r) \in \tau(\mathcal{D}) : \tilde{c} = c\}$. In Equation 7, we used that $p(x|c) = 0$ if $c \neq \xi(x)$. Since this equation is essentially $\mathbb{E}_{p(x|c)}[\pi(a|x)]$, we approximate this value by averaging $\pi(a|x)$ over all contexts inside the given cluster.

The second approximation needed involves the reward being independent of the context given the action and the cluster, i.e., $q(a, c, x) = q(a, c)$. Following a similar approach than in the previous case, for a particular (given) action a and cluster c , we observe that $q(a, c) = \mathbb{E}_{p(x|c)}[\pi(r|a, c, x)]$, which motivates the idea of an *average reward* per cluster. In our synthetic experiments, the reward is binary, therefore we will assume that the observations inside a cluster are observations in a Bernoulli process (i.e., $R_c \sim \text{Ber}(\theta)$) and estimate this average reward using two different approaches:

- **Maximum Likelihood (ML)** In which we just average the rewards observed within a cluster c for each action a as $\hat{r}_{\text{mean}}(a, c) = \frac{1}{|R_c|} \sum_{R_c} r_k$ with $R_c := \{r_k : (x_k, c_k, a_k, r_k) \in \mathcal{D}_c\}$.
- **Maximum A Posteriori (MAP)**. In this setting, estimating the average reward is equivalent to estimating the most probable θ using a beta prior, where we obtain:

$$\hat{r}_{\text{bayes}}(\alpha, \hat{\beta}; c) = \frac{(\alpha - 1) + \sum_{R_c} r_k}{\alpha + \hat{\beta} + |R_c| - 2}$$

Where we denote $\alpha, \hat{\beta}$ as the parameters of the prior Beta distribution. In our experiments, we use non-informative priors ($\alpha = \hat{\beta}$) (Tuyl, Gerlach, and Mengersen, 2008; Kerman, 2011) and we explore the choosing of this parameter for arbitrary problems in Appendix D.4. Please refer to Appendix B for the complete derivations of the MAP and ML estimations.

4 EXPERIMENTS

4.1 Synthetic dataset

We compare CHIPS with other baseline estimators (IPS, DM, DR, SNIPS (Swaminathan and Joachims, 2015b), DRoS (Su et al., 2020), SNDR (Thomas and Brunskill, 2016), MR (Taufiq et al., 2023)) in estimating the evaluation policy value in a cluster-based synthetic dataset in which we can control the difficulty of the OPE problem. A description of all hyperparameters used for generation (e.g., a_{num} , c_{exp} ...) can be found in Appendix C. We start by generating cluster centers $\mathcal{C} := \{c_k\}_{k=1}^m$ inside a

d_x -dimensional ball $B(0, c_{\text{exp}}) := \{x \in \mathbb{R}^{d_x} : \|x\|^2 < c_{\text{exp}}\}$ using a variation of the Box-Muller transformation (Box and Muller, 1958):

$$c_k = \frac{c_{\text{exp}} \cdot u_k^{-d_x} \cdot z_k}{\|z_k\|},$$

where $U := \{u_k\}_{k=1}^m \sim U[0, 1]$ and $Z := \{z_k\}_{k=1}^m \sim \mathcal{N}(0, \mathbb{I}_{d_x})$. We sample $S := \{s_k\}_{k=1}^m \sim U[0, 1]$, and use the softmax transformation $\phi(S)$ to define $p(c_i) = \phi(S)_i$. Then, we sample cluster centers according to this distribution $w = \{w_i\}_{i=1}^{x_{\text{num}}} \sim \phi(S)$, and, for each center c_i , we uniformly sample points belonging to the n -ball centered on c_i , using the same variation of the Box-Muller transform that we used previously:

$$\mathcal{X}_i = (x_i^1, \dots, x_i^{h_i}) \sim U[B(c_i, c_{\text{rad}})]$$

Note here that $h_i = \sum_{i=1}^{x_{\text{num}}} \mathbb{1}_{\{c_i = w_i\}}$. We define the context space as the union of these generated points $\mathcal{X} = \bigcup_{i=1}^m \mathcal{X}_i = \{x_i\}_{i=1}^{x_{\text{num}}}$. We sample $\mathcal{V} = \{v_i\}_{i=1}^{x_{\text{num}}} \sim \mathcal{N}(0, 1)$ and define $p(x_i) = \phi(\mathcal{V})_i$ using the ϕ softmax transformation again. We then use these probabilities to sample the logging (\mathcal{X}_{log}) and evaluation ($\mathcal{X}_{\text{eval}}$) data, with $|\mathcal{X}_{\text{eval}}| = e_{\text{len}}$ and $|\mathcal{X}_{\text{log}}| = b_{\text{len}}$. To generate the policies, we sample $y_i = \{y_i^j\}_{j=1}^{a_{\text{num}}} \sim \mathcal{N}(0, 1)$ for every cluster c_i (where a_{num} is the number of actions) and $z = \{z_k\}_{k=1}^{x_{\text{num}}} \sim \mathcal{N}(0, 1)$ to define the policies for every context in cluster c_i as:

$$\begin{aligned}
 \pi(a_j|c_i, x_k) &= \frac{e^{y_i^j + \sigma z_k}}{\sum_{m=1}^{a_{\text{num}}} e^{y_i^m + \sigma z_k}} \\
 \pi_0(a_j|c_i, x_k) &= \frac{e^{\beta(y_i^j + \sigma z_k)}}{\sum_{m=1}^{a_{\text{num}}} e^{\beta(y_i^m + \sigma z_k)}}, \quad -1 \leq \beta \leq 1
 \end{aligned}$$

Given a context x_k , both policies are determined by a term that depends on the cluster and the action (u_i^j), and a term that depends on the context itself (x_k). Here $0 \leq \sigma \leq 1$ controls how independent a policy is from the context and β how close the logging and evaluation policies are. For obtaining the actions, we sample $\mathcal{A}_{\text{log}} \sim \pi_0$ and $\mathcal{A}_{\text{eval}} \sim \pi$. For generating the rewards, we create a misspecified reward setting by defining:

$$r(a_i, c_i, x_i) = \mathbb{1} \left\{ u_i < \pi(a_i|c_i, x_i) \cdot \frac{\|x_i\|_1}{c_{\text{exp}} d_x} \right\},$$

where $u_i \sim U[0, 1]$. The reward depends on two factors; the first one is the Manhattan norm of the context; the further from 0, the more likely it is to observe a positive reward. The second factor is the evaluation policy $\pi(a_i|c_i, x_i)$, which makes this a misspecified reward setting when the logging and evaluation policies are different enough. In this case, the (a_i, c_i, x_i) triplets having the highest probability of observation under the evaluation policy are more likely to observe positive rewards, resulting in a significant difference with respect to the observed rewards under the logging policy for such triplets. We sample

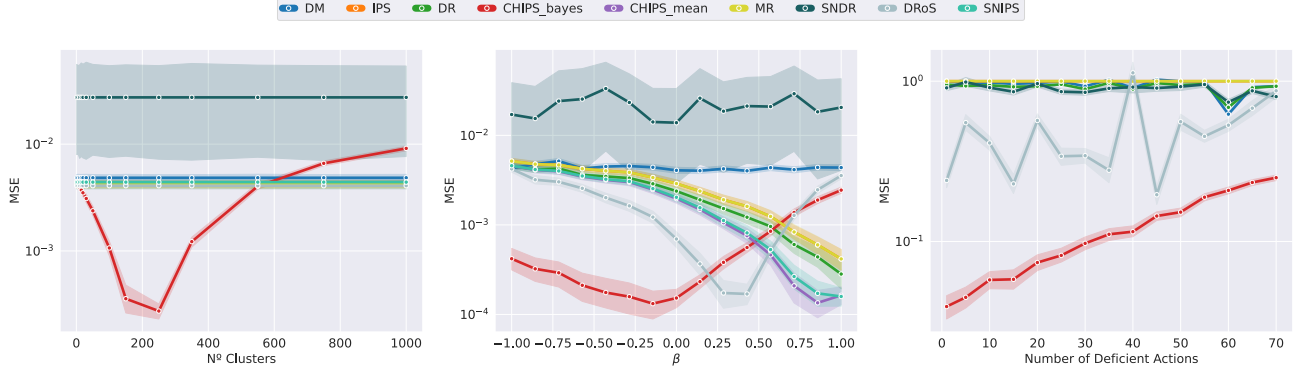


Figure 1: From left to right, the mean square error in the synthetic dataset experiments varying the number of clusters, the distributional shift between logging and evaluation policy (β), and the number of deficient actions in the logging data (normalized w.r.t. IPS).

rewards using this method for the logging (\mathcal{R}_{\log}) and evaluation ($\mathcal{R}_{\text{eval}}$) data to obtain $\mathcal{D}_{\log} := (\mathcal{X}, \mathcal{C}, \mathcal{A}_{\log}, \mathcal{R}_{\log})$ and $\mathcal{D}_{\text{eval}} := (\mathcal{X}, \mathcal{C}, \mathcal{A}_{\text{eval}}, \mathcal{R}_{\text{eval}})$. Finally, we select a subset for N samples from both sets. A representation of the generated structure can be found in Figure 20.

4.1.1 Synthetic results

In this section we analyze CHIPS performance while varying parameters of the synthetic dataset. In our experiments, the generation process for each parameter value is repeated 100 times with different random seeds. The final reported results are the average over all experiments, with the standard deviation corresponding to the lighter bands represented in all the figures. The basic configuration for the parameters used throughout the experiments can be found in Appendix C, along with the specifications of the hardware used. We use Random Forest (Breiman, 2001) to obtain $\hat{q}(x, a)$ in DM-based methods and mini-batch KMeans (Sculley, 2010) implementation in SciKit-Learn (Pedregosa et al., 2011) as the clustering method for CHIPS (alternative clustering methods and their performance are also discussed in Appendix D). We also use $\beta = -1$, maximizing the distributional shift between logging and evaluation policies.

Number of clusters. For this experiment, we vary the number of clusters the CHIPS estimator uses, with values ranging from 1 to 1000. Since $\beta = -1$, the implementation of CHIPS using ML reward estimation is unsuccessful (see Appendix D.3 for a further discussion). On the other hand, for the MAP case, we observe a v-shaped error graph (see Figure 1 (left)), suggesting that CHIPS performance is sensitive to effectiveness of clustering. In particular, we have a highly biased estimation when assuming insufficient or excessive clusters (see Figure 3). The reason for this bias in the first case might be an oversimplification of the structure of the cluster space. Conversely, we progressively gain bias when we select too many clusters according to Propo-

sition 3.8 as CHIPS converges to IPS. In this case, CHIPS is also vulnerable to reward misspecification, which causes an increase in variance. In practice, this parameter can be selected by considering the possible CHIPS estimates as a parametric family depending on the number of clusters and use the PAS-IF technique (Udagawa et al., 2023) to choose the optimal number of clusters.

Beta. This experiment examines the impact of the distributional policy shift between π and π_0 . Lower values in our range (i.e., $\pi_0 \longleftrightarrow \pi$) result in significant policy shifts that introduce bias in IPS estimates for large context-action spaces (Saito and Joachims, 2022; Sachdeva, Su, and Joachims, 2020). The CHIPS estimator mitigates this by treating all context-action samples within a cluster as if they share the same context. However, when β is low, these virtual extra samples may not suffice for accurate estimation, as the most relevant (x, a) pairs ($\pi(a|x)$ near 1) are underrepresented (see Appendix D.3). In such cases, ML estimation in CHIPS is ineffective, while MAP estimation provides some resistance by pushing reward estimates towards the posterior expectation, making it sensitive to prior choice. However, this resistance can be counterproductive when the distributional shift is small (β close to 1), as both ML estimates and IPS converge faster to more accurate estimations (see Figure 1 (center)).

Deficient actions. In this setting we explicitly set the probability (π_0) of observing a variable number of actions in the action space to 0 and evaluate CHIPS’ response in a space with 200 actions and $\beta = -1$. This setting is quite challenging as not only we have deficient actions but also a significant distributional shift between policies. The majority of baselines perform at a similar level than IPS with the exception of DRoS (Su et al., 2020), that performs slightly better but is still outperformed by CHIPS.

Additional experiments and discussions of results varying other parameters, different clustering methods, and a time

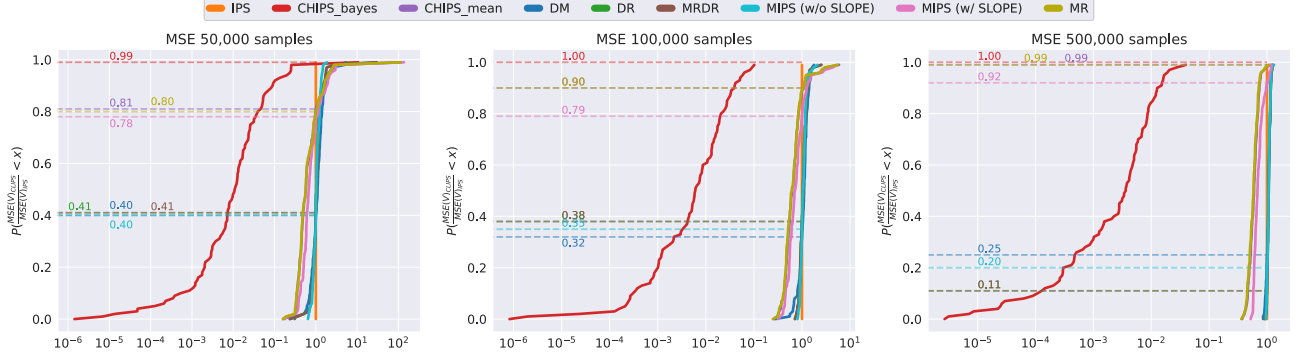


Figure 2: ECDF of the relative mean squared error with respect to IPS for the real dataset using 50000 (left), 100000 (center), and 500000 (right) logging samples.

complexity analysis can be found in Appendix G.

4.2 Real dataset

Following the literature, for assessing the capabilities of the CHIPS estimator in a real-world environment, we compare the performance in the Open Bandit Dataset (OBD) (Saito et al., 2020) of IPS, DM, DR, MRDR (Farajtabar, Chow, and Ghavamzadeh, 2018) and MIPS (Saito and Joachims, 2022), with and without SLOPE (Su, Srinath, and Krishnamurthy, 2020). The OBD dataset was gathered using two different policies during an A/B test: uniform random, which we consider as logging (i.e., π_0), and Thompson sampling (Thompson, 1933, 1935), which we consider as evaluation (i.e., π). The dataset is based on a recommendation system for fashion e-commerce. We observe user data as contexts x , items to recommend $a \in \mathcal{A}$ (with $|\mathcal{A}| = 240$) and rewards $r \in \{0, 1\}$ representing user interactions.

Following the experimental protocol of Saito and Joachims (2022) (see Appendix F), we experiment with the real dataset varying the number of logging samples available for the estimation using 50 000, 100 000, and 500 000 samples to compute the Empirical Cumulative Distribution Function (ECDF) of the normalized mean squared error with respect to IPS. We increase the number of clusters for CHIPS as more logging samples are available to try to maximize performance, following the intuition from our earlier experiments on the synthetic dataset (see Figure 12 (right)). We use 8 clusters for 100 000 samples as a reference from our results for 240 actions in the synthetic dataset (see Figure 12 (left)). Regarding the clustering method, we use again mini-batch KMeans.

We observe that the CHIPS estimator using the ML approximation is slightly better (+3%) than MIPS when few samples are available (see Figure 2, (left)). This performance gap widens (+11%) as the CHIPS estimator has more samples available (see Figure 2, (center)) and starts narrowing (+7%) as the number of samples is enough for MIPS to also start making more accurate estimations (see Figure 2,

(right)).

Using the MAP reward estimation for CHIPS provides a considerable advantage in all experiments since the real dataset present severe reward misspecification, as discussed in Appendix D.3. Similarly to the synthetic dataset, the partition structure of the cluster space and the α parameter in MAP are sensitive parameters. In particular, for the number of clusters, we observe that using an insufficient or excessive number of clusters can negatively impact performance (see Figure 14 (left)) as we discussed in section Section 4.1.1. Regarding the value of α for the Beta prior, following the results from the synthetic experiment studying the effect of this parameter conjointly with the distributional shift between logging and evaluation policies (see discussion in Appendix D.4 and Figure 13), we used $\alpha = 20$ as the logging policy is uniform (the equivalent of $\beta = 0$ in the synthetic dataset). Figure 14 (right) shows how choosing a lower or higher value for α deteriorates the performance of the CHIPS estimator, reaffirming the results observed in the synthetic dataset (see Figure 13).

5 CONCLUSIONS, LIMITATIONS AND FUTURE WORK

In this work we have explored an alternative approach to the OPE problem by clustering contexts instead of pooling information over actions to mitigate the problems arising in IPS when the Common Support condition does not hold. The proposed setup for the OPE problem using contexts led to the CHIPS estimator, which uses a similar approach to IPS applied over clusters instead of contexts. We have studied this estimator extensively from a theoretical and practical perspective, evaluating its performance for different configurations in a controlled synthetic dataset and a real-world example. The results obtained in the experiments for both cases demonstrate that the CHIPS estimator provides a significant improvement in estimation accuracy, outperforming existing estimators if the context space has a clus-

ter structure. The accuracy of CHIPS is also influenced by the accuracy of the clustering method and the homogeneity behaviour of contexts inside the same cluster. Additionally, choosing a balanced number of clusters to avoid over- and under-simplification of the cluster structure is an important part of the estimation process and opens the possibility of exploring if it is possible to estimate the optimal value for hyperparameters beyond empirical estimation or even if combining CHIPS with pure action-embedding methods like MIPS can improve general performance.

References

- Ankerst, M.; Breunig, M. M.; Kriegel, H.-P.; and Sander, J. 1999. OPTICS: Ordering Points to Identify the Clustering Structure. *SIGMOD Rec.*, 28(2): 49–60.
- Attias, H. 1999. A Variational Bayesian Framework for Graphical Models. In Solla, S.; Leen, T.; and Müller, K., eds., *Advances in Neural Information Processing Systems*, volume 12. MIT Press.
- Bastani, H.; and Bayati, M. 2019. Online Decision Making with High-Dimensional Covariates. *Operations Research*, 68.
- Bendada, W.; Salha, G.; and Bontempelli, T. 2020. Carousel Personalization in Music Streaming Apps with Contextual Bandits. In *Proceedings of the 14th ACM Conference on Recommender Systems*, RecSys '20, 420–425. New York, NY, USA: Association for Computing Machinery. ISBN 9781450375832.
- Beygelzimer, A.; and Langford, J. 2009. The Offset Tree for Learning with Partial Labels. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '09, 129–138. New York, NY, USA: Association for Computing Machinery. ISBN 9781605584959.
- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag. ISBN 0387310738.
- Blei, D.; and Jordan, M. 2006. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1.
- Bouneffouf, D.; Rish, I.; and Aggarwal, C. 2020. Survey on Applications of Multi-Armed and Contextual Bandits. In *2020 IEEE Congress on Evolutionary Computation (CEC)*, 1–8.
- Box, G. E. P.; and Muller, M. E. 1958. A Note on the Generation of Random Normal Deviates. *The Annals of Mathematical Statistics*, 29(2): 610 – 611.
- Breiman, L. 2001. Random Forests. *Machine Learning*, 45: 5–32.
- Comaniciu, D.; and Meer, P. 2002. Mean Shift: A Robust Approach Toward Feature Space Analysis. *IEEE Trans. Pattern Anal. Mach. Intell.*, 24(5): 603–619.
- Deng, L. 2012. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6): 141–142.
- Dua, D.; and Graff, C. 2017. UCI Machine Learning Repository.
- Dudík, M.; Langford, J.; and Li, L. 2011. Doubly Robust Policy Evaluation and Learning. In *International Conference on Machine Learning*.
- Dumitrescu, B.; Feng, K.; and Engelhardt, B. E. 2018. PG-TS: Improved Thompson Sampling for Logistic Contextual Bandits. In *Neural Information Processing Systems*.
- Ester, M.; Kriegel, H.-P.; Sander, J.; and Xu, X. 1996. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, KDD'96, 226–231. AAAI Press.
- Everitt, B. 1996. An introduction to finite mixture distributions. *Statistical Methods in Medical Research*, 5(2): 107–127. PMID: 8817794.
- Farajtabar, M.; Chow, Y.; and Ghavamzadeh, M. 2018. More Robust Doubly Robust Off-policy Evaluation. In Dy, J.; and Krause, A., eds., *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, 1447–1456. PMLR.
- Felicioni, N.; Dacrema, M. F.; Restelli, M.; and Cremonesi, P. 2022. Off-Policy Evaluation with Deficient Support Using Side Information. In *Neural Information Processing Systems*.
- Frey, B. J.; and Dueck, D. 2007. Clustering by Passing Messages Between Data Points. *Science*, 315(5814): 972–976.
- Guo, Y.; Liu, H.; Yue, Y.; and Liu, A. 2024. Distributionally Robust Policy Evaluation under General Covariate Shift in Contextual Bandits. *ArXiv*, abs/2401.11353.
- Horvitz, D. G.; and Thompson, D. J. 1952. A Generalization of Sampling Without Replacement from a Finite Universe. *Journal of the American Statistical Association*, 47: 663–685.
- Irpan, A.; Rao, K.; Bousmalis, K.; Harris, C.; Ibarz, J.; and Levine, S. 2019. Off-Policy Evaluation via Off-Policy Classification. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Kallus, N.; and Uehara, M. 2019. *Intrinsically Efficient, Stable, and Bounded off-Policy Evaluation for Reinforcement Learning*. Red Hook, NY, USA: Curran Associates Inc.

- Kerman, J. 2011. Neutral noninformative and informative conjugate beta and gamma prior distributions. *Electronic Journal of Statistics*, 5.
- Kiyohara, H.; Uehara, M.; Narita, Y.; Shimizu, N.; Yamamoto, Y.; and Saito, Y. 2023. Off-Policy Evaluation of Ranking Policies under Diverse User Behavior. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '23, 1154–1163. New York, NY, USA: Association for Computing Machinery. ISBN 9798400701030.
- Krizhevsky, A.; Nair, V.; and Hinton, G. 2009. CIFAR-100 (Canadian Institute for Advanced Research).
- Kuzborskij, I.; Vernade, C.; Gyorgy, A.; and Szepesvari, C. 2020. Confident Off-Policy Evaluation and Selection through Self-Normalized Importance Weighting. *ArXiv*, abs/2006.10460.
- Li, L.; Chu, W.; Bellevue, M.; Langford, J.; and Wang, X. 2011. An Unbiased Offline Evaluation of Contextual Bandit Algorithms with Generalized Linear Models. *Journal of Machine Learning Research*, 1.
- McNellis, R.; Elmachtoub, A. N.; Oh, S.; and Petrik, M. 2017. A Practical Method for Solving Contextual Bandit Problems Using Decision Trees. In Elidan, G.; Kersting, K.; and Ihler, A. T., eds., *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence, UAI 2017, Sydney, Australia, August 11-15, 2017*. AUAI Press.
- Metelli, A. M.; Russo, A.; and Restelli, M. 2021. Subgaussian and Differentiable Importance Sampling for Off-Policy Evaluation and Learning. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 8119–8132. Curran Associates, Inc.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; and Duchesnay, E. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12: 2825–2830.
- Peng, J.; Zou, H.; Liu, J.; Li, S.; Jiang, Y.; Pei, J.; and Cui, P. 2023. Offline Policy Evaluation in Large Action Spaces via Outcome-Oriented Action Grouping. In *Proceedings of the ACM Web Conference 2023*, WWW '23, 1220–1230. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394161.
- Sachdeva, N.; Su, Y.-H.; and Joachims, T. 2020. Off-policy Bandits with Deficient Support. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*.
- Sachdeva, N.; Wang, L.; Liang, D.; Kallus, N.; and McAuley, J. 2023. Off-Policy Evaluation for Large Action Spaces via Policy Convolution. *arXiv*:2310.15433.
- Saito, Y.; Aihara, S.; Matsutani, M.; and Narita, Y. 2020. Open Bandit Dataset and Pipeline: Towards Realistic and Reproducible Off-Policy Evaluation. In *NeurIPS Datasets and Benchmarks*.
- Saito, Y.; and Joachims, T. 2022. Off-Policy Evaluation for Large Action Spaces via Embeddings. In *International Conference on Machine Learning*.
- Saito, Y.; Ren, Q.; and Joachims, T. 2023. Off-Policy Evaluation for Large Action Spaces via Conjunct Effect Modeling. In Krause, A.; Brunskill, E.; Cho, K.; Engelhardt, B.; Sabato, S.; and Scarlett, J., eds., *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, 29734–29759. PMLR.
- Saito, Y.; Udagawa, T.; Kiyohara, H.; Mogi, K.; Narita, Y.; and Tateno, K. 2021. Evaluating the Robustness of Off-Policy Evaluation. In *Proceedings of the 15th ACM Conference on Recommender Systems*, RecSys '21, 114–123. New York, NY, USA: Association for Computing Machinery. ISBN 9781450384582.
- Sculley, D. 2010. Web-Scale k-Means Clustering. In *Proceedings of the 19th International Conference on World Wide Web*, WWW '10, 1177–1178. New York, NY, USA: Association for Computing Machinery. ISBN 9781605587998.
- Shi, J.; and Malik, J. 2000. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8): 888–905.
- Shimizu, T.; and Forastiere, L. 2023. Doubly Robust Estimator for Off-Policy Evaluation with Large Action Spaces. *arXiv*:2308.03443.
- Su, Y.; Dimakopoulou, M.; Krishnamurthy, A.; and Dudik, M. 2020. Doubly robust off-policy evaluation with shrinkage. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 9167–9176. PMLR.
- Su, Y.; Srinath, P.; and Krishnamurthy, A. 2020. Adaptive Estimator Selection for Off-Policy Evaluation. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 9196–9205. PMLR.
- Swaminathan, A.; and Joachims, T. 2015a. Counterfactual Risk Minimization: Learning from Logged Bandit Feedback. In Bach, F.; and Blei, D., eds., *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, 814–823. Lille, France: PMLR.
- Swaminathan, A.; and Joachims, T. 2015b. The Self-Normalized Estimator for Counterfactual Learning. In *Advances in Neural Information Processing Systems*, 28.

- Taufiq, M. F.; Doucet, A.; Cornish, R.; and Ton, J.-F. 2023. Marginal Density Ratio for Off-Policy Evaluation in Contextual Bandits. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Thomas, P. S.; and Brunskill, E. 2016. Data-efficient off-policy policy evaluation for reinforcement learning. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, 2139–2148. JMLR.org.
- Thompson, W. R. 1933. On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika*, 25: 285–294.
- Thompson, W. R. 1935. On the Theory of Apportionment. *American Journal of Mathematics*, 57(2): 450–456.
- Tuyl, F.; Gerlach, R.; and Mengersen, K. 2008. A Comparison of Bayes–Laplace, Jeffreys, and Other Priors. *American Statistician - AMER STATIST*, 62: 40–44.
- Udagawa, T.; Kiyohara, H.; Narita, Y.; Saito, Y.; and Tateno, K. 2023. Policy-Adaptive Estimator Selection for Off-Policy Evaluation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37: 10025–10033.
- Varatharajah, Y.; and Berry, B. 2022. A Contextual-Bandit-Based Approach for Informed Decision-Making in Clinical Trials. *Life*, 12(8).
- Voloshin, C.; Le, H.; Jiang, N.; and Yue, Y. 2021. Empirical Study of Off-Policy Policy Evaluation for Reinforcement Learning. In Vanschoren, J.; and Yeung, S., eds., *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran.
- Wang, Y.-X.; Agarwal, A.; and Dudík, M. 2017. Optimal and Adaptive Off-Policy Evaluation in Contextual Bandits. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, 3589–3597. JMLR.org.
- Ward, J. H. 1963. Hierarchical Grouping to Optimize an Objective Function. *Journal of the American Statistical Association*, 58(301): 236–244.
- Zhan, R.; Hadad, V.; Hirshberg, D. A.; and Athey, S. 2021. Off-Policy Evaluation via Adaptive Weighting with Data from Contextual Bandits. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*.
- Zhang, T.; Ramakrishnan, R.; and Livny, M. 1996. BIRCH: An Efficient Data Clustering Method for Very Large Databases. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*, SIGMOD '96, 103–114. New York, NY, USA: Association for Computing Machinery. ISBN 0897917944.

Checklist

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model.
[Yes]. Justification: The mathematical setting of both synthetic and real experiments is detailed in Section 4, the theoretical assumptions, its analysis when they don't hold, and the main algorithm are detailed in Section 3.
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm.
[Yes]. Justification: A complete analysis on the complexity of the method can be found in Appendix G, while the analysis of properties can be found in Section 3.
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes/No/Not Applicable]
[Yes]. Justification: The code, an explanation on how to execute it and a Poetry environment to take care of the dependencies are included in the supplemental materials.
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results.
[Yes]. Justification: The full set of assumptions, relaxation of them and all the theoretical analysis of the method can be found in Section 3.
 - (b) Complete proofs of all theoretical results.
[Yes]. Justification: For every proposition in the theoretical analysis (Section 3), the proof is referenced and can be found in Appendix A.
 - (c) Clear explanations of any assumptions.
[Yes]. Justification: For every proposition in the theoretical analysis (Section 3), the assumptions used for the results are detailed in the introduction of the proposition.
3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL).
[Yes]. Justification: The code contains the necessary scripts to reproduce every experiment in the paper, as well as instructions on how to execute it.
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen).

[Yes]. **Justification:** The training details and parameter variations can be found in Section 4. Extra experiments varying more parameters, 2 parameters at the same time or different clustering options can be found in Appendix D. Additionally, a table with the base value of every parameter can be found in Appendix C.

- (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times).

[Yes]. **Justification:** All deviation bars obtained after 100 runs with different seeds for every parameter are included in the result figures as explained in Section 4.

- (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider).

[Yes]. **Justification:** The table with the specifications of the computing infrastructure used for the experiments can be found in Appendix C.

- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:

- (a) Citations of the creator If your work uses existing assets.

[Yes]. **Justification:** The evaluation pipeline that we use for the real experiments, the systems that we compare our method with, and the code implementations of known algorithms for different libraries that we use are referenced in the paper.

- (b) The license information of the assets, if applicable.

[Yes]. **Justification:** The license included in the code does not conflict with the license of the used resources.

- (c) New assets either in the supplemental material or as a URL, if applicable.

[Yes]. **Justification:** We provide all the code with our pipeline and implementations for common open algorithms in the supplemental material.

- (d) Information about consent from data providers/curators.

[Yes]. **Justification:** The code includes a licence to use.

- (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content.

[Not Applicable]. **Justification:** There is not personally identifiable information or offensive content of any kind in our data or experiments.

- 5. If you used crowdsourcing or conducted research with human subjects, check if you include:

- (a) The full text of instructions given to participants and screenshots.

[Not Applicable]. **Justification:** Our research did not use crowdsourcing or was conducted with human subjects.

- (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable.

[Not Applicable]. **Justification:** Our research did not use crowdsourcing or was conducted with human subjects.

- (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation.

[Not Applicable]. **Justification:** Our research did not use crowdsourcing or was conducted with human subjects.

A THEORETICAL RESULTS PROOFS

A.1 Proposition 3.3

Given a policy π , if both Assumption 3.1 and 3.2 hold, from the refinement of the policy value definition in a cluster-based bandits process (introduced in Section 3), we have that:

$$\begin{aligned} V(\pi) &:= \mathbb{E}_{p(x)p(c|x)\pi(a|x)p(r|a,c,x)} [r] \\ &= \mathbb{E}_{p(c)p(x|c)\pi(a|x)} [q(a, c, x)] \end{aligned} \quad (8)$$

$$= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} p(a|c) q(a, c) dx \right] \quad (9)$$

$$= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} \sum_{a \in \mathcal{A}} p(x|c) \pi(a|x) q(a, c) dx \right]$$

$$= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) \int_{\mathcal{X}} p(x|c) \pi(a|x) dx \right]$$

$$= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} p(a|c, \pi) q(a, c) \right] \quad (10)$$

$$= \mathbb{E}_{p(c)p(a|c,\pi)} [q(a, c)]$$

Where in Equation 8 we used the Bayes Theorem, in Equation 9 the fact that under Assumption 3.2 $q(a, c, x) = q(a, c)$, and the definition of $p(a|c, \pi)$ in Equation 10.

A.2 Proposition 3.4

Given a policy π and under Assumptions 3.1 and 3.2 we have that:

$$\mathbb{E}_{\mathcal{D}} [\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D})] = \mathbb{E}_{\mathcal{D}} [w(a, c)r] \quad (11)$$

$$\begin{aligned} &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)p(r|a,c,x)} [w(a, c)r] \\ &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)} [w(a, c)q(a, c)] \end{aligned} \quad (12)$$

$$= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi_0(a|x) w(a, c) q(a, c) dx \right]$$

$$= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} \sum_{a \in \mathcal{A}} p(x|c) \pi_0(a|x) w(a, c) q(a, c) dx \right]$$

$$= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} w(a, c) q(a, c) \left(\int_{\mathcal{X}} p(x|c) \pi_0(a|x) dx \right) \right]$$

$$= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} \frac{p(a|c, \pi)}{p(a|c, \pi_0)} q(a, c) p(a|c, \pi_0) \right] \quad (13)$$

$$= \mathbb{E}_{p(c)p(a|c,\pi)} [q(a, c)]$$

$$\begin{aligned} &= \mathbb{E}_{p(x)p(c|x)\pi(a|x)p(r|a,c,x)} [r] \\ &= V(\pi) \end{aligned} \quad (14)$$

In Equation 11, we have used the linearity of expectation, in Equation 12 the definition of $q(a, c, x)$ and Assumption 3.2. Equation 13 is just using the definition of $p(a|c, \pi)$ while Equation 14 is a combination of Proposition 3.3 and the equivalence $q(a, c) = q(a, c, x)$ under the given assumptions.

A.3 Proposition 3.6

Given the logging data $\mathcal{D} = \{(x_i, a_i, r_i)\}$, a logging policy π_0 , and an evaluation policy π having common cluster support over it, we have that:

$$\text{Bias}(\hat{V}_{\text{CHIPS}}(V; \mathcal{D})) = \mathbb{E}_{\mathcal{D}} [w(c, a)r] - V(\pi)$$

$$\begin{aligned}
 &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)p(r|a,c,x)} [w(a, c)r] - V(\pi) \\
 &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)} [w(a, c)q(a, c, x)] - V(\pi) \\
 &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)} [w(a, c)q(a, c, x)] - \mathbb{E}_{p(x)p(c|x)\pi(a|x)} [q(a, c, x)] \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)} [w(a, c)q(a, c, x)] - \mathbb{E}_{p(c)p(x|c)\pi(a|x)} [q(a, c, x)] \\
 &= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi_0(a|x) w(c, a) q(a, c, x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi(a|x) q(a, c, x) dx \right]
 \end{aligned} \tag{15}$$

Under Assumption 3.5 we have that $\delta_{\pi, c, a}^- \leq \frac{\pi(a|c, x)}{p(a|c, \pi)} \leq \delta_{\pi, c, a}^+$, $\delta_{\pi_0, c, a}^- \leq \frac{\pi_0(a|c, x)}{p(a|c, \pi_0)} \leq \delta_{\pi_0, c, a}^+$ $\forall x \in \mathcal{X}$ given an action $a \in \mathcal{A}$ and a context $c \in \mathcal{C}$. We denote then $\delta_{c, a}^+ = \max\{\delta_{\pi, c, a}^+, \delta_{\pi_0, c, a}^+\}$, $\delta_{c, a}^- = \min\{\delta_{\pi, c, a}^-, \delta_{\pi_0, c, a}^-\}$, $\Delta_{a, c} = \delta_{c, a}^+ - \delta_{c, a}^-$, and we can give an upper bound as follows:

$$\begin{aligned}
 &\mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi_0(a|x) w(c, a) q(a, c, x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi(a|x) q(a, c, x) dx \right]
 \end{aligned} \tag{16}$$

$$\begin{aligned}
 &\leq \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \delta_{c, a}^+ p(a|c, \pi_0) w(c, a) q(a, c, x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \delta_{c, a}^- p(a|c, \pi) q(a, c, x) dx \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} p(a|c, \pi) \int_{\mathcal{X}} p(x|c) \delta_{c, a}^+ \frac{p(a|c, \pi)}{p(a|c, \pi_0)} p(a|c, \pi_0) q(a, c, x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} p(a|c, \pi) \int_{\mathcal{X}} p(x|c) \delta_{c, a}^- q(a, c, x) dx \right] \\
 &= \mathbb{E}_{p(c)p(a|c, \pi)} \left[\delta_{c, a}^+ \int_{\mathcal{X}} p(x|c) q(a, c, x) dx \right] - \mathbb{E}_{p(c)p(a|c, \pi)} \left[\delta_{c, a}^- \int_{\mathcal{X}} p(x|c) q(a, c, x) dx \right] \\
 &= \mathbb{E}_{p(c)p(a|c, \pi)} [\mathbb{E}_{p(x|c)} [q(a, c, x)] (\delta_{c, a}^+ - \delta_{c, a}^-)] \\
 &= \mathbb{E}_{p(c)p(a|c, \pi)} [\mathbb{E}_{p(x|c)} [q(a, c, x)] \Delta_{a, c}]
 \end{aligned} \tag{17}$$

Note that in Equation 17 we can follow an analogous path to establish a lower bound:

$$\begin{aligned}
 &\mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi_0(a|x) w(c, a) q(a, c, x) dx \right] - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi(a|x) q(a, c, x) dx \right] \\
 &\geq \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \delta_{c, a}^- p(a|c, \pi_0) w(c, a) q(a, c, x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \delta_{c, a}^+ p(a|c, \pi) q(a, c, x) dx \right] \\
 &= -\mathbb{E}_{p(c)p(a|c, \pi)} [\mathbb{E}_{p(x|c)} [q(a, c, x)] \Delta_{a, c}]
 \end{aligned}$$

From which we have:

$$|\text{Bias}(\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D}))| \leq |\mathbb{E}_{p(c)p(x|c)p(a|c, \pi)} [q(a, c, x) \cdot \Delta_{a, c}]|$$

A.4 Proposition 3.7

Since the observations are independent we have that

$$\begin{aligned}
 &N (\text{MSE}(\hat{V}_{\text{IPS}}(\pi)) - \text{MSE}(\hat{V}_{\text{CHIPS}}(\pi))) \\
 &= \mathbb{V}_{x, a, r} [\omega(x, a)r] - \mathbb{V}_{c, a, r} [\omega(a, c)r] - N \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi))^2
 \end{aligned}$$

We now analyze the difference in variance:

$$\begin{aligned}
 & V_{p(c)p(x|c)\pi_0(a|x)p(r|a,c,x)}[\omega(x, a)r] - V_{p(c)p(x|c)\pi_0(a|x)p(r|a,c,x)}[\omega(a, c)r] \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)p(r|a,c,x)}[\omega(x, a)r^2] - V(\pi)^2 \\
 &\quad - \left(\mathbb{E}_{p(c)p(x|c)\pi_0(a|x)p(r|a,c,x)}[\omega(a, c)^2 \cdot r^2] - (V(\pi) + \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi)))^2 \right) \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)}[(\omega(x, a)^2 - \omega(a, c)^2) \mathbb{E}_{p(r|a,c,x)}[r^2]] \\
 &\quad + 2V(\pi) \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi)) + \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi))^2
 \end{aligned}$$

This implies that

$$\begin{aligned}
 & N(\text{MSE}(\hat{V}_{\text{IPS}}(\pi)) - \text{MSE}(\hat{V}_{\text{CHIPS}}(\pi))) \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)}[(\omega(x, a)^2 - \omega(a, c)^2) \mathbb{E}_{p(r|a,c,x)}[r^2]] \\
 &\quad + 2V(\pi) \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi)) + (1 - N) \text{Bias}(\hat{V}_{\text{CHIPS}}(\pi))^2
 \end{aligned}$$

A.5 Proposition 3.8

Given the logging policy π_0 and some evaluation policy π , the absolute bias of the CHIPS estimator when Assumption 3.2, we have that:

$$\begin{aligned}
 \text{Bias}(\hat{V}_{\text{CHIPS}}(V; \mathcal{D})) &= \mathbb{E}_{\mathcal{D}}[w(c, a)r] - V(\pi) \\
 &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)p(r|a,c,x)}[w(a, c)r] - V(\pi) \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)}[w(a, c)q(a, c)] - \mathbb{E}_{p(c)p(x|c)\pi(a|x)}[w(a, c)q(a, c)] \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)}[w(a, c)q(a, c)] - \mathbb{E}_{p(c)p(x|c)\pi(a|x)}[w(a, c)q(a, c)] \\
 &= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi_0(a|x) w(c, a) q(a, c) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi(a|x) q(a, c) dx \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} w(c, a) q(a, c) \int_{\mathcal{X}} p(x|c) \pi_0(a|x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) \int_{\mathcal{X}} p(x|c) \pi(a|x) dx \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} w(c, a) q(a, c) p(a|c, \pi_0) \right] - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a|c, \pi) \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi_0)^c} w(c, a) q(a, c) p(a|c, \pi_0) \right] - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a|c, \pi) \right] \tag{18} \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi_0)^c} \frac{p(a|c, \pi_0)}{p(a|c, \pi_0)} p(a|c, \pi_0) q(a, c) \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a|c, \pi) \right] \\
 &= \mathbb{E}_{p(c)} \left[- \sum_{a \in \mathcal{U}(c, \pi_0)} p(a|c, \pi) q(a, c) \right]
 \end{aligned}$$

Where in Equation 18 we note that $p(a|c, \pi_0) = 0$ if $a \in \mathcal{U}(c, \pi_0)$. Following an analogous procedure we can give an expression for the bias of IPS in a cluster bandits setup:

$$\begin{aligned}
 \text{Bias}(\hat{V}_{\text{IPS}}(V; \mathcal{D})) &= \mathbb{E}_{\mathcal{D}}[w(a, x)r] - V(\pi) \\
 &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)p(r|a,c,x)}[w(a, x)r] - V(\pi) \\
 &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)}[w(a, x)q(a, c)] - \mathbb{E}_{p(c)p(x|c)\pi(a|x)}[q(a, c)]
 \end{aligned}$$

$$\begin{aligned}
 &= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{U}(c, x, \pi_0)^c} \pi_0(a|x) w(a, x) q(a, c) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi(a|x) q(a, c) dx \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, x, \pi_0)^c} q(a, c) \int_{\mathcal{X}} p(x|c) \frac{\pi(a|x)}{\pi_0(a|x)} \pi_0(a|x) dx \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) \int_{\mathcal{X}} p(x|c) \pi(a|x) dx \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, x, \pi_0)^c} q(a, c) p(a|c, \pi) \right] - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a|c, \pi) \right] \\
 &= \mathbb{E}_{p(c)} \left[- \sum_{a \in \mathcal{U}(c, x, \pi_0)} p(a|c, \pi) q(a, c) \right]
 \end{aligned}$$

Since $q(a, c) \geq 0$ in the binary reward setting, it follows that $|\text{Bias}(\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D}))| = \mathbb{E}_{p(c)} [\sum_{\mathcal{U}(c, \pi_0)} p(a|\pi, c) q(a, c)]$ and $|\text{Bias}(\hat{V}_{\text{IPS}}(\pi; \mathcal{D}))| = \mathbb{E}_{p(c)} [\sum_{\mathcal{U}(c, x, \pi_0)} p(a|\pi, c) q(a, c)]$ and consequently we have that:

$$\begin{aligned}
 |\text{Bias}(\hat{V}_{\text{IPS}}(\pi; \mathcal{D}))| - |\text{Bias}(\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D}))| &= \mathbb{E}_{p(c)} \left[\sum_{\mathcal{U}(c, x, \pi_0)} p(a|\pi, c) q(a, c) \right] \\
 &\quad - \mathbb{E}_{p(c)} \left[\sum_{\mathcal{U}(c, \pi_0)} p(a|\pi, c) q(a, c) \right] \\
 &= \mathbb{E}_{p(c)} \left[\sum_{\mathcal{U}(c, x, \pi_0) \setminus \mathcal{U}(c, \pi_0)} p(a|\pi, c) q(a, c) \right]
 \end{aligned}$$

A.6 Proposition 3.10

Assuming that we have a set of embeddings $e \in \mathcal{E} \subset \mathbb{R}^{d_e}$ associated with the actions $a \in \mathcal{A}$ and an approximation $f_{\phi^*}(r)$ to the importance weights $w(a, x)$:

$$\begin{aligned}
 f_{\phi^*}(r) &:= \text{argmin}_{\phi} \mathbb{E}_{\phi} [(w(a, x) - f_{\phi}(r))^2] \\
 f_{\phi} &\in \{f_{\phi} : \mathbb{R} \rightarrow \mathbb{R} \mid \phi \in \Phi\}
 \end{aligned} \tag{19}$$

Then if we assume that $f_{\phi^*}(r) = w(a, x) + \epsilon$ for some $\epsilon \in \mathbb{R}$ we have that

$$\begin{aligned}
 &|\text{Bias}(\hat{V}_{\text{MR}}; \mathcal{D})| - |\text{Bias}(\hat{V}_{\text{CHIPS}}; \mathcal{D})| \\
 &= -\mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi_0)} p(a | \pi, c) q(a, c) \right] + \text{Bias}(\hat{V}_{\text{IPS}}; \mathcal{D}) + \mathbb{E}_{\mathcal{D}} [f_{\phi^*}(r)r] - \mathbb{E}_{\mathcal{D}} [w(a, x)] \\
 &= -\mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi)} p(a | \pi, c) q(a, c) \right] - V(\pi) + \mathbb{E}_{\mathcal{D}} [f_{\phi^*}(r)r] \\
 &= -\mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi_0)} p(a | \pi, c) q(a, c) \right] - V(\pi) + \mathbb{E}_{\mathcal{D}} [w(a, x)r] + \epsilon \mathbb{E}_{\mathcal{D}} [r] \\
 &= -\mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi)} p(a | \pi, c) q(a, c) \right] + \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(x, c, \pi_0)} q(a, c) \underbrace{\int_{x \in x} p(x | c) \pi(a | x) dx}_{p(a|\pi, c)} \right]
 \end{aligned}$$

$$\begin{aligned}
 & -\mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) \underbrace{\int_{x \in \mathcal{X}} p(x | c) \pi(a | x) dx}_{p(a | \pi, c)} \right] + \varepsilon \mathbb{E}_D[r] \\
 & = \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(x, c, \pi_0) \setminus \mathcal{U}(c, \pi_0)} q(a, c) p(a | \pi, c) \right] - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a | \pi, c) \right] + \varepsilon \mathbb{E}_D[r] \\
 & = -\mathbb{E}_{p(c)} \left[\sum_{a \in (\mathcal{U}(x, c, \pi_0) \setminus \mathcal{U}(c, \pi_0))^c} q(a, c) p(a | \pi, c) \right] + \varepsilon \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} q(a, c) p(a | \pi_0, c) \right]
 \end{aligned}$$

for the MIPS case, we note that MIPS' bias can also be expressed similarly to CHIPS':

$$\begin{aligned}
 \text{Bias}(\hat{V}_{\text{MIPS}}; D) & = \\
 & = \mathbb{E}_D[w(x, e)r] - V(\pi) \\
 & = \mathbb{E}_{p(x)\pi_0(a|x)p(e|x,a)p(r|x,a,e)}[w(x, e)r] - V(\pi) \\
 & = \mathbb{E}_{p(x)} \left[\sum_{a \in \mathcal{A}} \pi_0(a | x) \sum_{e \in \mathcal{E}} p(e | x, a) w(x, e) q(x, e) \right] \\
 & \quad - \mathbb{E}_{p(x)} \left[\sum_{a \in \mathcal{A}} \pi(a | x) \sum_{e \in \mathcal{E}} p(e | x, a) w(x, e) q(x, e) \right] \\
 & = \mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{E}} q(x, e) \left(\sum_{a \in \mathcal{A}} \pi_0(a | x), p(e | x, a) \right) \right] \\
 & \quad - \mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{E}} q(x, e) \left(\sum_{a \in \mathcal{A}} \pi(a | x) p(e | x, a) \right) \right] \\
 & = \mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{U}(e, \pi_0)^c} p(e | x, \pi_0) q(x, e) \frac{p(e | x, \pi)}{p(e | x, \pi_0)} \right] \\
 & \quad - \mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{E}} q(x, e) p(e | x, \pi) \right] \\
 & = -\mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{U}(e, \pi_0)} p(e | x, \pi) q(x, e) \right]
 \end{aligned}$$

Therefore the difference in bias is:

$$\begin{aligned}
 & |\text{Bias}(\hat{V}_{\text{MIPS}}; \mathcal{D})| - |\text{Bias}(\hat{V}_{\text{CHIPS}}; \mathcal{D})| \\
 & = \mathbb{E}_{p(x)} \left[\sum_{e \in \mathcal{U}(e, \pi_0)} p(e | x, \pi) q(x, e) \right] - \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{U}(c, \pi_0)} p(a | c, \pi) q(a, c) \right]
 \end{aligned}$$

A.7 Proposition 3.11

Lemma A.1. *Given a policy π , under Assumption 3.1 we have the transformation:*

$$w(a, c) = \mathbb{E}_{\pi_0(x|a,c)}[w(a, x)]$$

Proof:

Given a logging policy π_0 and an evaluation policy π , in the cluster setting of the bandits problem we have that:

$$w(a, c) = \frac{p(a | \pi, c)}{p(a | \pi_0, c)}$$

$$= \frac{\int_{\mathcal{X}} \pi(a|x) p(x|c)}{p(a|\pi_0, c)} \quad (20)$$

$$= \frac{\overbrace{p(a|c, \pi_0)}^{\frac{\pi(a|x)}{\pi_0(a|x)}} \int_{\mathcal{X}} \frac{\pi(a|x)}{\pi_0(a|x)} \pi_0(x|a, c)}{\overbrace{p(a|\pi_0, c)}^{\frac{\pi(a|x)}{\pi_0(a|x)}}} \quad (21)$$

$$= \mathbb{E}_{\pi_0(x|a, c)} [w(a, x)]$$

Where we have used the definition $p(a|\pi, c) = \int_{\mathcal{X}} \pi(a|x) p(x|c)$ in Equation 20, and that $\pi_0(x|a, c) = \frac{p(x|c) \pi_0(a|x)}{p(a|c, \pi_0)}$ in Equation 21.

Given a logging policy π_0 and an evaluation policy π , under Assumption 3.1 and Assumption 3.2 we have that

$$\begin{aligned} & N (\mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{IPS}}(\pi; \mathcal{D})] - \mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D})]) \\ &= N \left(\mathbb{V}_{\mathcal{D}} \left[\frac{1}{N} \sum_{i=1}^N \frac{\pi(a_i|x_i)}{\pi_0(a_i|x_i)} r_i \right] - \mathbb{V}_{\mathcal{D}} \left[\frac{1}{N} \sum_{i=1}^N \frac{p(a_i|c_i, \pi)}{p(a_i|c_i, \pi_0)} r_i \right] \right) \\ &= \mathbb{V}_{\mathcal{D}} \left[\frac{\pi(a|x)}{\pi_0(a|x)} r \right] - \mathbb{V}_{\mathcal{D}} \left[\frac{p(a|c, \pi)}{p(a|c, \pi_0)} r \right] \end{aligned} \quad (22)$$

$$= \left(\mathbb{E}_{\mathcal{D}} [w(a, x)^2 r^2] - \underbrace{\mathbb{E}_{\mathcal{D}} [w(a, x) r]^2}_{V(\pi)} \right) - \left(\mathbb{E}_{\mathcal{D}} [w(a, c)^2 r^2] - \underbrace{\mathbb{E}_{\mathcal{D}} [w(a, c) r]^2}_{V(\pi)} \right) \quad (23)$$

$$\begin{aligned} &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)} [w(a, x)^2 \mathbb{E}_{p(r|a, c, x)} [r^2]] - \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)} [w(a, c)^2 \mathbb{E}_{p(r|a, c, x)} [r^2]] \\ &= \mathbb{E}_{p(x)p(c|x)\pi_0(a|x)} [(w(a, x)^2 - w(a, c)^2) \mathbb{E}_{p(r|a, c, x)} [r^2]] \\ &= \mathbb{E}_{p(c)p(x|c)\pi_0(a|x)} [(w(a, x)^2 - w(a, c)^2) \mathbb{E}_{p(r|a, c, x)} [r^2]] \\ &= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x|c) \sum_{a \in \mathcal{A}} \pi_0(a|x) (w(a, x)^2 - w(a, c)^2) \mathbb{E}_{p(r|a, c, x)} [r^2] dx \right] \\ &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} \int_{\mathcal{X}} p(x|c) \pi_0(a|x) (w(a, x)^2 - w(a, c)^2) \mathbb{E}_{p(r|a, c, x)} [r^2] dx \right] \\ &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} \int_{\mathcal{X}} \frac{\pi_0(x|a, c) p(a|c, \pi_0)}{\overbrace{\pi_0(a|x)}^{\frac{\pi(a|x)}{\pi_0(a|x)}}} \overbrace{\pi_0(a|x)}^{\frac{\pi(a|x)}{\pi_0(a|x)}} (w(a, x)^2 - w(a, c)^2) \mathbb{E}_{p(r|a, c, x)} [r^2] dx \right] \quad (24) \\ &= \mathbb{E}_{p(c)} \left[\sum_{a \in \mathcal{A}} p(a|c, \pi_0) \int_{\mathcal{X}} \pi_0(x|a, c) (w(a, x)^2 - w(a, c)^2) \mathbb{E}_{p(r|a, c, x)} [r^2] dx \right] \\ &= \mathbb{E}_{p(c)p(a|c, \pi_0)} \left[\left(\int_{\mathcal{X}} \pi_0(x|a, c) w(a, x)^2 dx - w(a, c)^2 \int_{\mathcal{X}} \overbrace{\pi_0(x|a, c)}^{\frac{\pi(x|a, c)}{\pi_0(a|x, c)}} dx \right) \mathbb{E}_{p(r|a, c, x)} [r^2] \right] \\ &= \mathbb{E}_{p(c)p(a|c, \pi_0)} \left[(\mathbb{E}_{\pi_0(x|a, c)} [w(a, x)^2] - \mathbb{E}_{\pi_0(x|a, c)} [w(a, x)]^2) \mathbb{E}_{p(r|a, c, x)} [r^2] \right] \quad (25) \\ &= \mathbb{E}_{p(c)p(a|c, \pi_0)} [\mathbb{V}_{\pi_0(x|a, c)} [w(a, x)] \mathbb{E}_{p(r|a, c, x)} [r^2]] \geq 0 \end{aligned}$$

Note in Equation 22 we used that the samples in \mathcal{D} are i.i.d, in particular the linearity of variance under this condition. The cancellation of terms in Equation 23 results from IPS and CHIPS being unbiased under Assumptions 3.1 and 3.2. In Equation 24 we used that $\pi_0(x|a, c) = \frac{p(x|c) \pi_0(a|x)}{p(a|c, \pi_0)}$, while Equation 25 uses Lemma A.1.

A.8 Proposition 3.12

The first thing we need to note is that CHIPS and MIPS are in different spaces regarding the contextual bandits generating process. MIPS assumes the existence of an action embedding space $e \in \mathcal{E} \subseteq \mathbb{R}^{d_e}$ and CHIPS assumes the existence of a partition of the context space $\mathcal{C} := \{\mathcal{C}_i\}_{i=1}^K$ with $\mathcal{C}_i \subset \mathcal{X}$ and $c_i \cap c_j = \emptyset$. For joining this spaces, we assume that given a policy π , at every iteration of the data generation process, apart from the classical context ($x \in \mathcal{X}$), action ($a \in \mathcal{A}$) and reward ($r \in [0, r_{max}] \subset \mathbb{R}$), we observe a cluster $c \sim p(c|x)$ and an action embedding $e \sim p(e|a, c, x)$. Given a policy π the policy value $V(\pi)$ equation can be then refined to:

$$V(\pi) := \mathbb{E}_{p(x)p(c|x)\pi(a|x)p(e|a, c, x)p(r|e, a, c, x)} [r]$$

$$= \mathbb{E}_{p(x)p(c|x)\pi(a|x)p(e|a,c,x)q(e,a,c,x)}$$

Here $q(e, a, c, x) := \mathbb{E}_{p(r|e,a,c,x)}[r]$. Note that as in MIPS and CHIPS case, the refinement does not contradict the classical policy value definition.

We also need to refine $p(a | c, \pi)$ (from CHIPS) and $p(e | a, \pi)$ (from MIPS) in the joint space:

$$\begin{aligned} p(a | c, \pi) &= \sum_{e \in \mathcal{E}} \int_{\mathcal{X}} p(e | a, c, x) p(x | c) \pi(a | x) \\ p(e | x, \pi) &:= \sum_{c \in \mathcal{C}} \sum_{a \in \mathcal{A}} p(e | a, c, x) p(c | x) \pi(a | x) \end{aligned}$$

It is important to note that after joining the context space, to make a fair comparison between MIPS and CHIPS, there are some dependencies that we want to eliminate to prevent information from passing between variables that were not originally in the definition of MIPS and CHIPS. In particular, for CHIPS, we eliminate the dependency of the cluster with respect to the embedding given the context and the action (i.e., $c \perp e | (x, a)$), and for MIPS, the dependency of the action with respect to the cluster given the embedding and the context (i.e., $a \perp c | (x, e)$). From Proposition 3.11 we know that $\mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{IPS}}(\pi; \mathcal{D})] \geq \mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D})]$ and from MIPS Theorem 3.6 we know that $\mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{IPS}}(\pi; \mathcal{D})] \geq \mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{MIPS}}(\pi; \mathcal{D})]$. Therefore, we need to make a comparison between $\mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{MIPS}}(\pi; \mathcal{D})]$ and $\mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D})]$.

To follow the structure of Proposition 3.11, we are going to assume that Assumptions 3.1 and 3.2 hold as well as their counterparts from MIPS. The following identities hold under these conditions:

$$\begin{aligned} p(x | c) \pi(a | x) &= \frac{p(e | x, \pi) p(c | x, a) \pi_0(a | x) p(x)}{p(e | x, \pi_0) p(c)} \\ p(e | x, a, c) &= \frac{p(x | e, a, c) p(e | a, c) p(a | c, \pi_0)}{p(c | x, a) \pi_0(a | x) p(x)} \end{aligned}$$

Now, under these conditions, we need a relation between the weights of MIPS and CHIPS:

$$\begin{aligned} \omega(a, c)^2 &= \frac{p(a | c, \pi)}{p(a | c, \pi_0)} \\ &= \frac{\int_{\mathcal{X}} p(x | c) \sum_{e \in \mathcal{E}} \pi(a | x) p(e | c, a, x)}{p(a | c, \pi_0)} \\ &= \frac{\int_{\mathcal{X}} \sum_{e \in \mathcal{E}} w(e, x) p(x | e, a, c) p(e | a, c) p(a | c, \pi_0)}{p(a | c, \pi_0)} \\ &= \sum_{e \in \mathcal{E}} p(e | a, c) \int_{\mathcal{X}} p(x | e, a, c) \omega(e, x) \\ &= \mathbb{E}_{p(e|a,c)p(x|e,a,c)}[\omega(e, x)] \end{aligned}$$

Therefore the scaled difference in variance can be expressed as:

$$\begin{aligned} &N (\mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{MIPS}}(\pi; \mathcal{D})] - \mathbb{V}_{\mathcal{D}} [\hat{V}_{\text{CHIPS}}(\pi; \mathcal{D})]) \\ &= N \left(\mathbb{V}_{\mathcal{D}} \left[\frac{1}{N} \sum_{i=1}^N \frac{\pi(e_i | x_i)}{\pi_0(e_i | x_i)} r_i \right] - \mathbb{V}_{\mathcal{D}} \left[\frac{1}{N} \sum_{i=1}^N \frac{p(a_i | c_i, \pi)}{p(a_i | c_i, \pi_0)} r_i \right] \right) \\ &= \mathbb{V}_{\mathcal{D}} \left[\frac{\pi(e | x)}{\pi_0(e | x)} r \right] - \mathbb{V}_{\mathcal{D}} \left[\frac{p(a | c, \pi)}{p(a | c, \pi_0)} r \right] \\ &= (\mathbb{E}_{\mathcal{D}} [\omega(e, x)^2 r^2] - \underbrace{\mathbb{E}_{\mathcal{D}} [\omega(e, x) r]^2}_{V(\pi)}) - (\mathbb{E}_{\mathcal{D}} [\omega(a, c)^2 r^2] - \underbrace{\mathbb{E}_{\mathcal{D}} [\omega(a, c) r]^2}_{V(\pi)}) \\ &= \mathbb{E}_{p(c)} \left[\int_{\mathcal{X}} p(x | c) \sum_{a \in \mathcal{A}} \pi_0(a | x) \sum_{e \in \mathcal{E}} p(e | a, c, x) (\omega(e, x)^2 - \omega(a, c)^2) \mathbb{E}_{p(r|a,c)} [r^2] dx \right] \end{aligned}$$

$$\begin{aligned}
 &= \mathbb{E}_{p(c)} \left[\sum_{a \in A} p(a | c, \pi_0) \sum_{e \in \mathcal{E}} p(e | a, c) \int_{\mathcal{X}} p(x | e, a, c) (\omega(e, x)^2 - \omega^2(a, c)) \mathbb{E}_{p(r|a,c)} [r^2] dx \right] \\
 &= \mathbb{E}_{p(c)p(a|c, \pi_0)} \left[\mathbb{E}_{p(r|a,c)} [r^2] \left(\mathbb{E}_{p(e|a,c)p(x|e,a,c)^2} [\omega(e, x)^2] \right) \right. \\
 &\quad \left. - \mathbb{E}_{p(c)p(a|c, \pi_0)} \left[\mathbb{E}_{p(r|a,c)} [r^2] \left(\omega(a, c)^2 \sum_{e \in \mathcal{E}} p(e | a, c) \int_{\mathcal{X}} p(x | e, a, c) dx \right) \right] \right] \\
 &= \mathbb{E}_{p(c)p(a|c, \pi_0)} \left[\mathbb{E}_{p(r|a,c)} [r^2] \left(\mathbb{E}_{p(e|a,c)p(x|e,a,c)} [\omega(e, x)^2] - \mathbb{E}_{p(e|a,c)p(x|e,a,c)} [\omega(e, x)]^2 \right) \right] \\
 &= \mathbb{E}_{p(c)p(a|c, \pi_0)} \left[\mathbb{E}_{p(r|a,c)} [r^2] \mathbb{V}_{p(e|a,c)p(x|e,a,c)} [\omega(e, x)] \right] \geq 0
 \end{aligned}$$

This implies that under Assumptions 3.1 and 3.2 (and their counterparts in MIPS), the variance of CHIPS is lower than the variance of MIPS, proving the proposition.

B REWARD ESTIMATES DERIVATION

B.1 MAP

From the setting in Subsection 3.2 we denote $R_c := \{r_i\}_{i=1}^M$ as the rewards observed in cluster c from the logging data². We consider R_c as independent trials of a Bernoulli random variable with parameter θ (i.e., $R_c \stackrel{i.i.d.}{\sim} \text{Ber}(\theta)$). Therefore, we have that the likelihood can be expressed as:

$$\begin{aligned}
 p(R_c | \theta) &= \prod_{i=1}^M p(r_i | \theta) \\
 &= \prod_{i=1}^M \theta^{r_i} (1 - \theta)^{1-r_i} \\
 &= \theta^{\sum_{i=1}^M r_i} (1 - \theta)^{M - \sum_{i=1}^M r_i}
 \end{aligned}$$

Using a Beta distribution as a prior we have that:

$$p(\theta) = \text{Beta}(\theta | \alpha, \hat{\beta}) = \frac{1}{\mathcal{B}(\alpha, \hat{\beta})} \theta^{\alpha-1} (1 - \theta)^{\hat{\beta}-1}$$

Where $\mathcal{B}(\alpha, \hat{\beta}) = \frac{\Gamma(\alpha)\Gamma(\hat{\beta})}{\Gamma(\alpha+\hat{\beta})}$ and $\Gamma(\cdot)$ is the Gamma function. The posterior probability can then be expressed as:

$$\begin{aligned}
 p(\theta | R_c) &\propto p(R_c | \theta) p(\theta) \\
 &\propto \theta^{\sum_{i=1}^M r_i} (1 - \theta)^{M - \sum_{i=1}^M r_i} \frac{1}{\mathcal{B}(\alpha, \hat{\beta})} \theta^{\alpha-1} (1 - \theta)^{\hat{\beta}-1} \\
 &\propto \theta^{\alpha-1 + \sum_{i=1}^M r_i} (1 - \theta)^{\hat{\beta}-1 + M - \sum_{i=1}^M r_i} \\
 &\propto \text{Beta} \left(\theta \mid \alpha + \sum_{i=1}^M r_i, \hat{\beta} + M - \sum_{i=1}^M r_i \right)
 \end{aligned}$$

The MAP estimator of θ is the mode of the resulting Beta distribution, i.e.

$$\hat{\theta}_{\text{MAP}} = \frac{(\alpha - 1) + \sum_{i=1}^M r_i}{\alpha + \hat{\beta} + M - 2}$$

B.2 ML

Using the same setting as in the previous section ($R_c \stackrel{i.i.d.}{\sim} \text{Ber}(\theta)$) we have that the maximum likelihood estimation can be expressed as

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in \Theta} \left\{ \prod_{i=1}^M \theta^{r_i} (1 - \theta)^{1-r_i} \right\}$$

²Here we refer to the already transformed version using clusters. See definition of τ in Subsection 3.2

$$= \arg \max_{\theta \in \Theta} \left\{ \underbrace{\log(\theta) \cdot \sum_{i=1}^M r_i + \log((1-\theta)) \cdot \sum_{i=1}^M (1-r_i)}_{l(\theta)} \right\}$$

We now search for local maxima by setting the differential to 0:

$$\begin{aligned} \frac{\partial l(\theta)}{\partial \theta} = 0 &\implies \frac{\sum_{i=1}^M r_i}{\theta} + \frac{\sum_{i=1}^M (1-r_i)}{(1-\theta)} = 0 \\ &\implies \sum_{i=1}^M r_i - \theta \sum_{i=1}^M r_i = \theta \sum_{i=1}^M (1-r_i) \\ &\implies \hat{\theta}_{\text{ML}} = \frac{1}{M} \sum_{i=1}^M r_i \end{aligned}$$

C EXPERIMENTAL PARAMETERS AND HARDWARE

Parameter	Value	Description
c_{exp}	10	Radius of the n-dimensional ball for context space generation.
c_{rad}	1	Cluster generation radius.
d_x	2	Dimension of context vectors.
x_{num}	1.000	No. of different context vectors in the experiment.
a_{num}	10	No. of actions in the experiment.
c_{num}	10	No. of clusters in the experiment.
$n_{samples}$	50.000	No. of logged samples to use in the experiment.
$emp_{c_{num}}$	100	No. of clusters to use empirically by the clustering method.
e_{len}	1.000.000	No. of samples extracted from the dataset for the evaluation policy
b_{len}	1.000.000	No. of samples extracted from the dataset for the evaluation policy
σ	0.2	Context-specific behaviour deviation from cluster behaviour.
β	-1	Deviation between evaluation and logging policies.
α	20	Parameter from beta distribution in Bayesian inference
$\hat{\beta}$	20	Parameter from beta distribution in Bayesian inference

Table 1: Parameters used in the basic configuration for experiments for generation and estimation.

CPU	AMD Ryzen Threadripper PRO 3975WX
RAM	256 GB
Cores	64
GPU	2x Nvidia A100 160GB

Table 2: Specifications of the machine in which the experiments were executed.

D ADDITIONAL EXPERIMENTS

D.1 Synthetic Experiments

Number of actions. From the fixed basic configuration that uses 100 clusters for CHIPS’ estimates, we observe a progressive deterioration in the estimator capabilities when increasing the number of actions (see Figure 4). We theorize that this behaviour might be a consequence of the violation of Assumption 3.1 when trying to group contexts using an excessive number of clusters in a large action space, resulting in deficient actions inside the clusters. This problem can be mitigated by decreasing the number of clusters used in the clustering method for the CHIPS estimation (see Figure 12 (left)).

Number of samples. We observe an approximation to the performance of IPS as we increase the number of samples in the logged data that we identify as an effect of reducing the number of observed deficient action-context pairs in IPS,

converging to an unbiased estimator under Assumption 2.1 (see Figure 1 (right)). In this case, the clustering effects under CHIPS become less noticeable according to Corollary 3.9 since $\mathcal{U}(c, x, \pi_0) \setminus \mathcal{U}(c, \pi_0) \rightarrow \emptyset$. It is worth mentioning that increasing the number of clusters when enough samples are available, as well as reducing it in the opposite case, can improve the performance of the CHIPS estimates, as shown in Figure 12 (right).

Cluster radius. Increasing the cluster radius in the generation process affects the separability of the cluster space and complicates the partitioning in clusters complying with Assumption 3.2. In this case, we could find significant differences in context behaviour for both actions and rewards within a cluster, resulting in increased bias from the empirical approximations. Therefore, we observe a convergence to IPS’ performance as cluster radius increases since the context space becomes less separable (see Figure 6).

Sigma. Increasing context-specific noise in the generation process produces a similar effect as in the cluster radius case. In particular, the larger the noise, the more common it is to observe inconsistent behaviour in actions and rewards for contexts within a cluster, complicating the approximation of a homogeneous cluster-wise behaviour and resulting again in a bias increase (see Figure 8).

Alpha (prior). In this experiment, we vary the alpha parameter of the Beta prior maintaining all other settings fixed. Like in the number of clusters case, we observe a similar v-shaped graph indicating that, as expected from the previous β analysis (see Section 4.1.1), the CHIPS (MAP) estimator is sensitive to the prior. In particular, lower values push the expected reward of each cluster to the ML’s estimate, while higher values push it to the prior’s expected value, decreasing performance in both cases (see Figure 1 center). For different values of distributional shift (β), the optimal value will depend on the *resistance* MAP offers to converge to the ML estimate, favouring lower values as β becomes larger (see Figure 13).

Clustering Method. In this experiment, we evaluate the performance of the CHIPS (MAP) estimator using different clustering methods while varying the clustering radius in the synthetic generation process. In Figure 10, we observe that using Mean Shift (Comaniciu and Meer, 2002) or Bayesian Gaussian Mixture (Bishop, 2006; Attias, 1999; Blei and Jordan, 2006) fails to separate the context space resulting in the same performance as IPS. DBSCAN (Ester et al., 1996) mitigates IPS’ increase in mean squared error when the context space is easier to separate (i.e., lower radii values) but converges to IPS when the context is complicated to separate (i.e., higher radii values). Affinity Propagation (Frey and Dueck, 2007) follows a similar behaviour to DBSCAN but still offers some improvement with respect to IPS when the space is difficult to separate. OPTICS (Ankerst et al., 1999) makes a general improvement to the Affinity Propagation performance, especially noticeable when the context space is separable. MiniBatch K-Means (Sculley, 2010), Gaussian Mixture (Everitt, 1996), Birch (Zhang, Ramakrishnan, and Livny, 1996), Spectral Clustering (Shi and Malik, 2000), and Agglomerative Clustering (Ward, 1963) have similar performance, outperforming Affinity Propagation for the separable case. We also note a general upward tendency in mean squared error for every clustering method as the space becomes more complicated to separate.

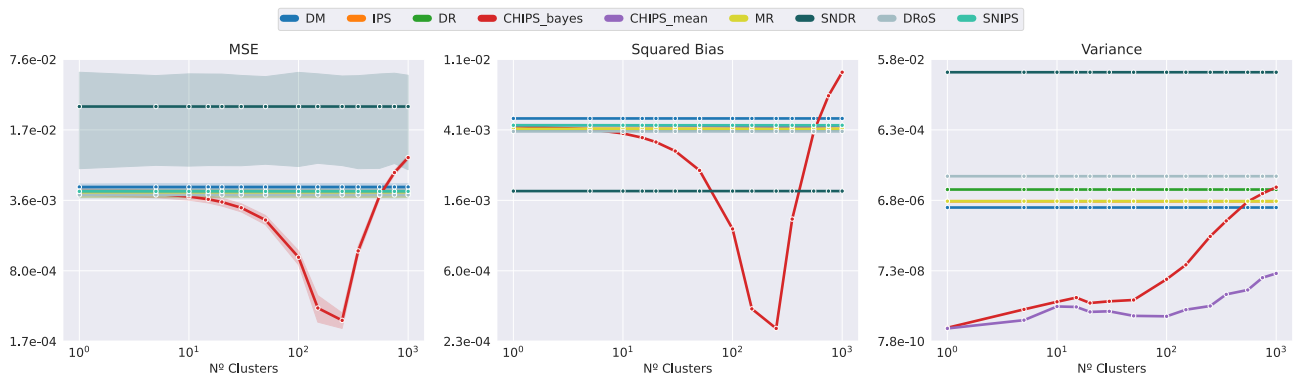


Figure 3: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the number of clusters.

³In this case we used a slightly different version of the configuration settings to make a more challenging environment in which we use 10.000 samples and consequently reduce the number of empirical cluster estimation to 30 to easily assess the role that similarity of logging and evaluation policies play in CHIPS capabilities.

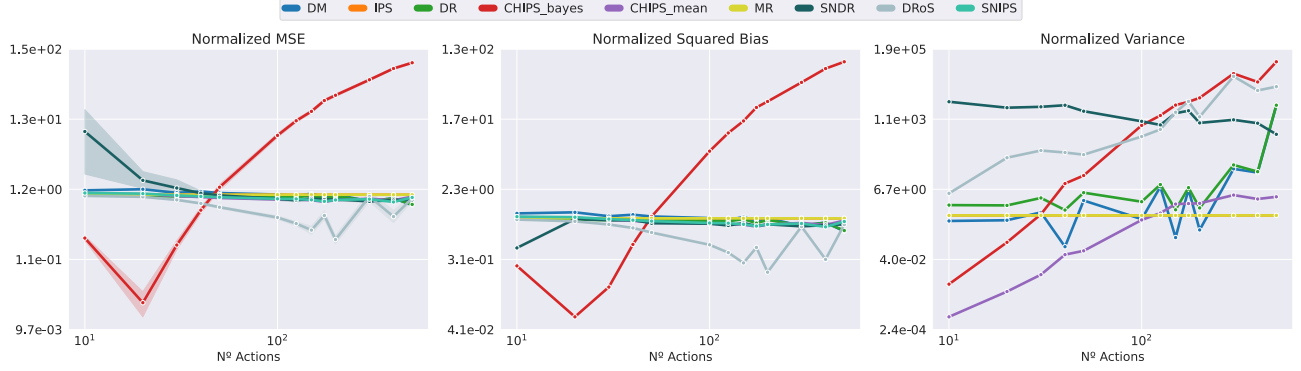


Figure 4: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the number of actions.

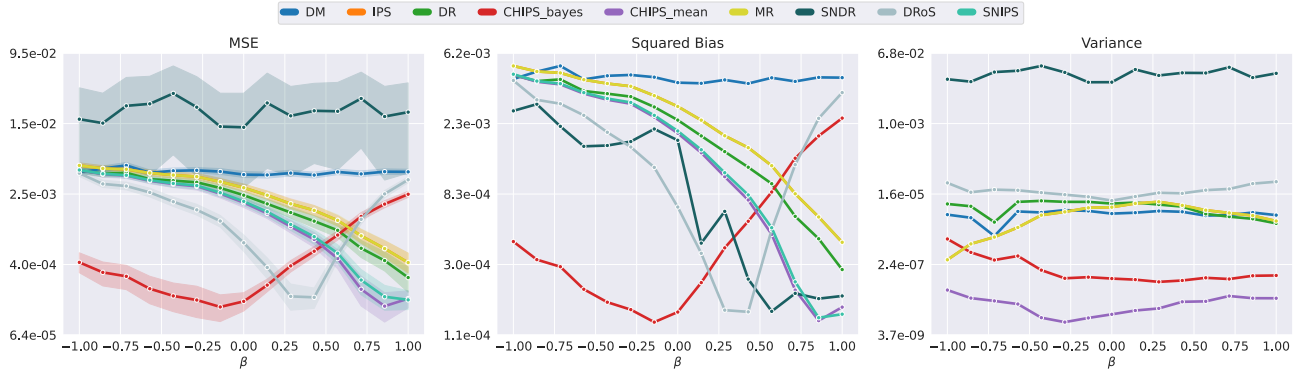


Figure 5: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying β values.³

D.1.1 Bi-parametric variations

The experiments varying single parameters described in the previous section indicate that increasing the number of actions in a fixed configuration progressively deteriorates CHIPS’ performance. This behaviour is expected since the larger the action space, the more likely it is to incur in a situation in which Assumption 3.1 does not hold with a fixed number of clusters. In this situation, we found that reducing the number of clusters can mitigate the performance decay by pooling information from broader contexts clusters while increasing it could be beneficial in reduced action spaces (see Figure 12 (a)). Similarly, the number of samples from the logging policy also conditions how significant the performance gap between CHIPS and IPS is. In particular, the higher the number of samples, the more beneficial it is to use a higher number of clusters to try to obtain a more detailed partition structure of the context space, while a reduced number of clusters has an edge on few-sample cases (see Figure 12 (b)).

We also study the effect of varying the α parameter in CHIPS’ (MAP) Beta prior, using different values of the distributional shift between policies (β). In Figure 13, we observe that mid values of α (30-50) offer better performance when there is a considerable distributional shift between logging and evaluation policy (i.e., $\beta \approx -1$) since the expected reward per cluster is pushed towards the prior’s expectation, creating some resistance from converging to the average observed rewards (i.e., mitigating the reward misspecification existing under this conditions). As the distributional gap closes, lower values of α are more favourable since the samples observed per cluster are better representatives of the real expected reward. However, higher values for α (80-100) result in excessive resistance that deteriorates CHIPS’ performance. It is also worth mentioning that as the distributional gap closes, CHIPS (MAP) loses its advantage with respect to IPS since the logging and evaluation policies are closer, and the ML estimates would offer better results, as previously shown in Figure 1.

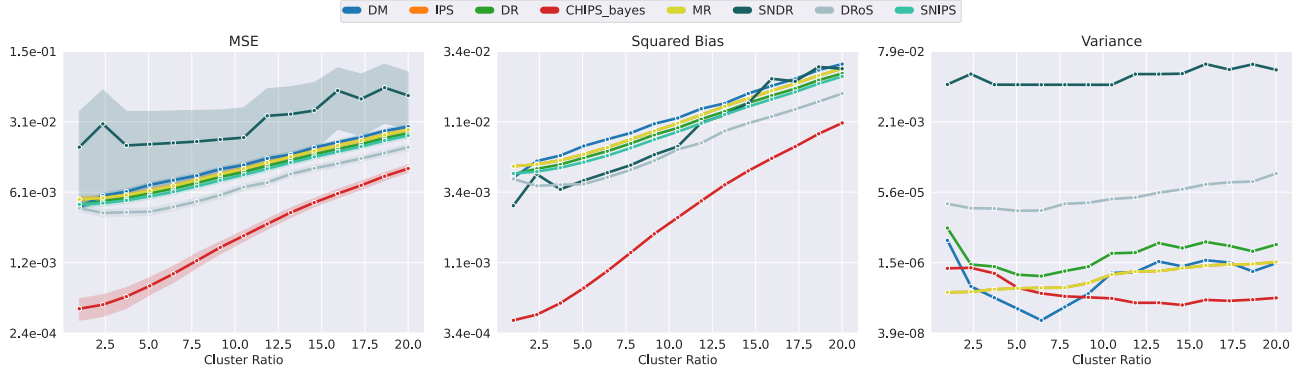


Figure 6: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the radius of the clusters generated.

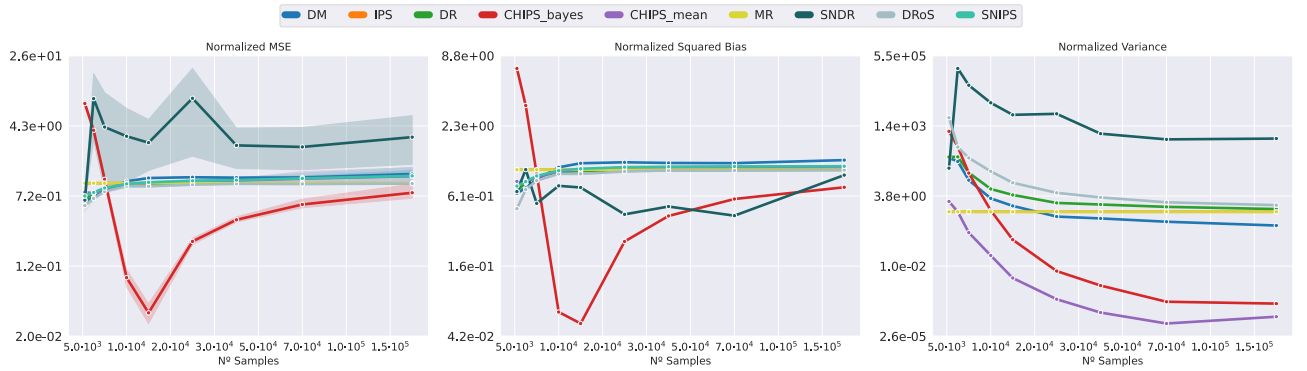


Figure 7: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the number of samples provided from the logging policy.

D.2 Real Experiments

D.3 MAP vs ML

In this section, we analyze the reason behind the jump in performance using the CHIPS estimator with the MAP estimate for the expected reward per cluster. For this purpose, we have conducted two experiments, one in the synthetic dataset and the other in the real dataset. For the synthetic experiment, given a distributional shift value β , we select the most relevant context-action pair (x^*, a^*) under the evaluation policy π (i.e., $(x^*, a^*) = \arg \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \pi(a|x)$). Then we analyze the mean squared error of the expected reward (given x^*) estimations made by CHIPS (MAP) (i.e., $w(a^*, c^*) \hat{r}_{\text{bayes}}(a^*, c^*)$) and CHIPS (ML) (i.e., $w(a^*, c^*) \hat{r}_{\text{mean}}(a^*, c^*)$) w.r.t IPS (where c^* is the cluster associated with x^*). We also compute the number of observations in c^* in which action a^* was selected. This process is repeated 100 times with different policies generated under different random seeds, and the results for the number of samples per cluster and squared errors are averaged. We repeat this for ten different values of β ranging from -1 to 1 and represent the moving averages for relative squared errors and samples in Figure 16. We observe that the number of samples per cluster increases with β as both policies become closer. This increase in the number of samples makes the ML estimates progressively more accurate since the extra samples push the estimated expected value to the real expected value. For lower values of β , when the gap between policies is more significant, although some samples are available in the cluster, the values for the rewards observed on them are non-informative of the real expected value (hence the difficulty of ML to make an accurate estimation and the difference between MAP and ML for misspecified reward settings as depicted in Figure 1). For the real dataset, we follow a similar procedure, but instead of the most relevant context-action pair, we select the top 15 and compare the conditional expected reward estimates MSE with respect to IPS' (see Figure 15). Since the logging policy for this dataset is uniform, the distributional shift between the logging and evaluation policies is not as significant as the one presented in the base configuration of the synthetic dataset ($\beta = -1$). In practice, this means that the CHIPS estimation of the expected reward per cluster using ML is more accurate than in the synthetic dataset but still far from the performance jump of the CHIPS

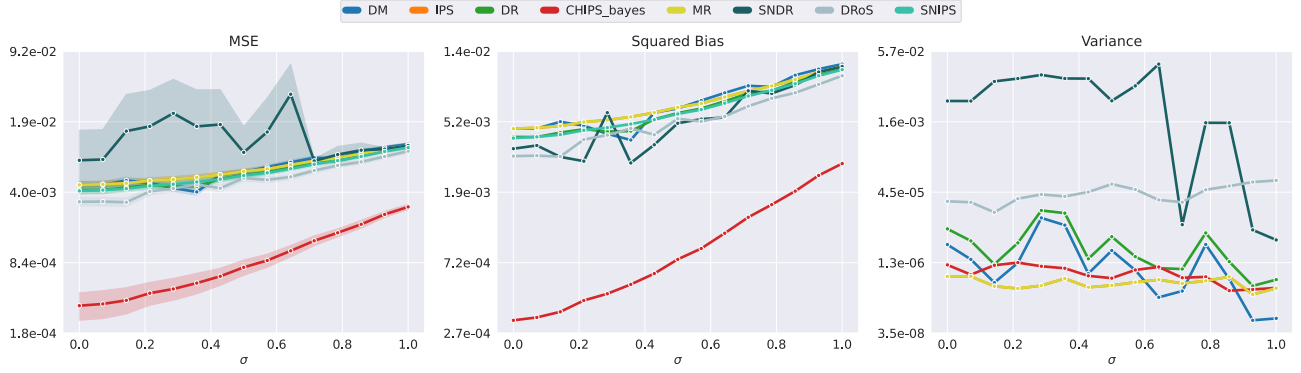


Figure 8: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the context-specific noise σ .

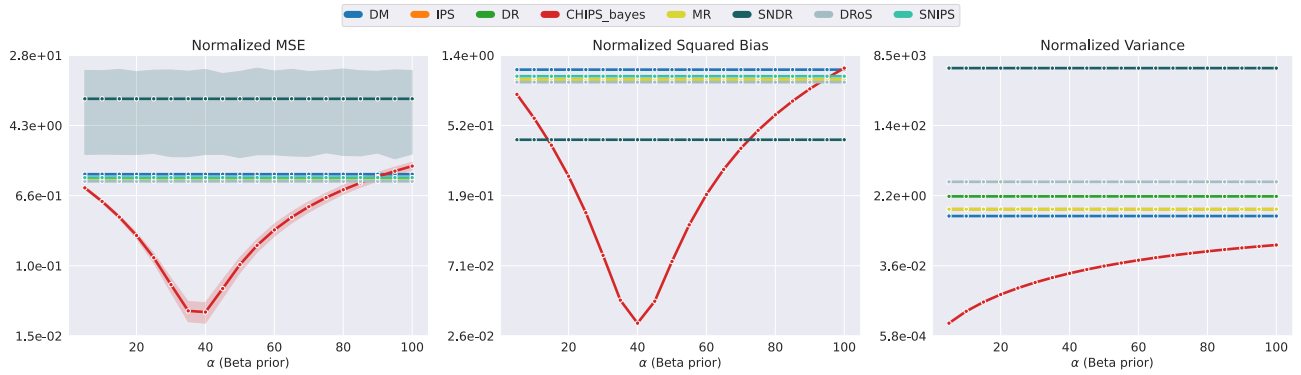


Figure 9: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the α parameter.

estimate using MAP, as we would expect from the results in Figure 2.

D.4 Choosing alpha in Arbitrary Problems

In Figure 14 (b), we observe that the hyperparameters of the MAP estimation process can heavily impact the performance of the method. As previously discussed in Appendix D.1.1, MAP hyperparameters control the resistance with which the expected reward per cluster is *pulled* towards the prior’s expectation. This resistance is particularly noticeable in smaller size clusters, in which estimating a reward based on observations alone is much more challenging. Since in these clusters the partitioning method cannot ensure high homogeneity at reward level, in our experimentation we decided to use a non-informative prior (i.e., $\alpha = \hat{\beta}$), to mitigate possible violations of Assumption 3.2 and reward misspecification. Intuitively, an optimal value for α under these conditions needs to balance the prior’s resistance to prevent reward misspecification without incurring into creating a quasi-uniform reward estimation (excessively large values of α). In Figure 17 we explore the optimal value of α for a given average number of datapoints per cluster-action. As expected, for small size clusters, lower values of *alpha* are favoured since the pull towards the prior’s expectation is soft, while on bigger clusters, the value of α (and consequently the resistance) needs to grow to effectively control reward misspecification (otherwise the expected reward value would be pulled towards the value of the observed samples).

To choose the value of α in an arbitrary problem, we propose the following selection process:

1. Determine the number of clusters to use depending on the number of clusters (reference in Figure 12 (a)).
2. Partition the context space \mathcal{X} in clusters c_1, c_2, \dots, c_n .
3. Generate synthetic data \hat{X}_{ev} using \mathcal{X}_{train} and π_e .

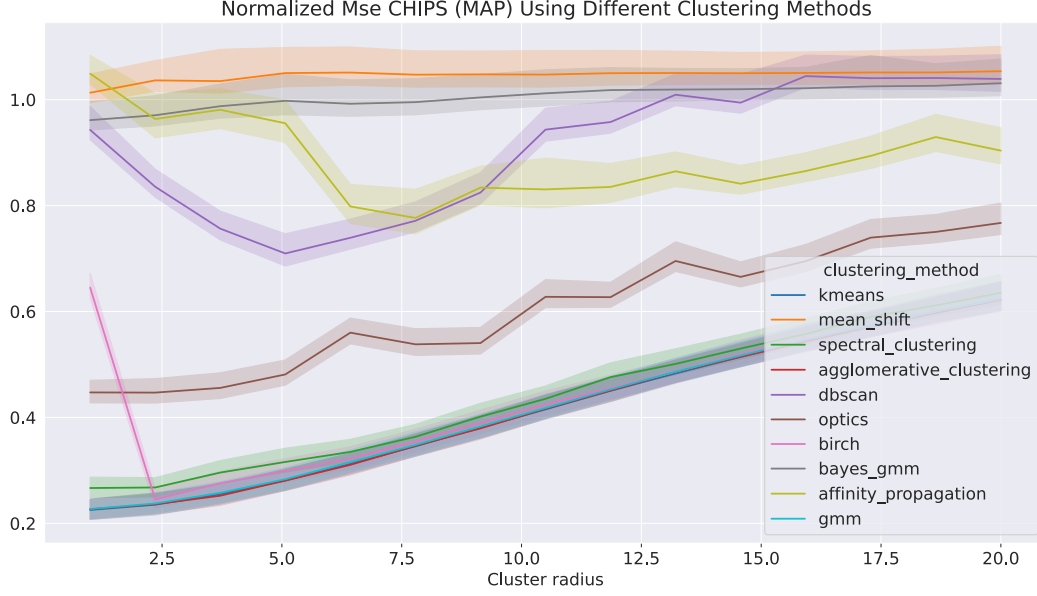


Figure 10: Normalized MSE of CHIPS (MAP) using different clustering methods with respect to IPS.

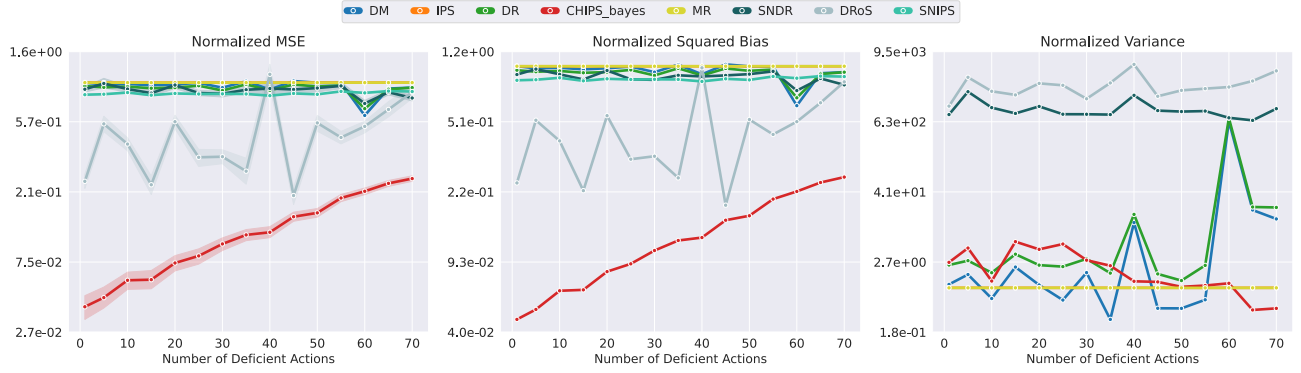


Figure 11: From left to right, MSE, Bias, and Variance of the CHIPS estimator compared to baselines while varying the number of deficient actions.

4. Estimate number of average data points per cluster-action from \hat{X}_{ev} .

5. Choose α from the reference in Figure 17.

For testing this selection process, following the experimental protocol of Taufiq et al. (2023), we transform five UCI datasets (Dua and Graff, 2017), MNIST (Deng, 2012), and CIFAR-100 (Krizhevsky, Nair, and Hinton, 2009) from multi-class classification problems into contextual bandits data (Dudík, Langford, and Li, 2011). The results (averaged 50 times) in Figure 18 show a consistent improvement with respect to existing methods, empirically proving the effectiveness of the α selection process.

Additionally, we perform an alternative experiment using the real dataset, in which instead of fixing α and vary the number of clusters according to the reference in Figure 12 (b) with 50000, 100000 and 500000 samples (see Figure 2, we follow the α selection process, fix the number of clusters and increase the value of α according to Figure 17. In Figure 19 we observe equivalent results as in our previous experiment confirming the equivalence of using a reference for the number of samples and varying the number of clusters with a fixed value for α , or varying α with a fixed number of clusters obtained by using a reference for the number of actions.

E CLUSTER STRUCTURE IN DATASETS

The generated synthetic dataset ensures that the expected reward inside a cluster is similar and that the best possible action is usually the same for all the context within the cluster (see Figure 20), mimicking real-world settings like e-commerce in which we can expect similar behaviour for close contexts.

F EXPERIMENTAL PROTOCOL

For evaluation in the real dataset, we follow (Saito and Joachims, 2022) protocol to evaluate estimators' accuracy given two sources of data. Given a logging policy π , a dataset collected under it \mathcal{D} , a logging policy π_0 , and the dataset collected under it \mathcal{D}_0 , we follow the following procedure:

1. Extract n independent bootstrap samples with replacement from the logging dataset $\mathcal{D}_0^* := \{(x_i, a_i, r_i)\}_{i=1}^n$.
2. Estimate the policy value of π using the sample \mathcal{D}_0^* . We denote this estimate as $\hat{V}(\pi; \mathcal{D}_0^*)$.
3. Compute the relative mean squared error with respect to IPS:

$$\mathcal{Z}(\hat{V}, \mathcal{D}_0^*) = \frac{(V(\pi) - \hat{V}(\pi; \mathcal{D}_0^*))^2}{(V(\pi) - \hat{V}_{\text{IPS}}(\pi; \mathcal{D}_0^*))^2}$$

Where $V(\pi) := \frac{1}{|\mathcal{D}|} \sum_{(\cdot, \cdot, r_i) \in \mathcal{D}} r_i$.

4. Repeat steps 1,2, and 3 $T = 100$ times and compute the Empirical Cumulative Distribution Function (ECDF) as:

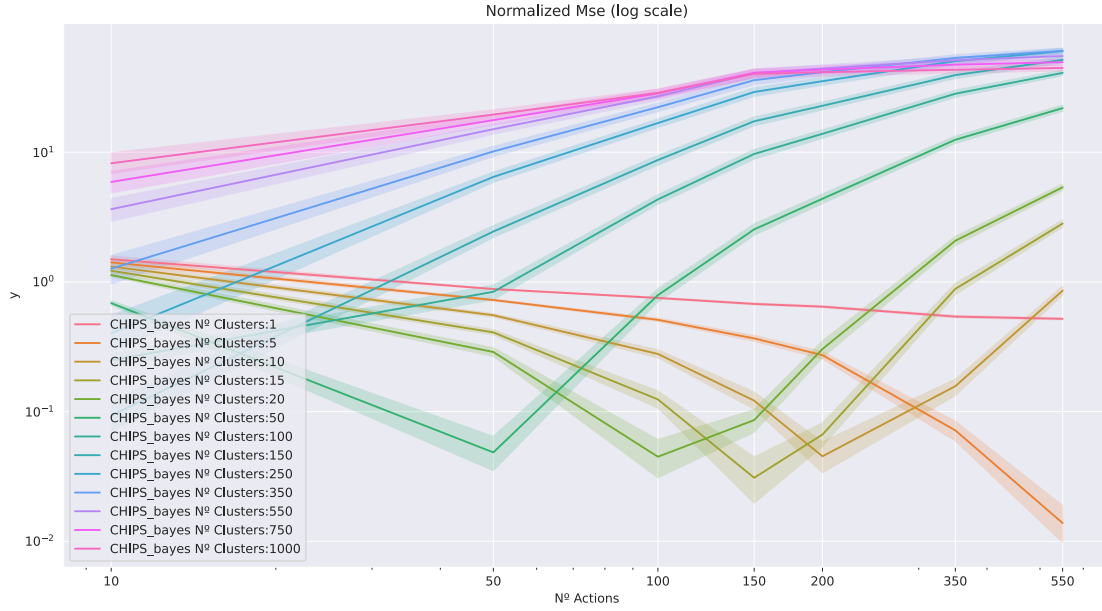
$$\hat{F}_{\mathcal{Z}}(x) := \frac{1}{T} \sum_{t=1}^T \mathbb{1}\{\mathcal{Z}_t(\hat{V}, \mathcal{D}_0^*) \leq x\}$$

G COMPLEXITY

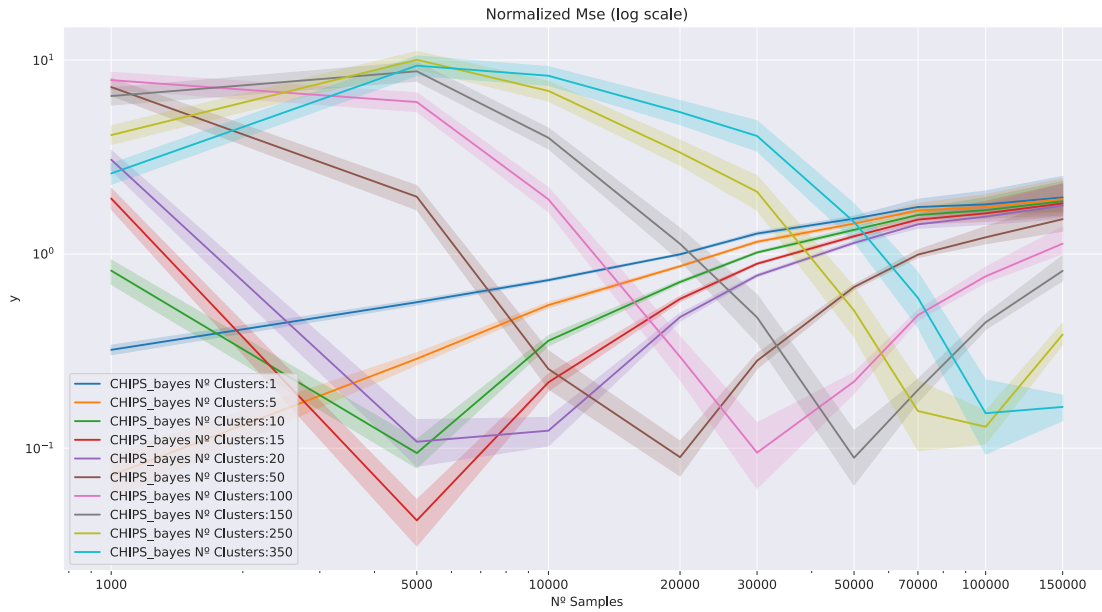
Algorithmically, since the CHIPS estimator can be regarded as performing the same procedure as IPS with different weights and rewards, the time complexity given n logging samples, and clustering method ξ can be expressed as:

$$\text{complexity}(\text{CHIPS}(n; \xi)) = \text{complexity}(\text{IPS}(n)) + \text{complexity}(\xi(n))$$

For example, since the time complexity of IPS is $\mathcal{O}(n)$, using DBSCAN ($\mathcal{O}(n \log n)$) as a clustering method, we would get a time complexity for CHIPS of $\mathcal{O}(n \log n)$. In our experiments, we used batch-Kmeans (Sculley, 2010) as clustering method, that has a time complexity of $\mathcal{O}(mkd_x t)$ where m is the batch size, k is the number of clusters, d_x is the dimension of the features and t is the number of iterations. In the implementation used, we fixed $m = 1024$ and $t = 100$, therefore, in this case, the time complexity of the CHIPS method is $\mathcal{O}(kd_x) + \mathcal{O}(n)$. The time complexity of the MIPS estimator can be estimated similarly as $\mathcal{O}(nd_e) + \mathcal{O}(n) = \mathcal{O}(nd_e)$, where the $\mathcal{O}(nd_e)$ term comes from the logistic regression used to estimate $\pi_0(a|x, e)$ (being e an action embedding) and d_e is the action embedding dimension. The methods using a supervised classifier (DM, DR, and MRDR) get their dominant term in time complexity from the training process of the classifier, in our case $\mathcal{O}(nds \log n)$ with s being the number of trees. In practice, this means that DM, , and MRDR will have significantly higher execution times (see Figure 21 (a)), and CHIPS will generally be faster than MIPS since $k \ll n$ to leverage the cluster structure, as we can appreciate in Figure 21 (b). The space complexity of CHIPS can be estimated following a similar approach, for example when using the KMeans algorithm the space complexity is $\mathcal{O}(n(k+d)) + \mathcal{O}(n) = \mathcal{O}(n(k+d))$.



(a) Normalized performance of CHIPS (MAP) with respect to IPS using different number of clusters and actions.



(b) Normalized performance of CHIPS (MAP) with respect to IPS using different number of clusters and logging samples.

Figure 12: Bi parametric experiments results using different number of clusters for analyzing CHIPS capabilities when increasing actions (a) and logging samples (b).

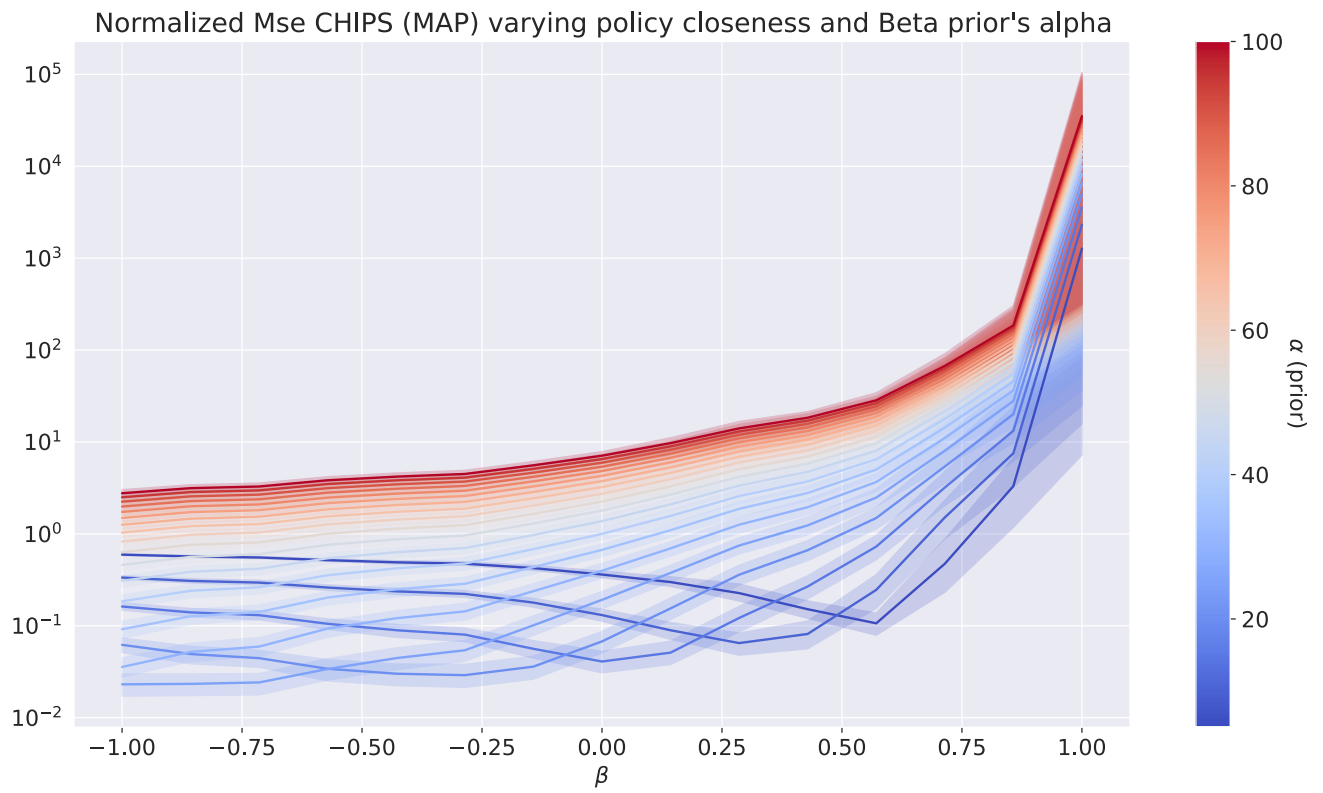
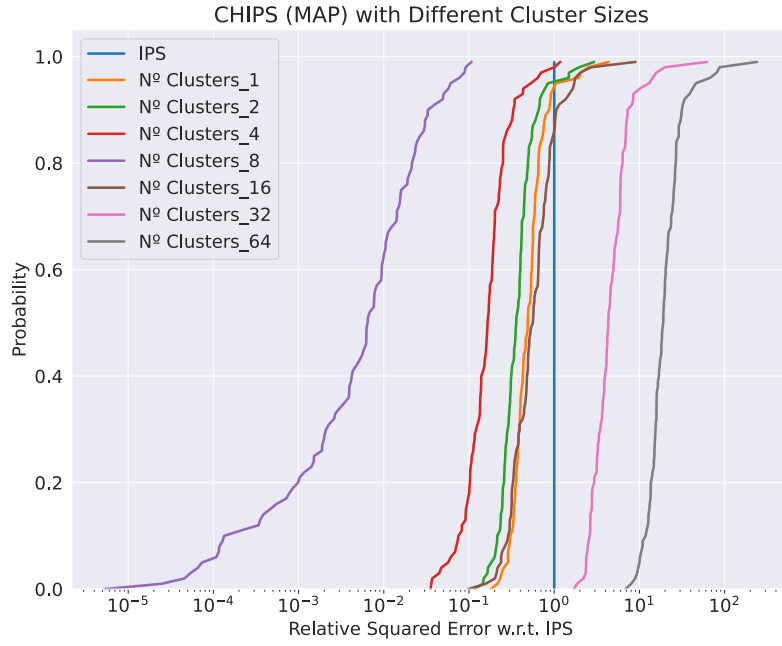
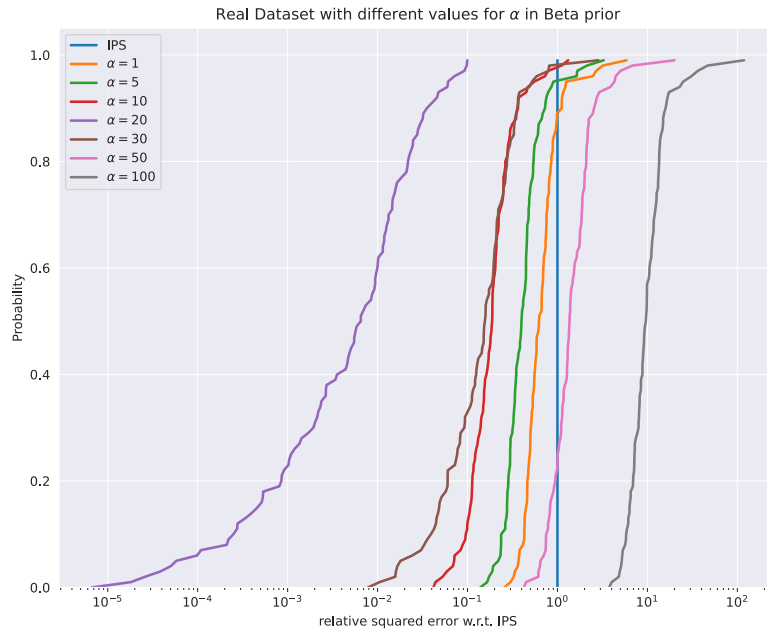


Figure 13: Normalized performance of CHIPS (MAP) with respect to IPS using different values for the α parameter in the Beta prior and distributional shift between logging and evaluation policies (β).



(a) ECDFs of CHIPS (MAP) using different number of clusters in the real dataset.



(b) ECDFs of CHIPS (MAP) using different values of α for the Beta prior in the real dataset.

Figure 14: Additional experiments varying the number of clusters and the α parameter in the Beta prior for CHIPS (MAP) in the real dataset (using 100000 samples).

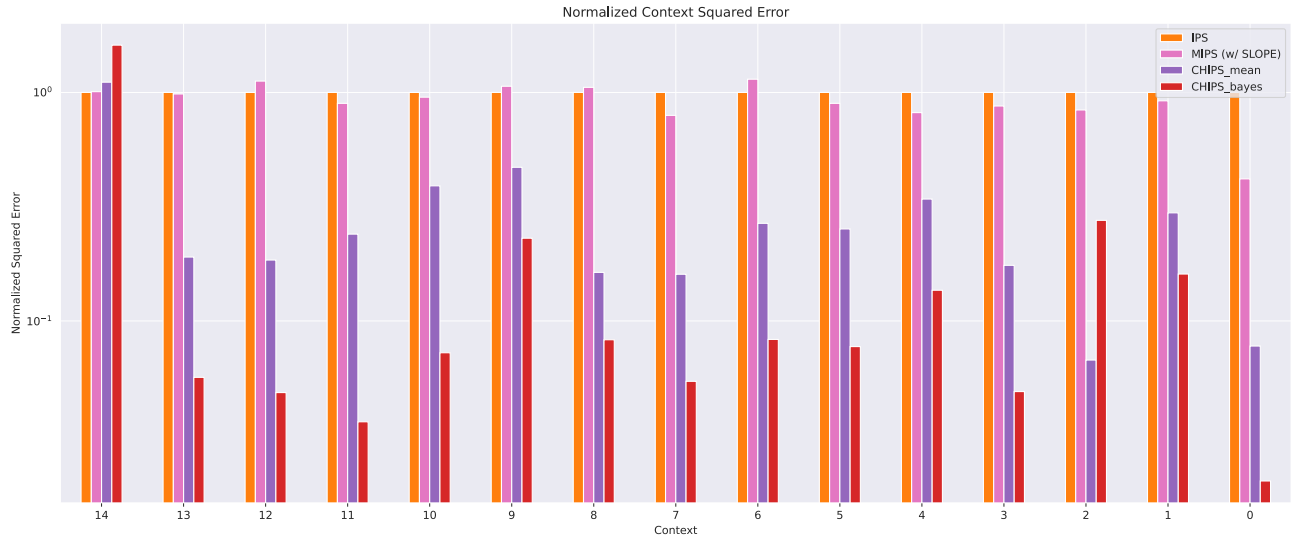


Figure 15: Normalized MSE with respect to IPS of the expected rewards for the 15 most common context-action pairs in the real logging dataset.

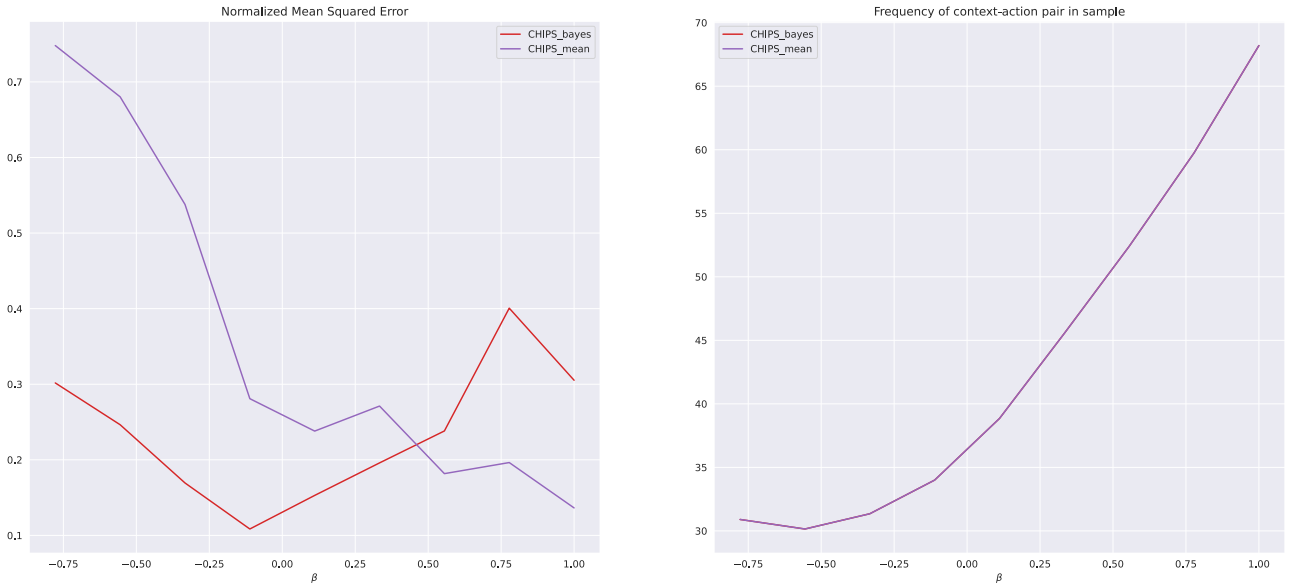


Figure 16: Normalized MSE of CHIPS with respect to IPS (left) and samples in the associated cluster (right) for the most common context-action pair in the evaluation policy while varying the distributional shift (β) in the synthetic dataset.

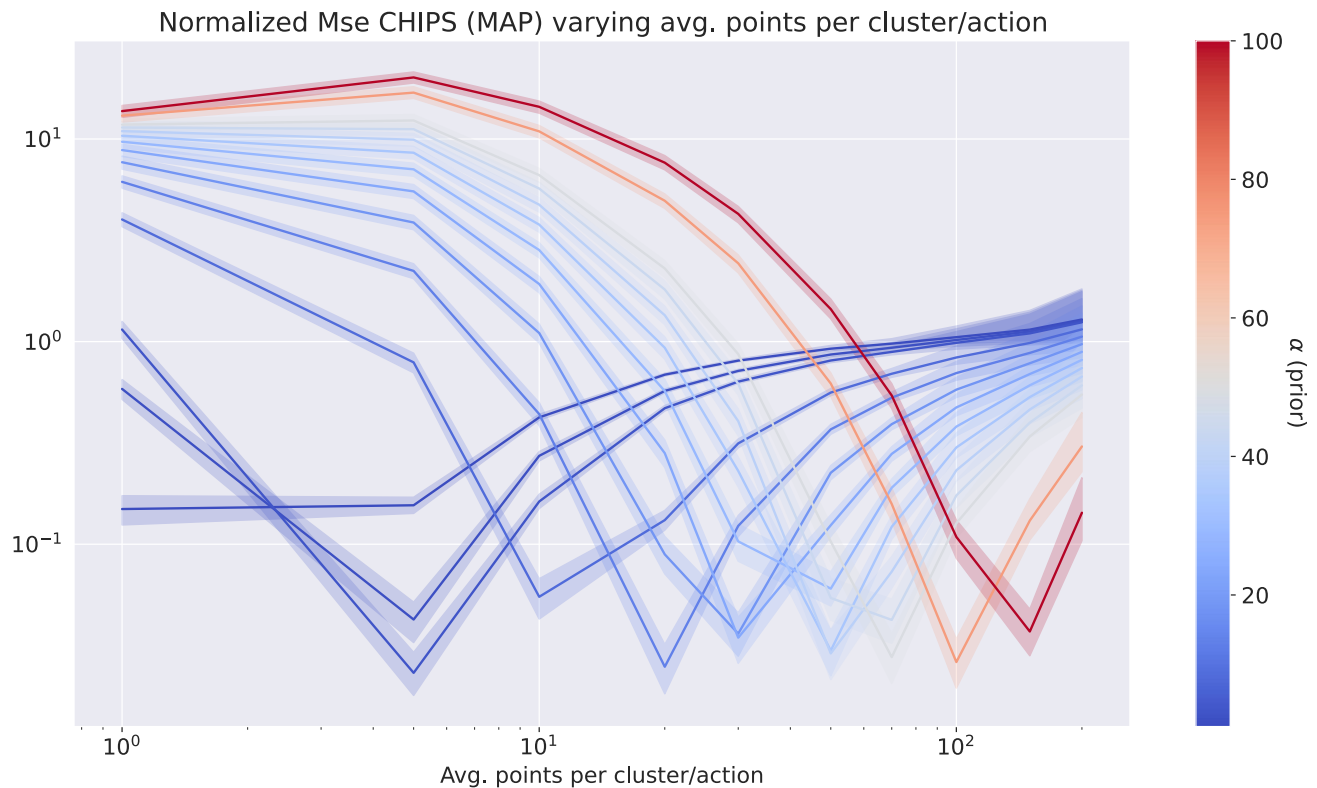


Figure 17: Normalized MSE of CHIPS (MAP) with respect to IPS using different values of α and number of expected data points per cluster-action.



Figure 18: MSE of CHIPS (MAP) using α selection policy with respect to IPS, DR, DM (Taufiq et al., 2023), DRos (Su et al., 2020) and SwitchDR (Wang, Agarwal, and Dudík, 2017).

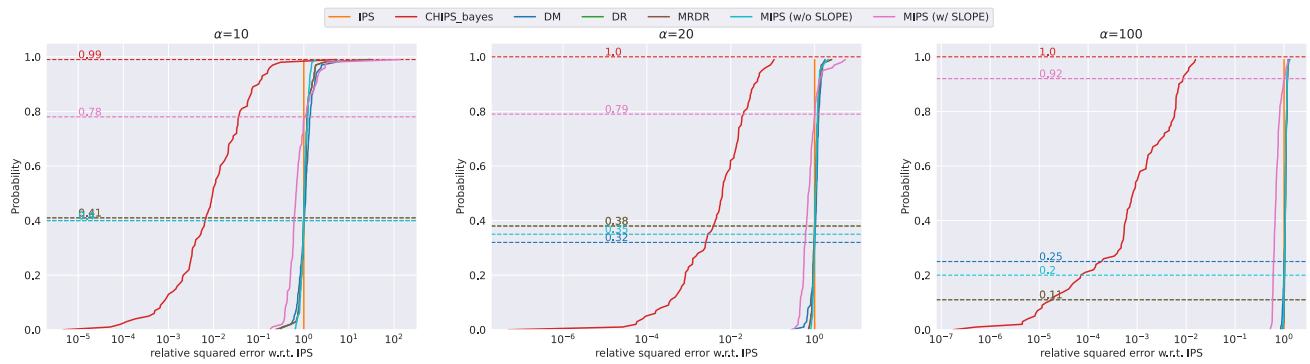


Figure 19: ECDF of the relative mean squared error with respect to IPS for the real dataset using 50000 (left), 100000 (center), and 500000 (right) logging samples and the α selection process.

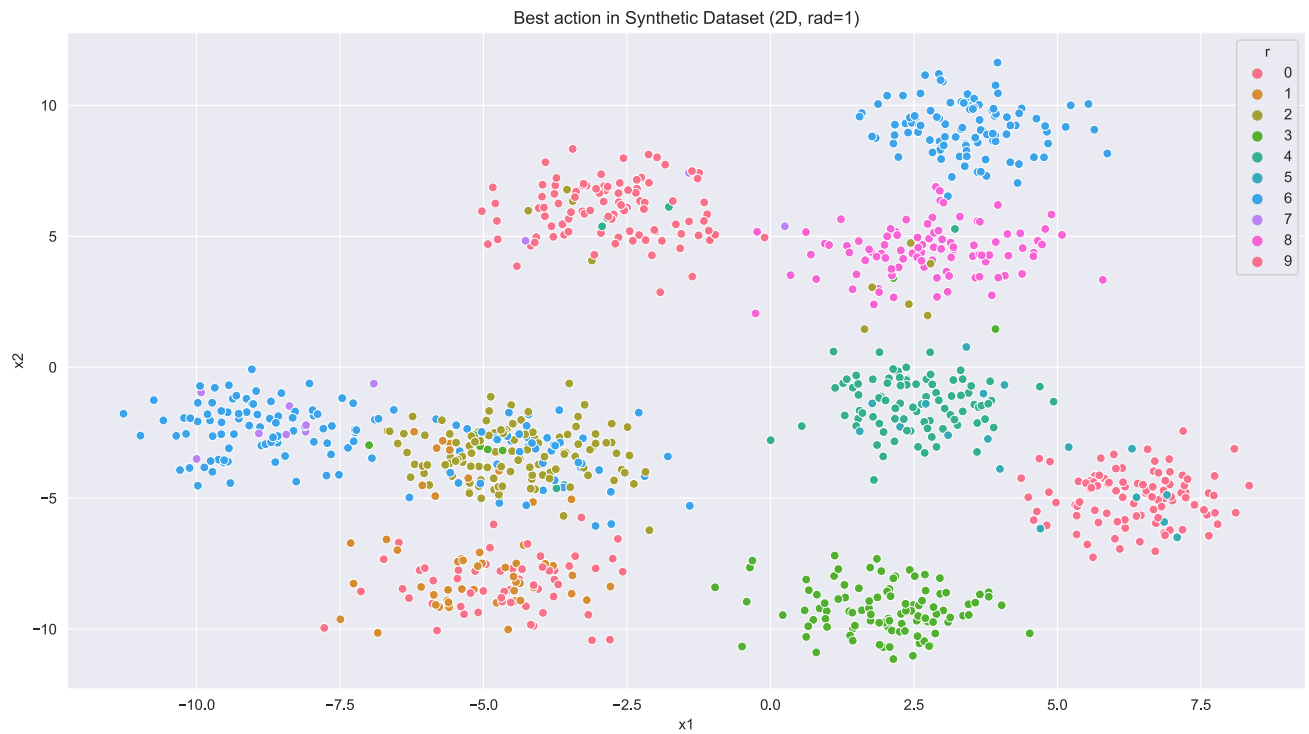
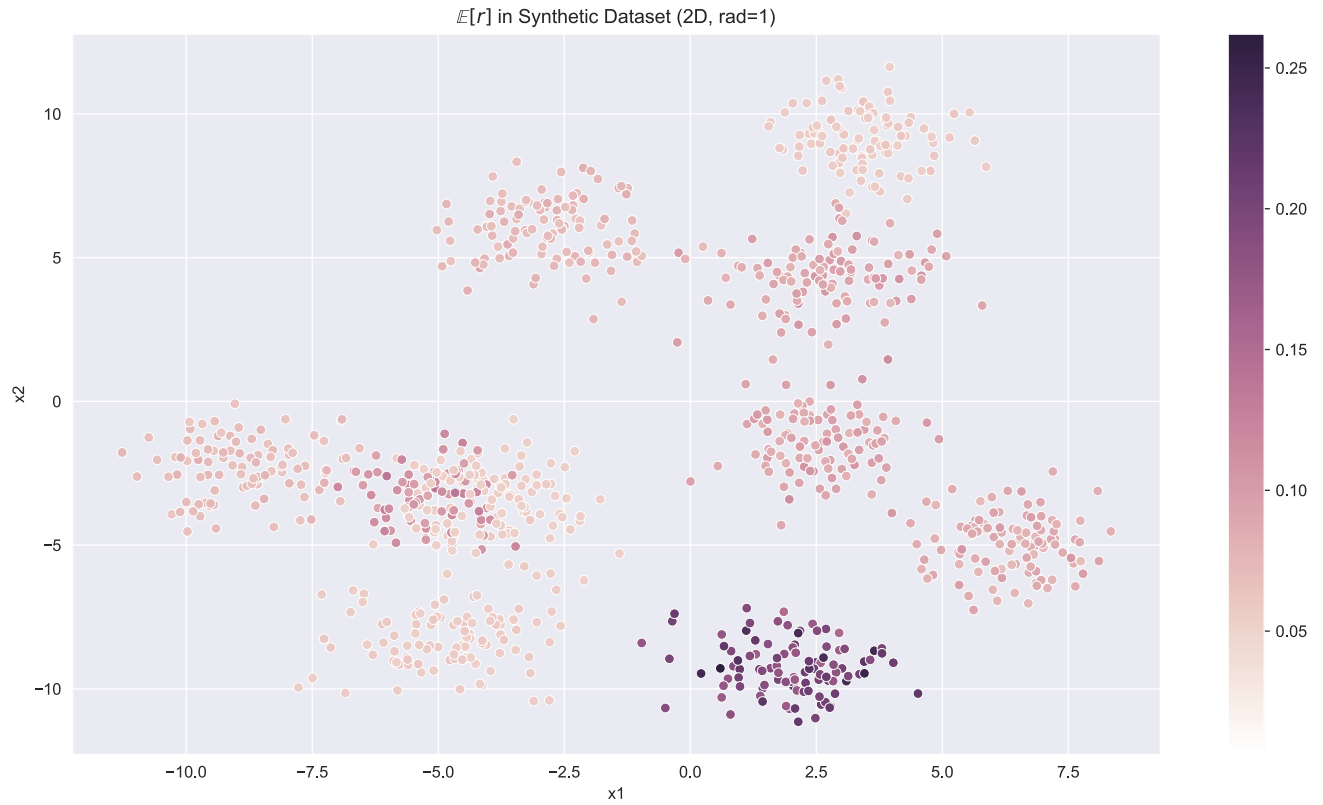
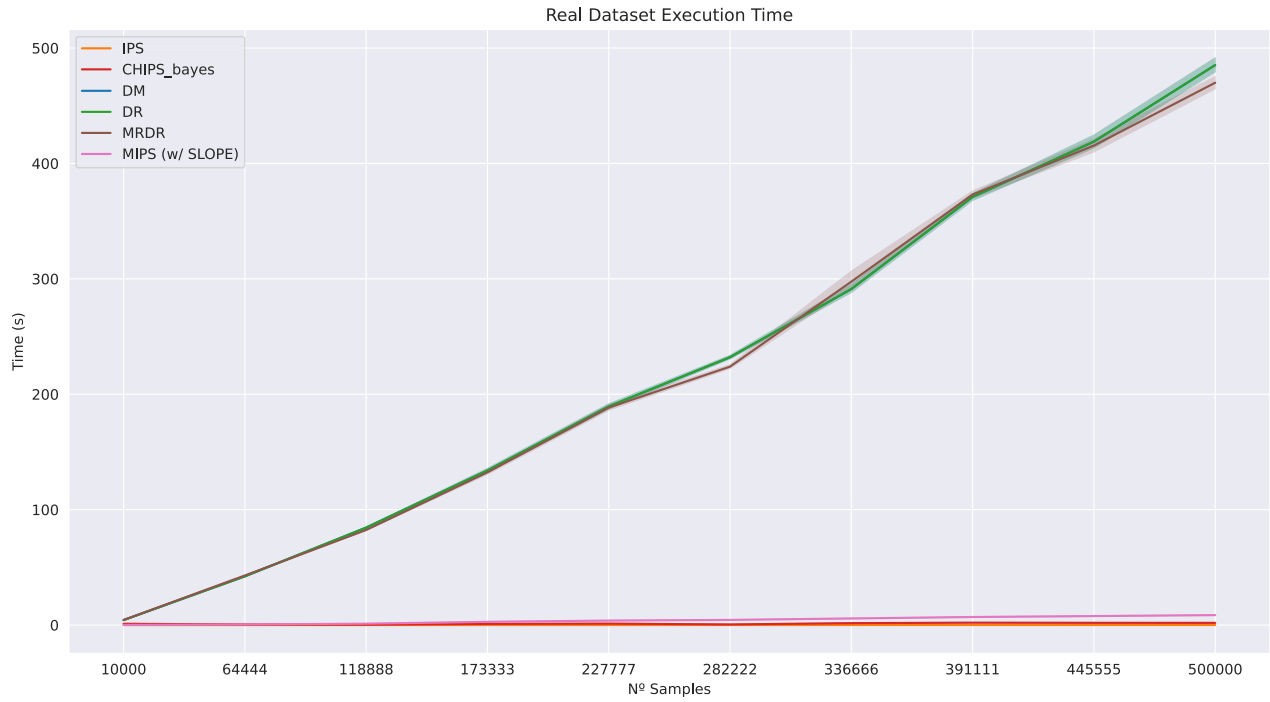
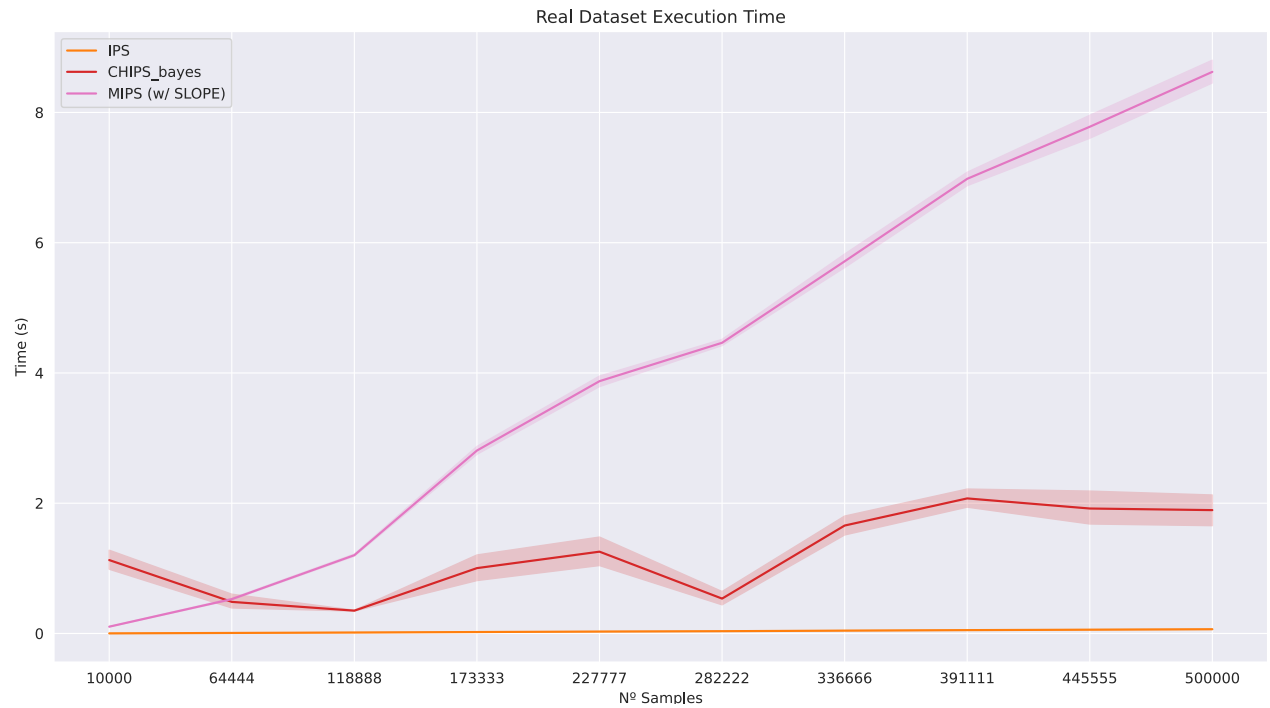


Figure 20: Representation of the synthetic dataset using 2-dimensional contexts.



(a) Execution times for the real dataset including DM, DR, and MRDR.



(b) Execution times for the real dataset for IPS, CHIPS, and MIPS.

Figure 21: Average execution times increasing the sample size (100 executions per sample size).