Concept-Based Off-Policy Evaluation

Ritam Majumdar, Jack Teversham, Sonali Parbhoo

Keywords: Off Policy Evalutaion, Interpretability, Concept Bottleneck Models, Reliable OPE

Summary

Evaluating off-policy decisions using batch data is challenging because of limited sample sizes which lead to high variance. Identifying and addressing the sources of this variance is crucial to improve off-policy evaluation in practice. Recent research on Concept Bottleneck Models (CBMs) shows that using human-explainable concepts can improve predictions and provide additional context for understanding decisions. In this paper, we propose incorporating an analogous notion of concepts into OPE to provide additional context that may help us identify specific areas where variance is high. We introduce a family of new concept-based OPE estimators and show that these estimators have two key properties when the concepts are known in advance: they remain unbiased whilst reducing variance of overall estimates. Since real-world applications often lack predefined concepts, we further develop an end-to-end algorithm to learn interpretable, concise, and diverse concepts optimized for variance reduction in OPE. Our experiments on synthetic and real-world datasets show that both known and learnt concept-based estimators significantly improve OPE performance. Crucially, our concept-based estimators offer two advantages over existing OPE methods. First, they are easily interpretable. Second, they allow us to isolate specific concepts contributing to variance. Upon performing targeted interventions on these concepts, we can further enhance the quality of OPE estimators.

Contribution(s)

- 1. We introduce a new family of IS estimators based on interpretable concepts. [Section 3] **Context:** Previous works perform IS in the state representations, we explicitly define what is a concept representation and tie the original definition of IS under concepts.
- We derive theoretical conditions ensuring lower variance compared to existing IS estimators. [Section 4]

Context: We compare the variance of the Concept-OPE estimators with traditional IS/PDIS and MIS estimators and devise conditions under which the variance is reduced.

- 3. We propose an end-to-end algorithm for optimizing parameterized concepts when concepts are unknown, using OPE characteristics like variance. [Section 5] **Context:** Under real-world scenarios, the concepts are typically unknown or hard to define, which adds to the complexity of performing OPE. In this section, we propose a novel algorithm which learns concepts that satisfy the desidarata: Explainability, Conciseness, Diversity while optimizing for variance.
- We show, through synthetic and real experiments, that our estimators for both known and unknown concepts outperform existing ones. [Sections 4,5]
 Context: None
- 5. We interpret the learned concepts to explain OPE characteristics and suggest intervention strategies to further improve OPE estimates. [Section 6] Context: Interventions have been typically studied in the context of improving the CBM performance in a supervised learning regime, we instead use interpretations to explain where a concept-OPE estimator has high variance and intervene to reduce variance.

Concept-Based Off-Policy Evaluation

Ritam Majumdar, Jack Teversham, Sonali Parbhoo

{r.majumdar24, jack.teversham22, s.parbhoo}@imperial.ac.uk

Imperial College London

Abstract

Evaluating off-policy decisions using batch data is challenging because of limited sample sizes which lead to high variance. Identifying and addressing the sources of this variance is crucial to improve off-policy evaluation in practice. Recent research on Concept Bottleneck Models (CBMs) shows that using human-explainable concepts can improve predictions and provide additional context for understanding decisions. In this paper, we propose incorporating an analogous notion of concepts into OPE to provide additional context that may help us identify specific areas where variance is high. We introduce a family of new concept-based OPE estimators and show that these estimators have two key properties when the concepts are known in advance: they remain unbiased whilst reducing variance of overall estimates. Since real-world applications often lack predefined concepts, we further develop an end-to-end algorithm to learn interpretable, concise, and diverse concepts optimized for variance reduction in OPE. Our experiments on synthetic and real-world datasets show that both known and learnt concept-based estimators significantly improve OPE performance. Crucially, our concept-based estimators offer two advantages over existing OPE methods. First, they are easily interpretable. Second, they allow us to isolate specific concepts contributing to variance. Upon performing targeted interventions on these concepts, we can further enhance the quality of OPE estimators.

1 Introduction

In domains like healthcare, education, and public policy, where interacting with the environment can be risky, prohibitively expensive, or unethical (Sutton & Barto, 2018; Murphy et al., 2001; Mandel et al., 2014), estimating the value of a policy from batch data before deployment is essential for the practical application of RL. OPE aims to estimate the effectiveness of a specific policy, known as the evaluation or target policy, using offline data collected beforehand from a different policy, known as the behavior policy (e.g., Komorowski et al. (2018a); Precup et al. (2000); Thomas & Brunskill (2016); Jiang & Li (2016)).

Importance sampling (IS) methods are a popular class of methods for OPE which adjust for distributional mismatches between behavior and target policies by reweighting historical data, yielding generally unbiased and consistent estimates (Precup et al., 2000). Despite their desirable properties (Thomas & Brunskill, 2016; Jiang & Li, 2016; Farajtabar et al., 2018), IS methods often face high variance, especially with limited overlap between behavioral samples and evaluation targets or in data-scarce conditions. Evaluation policies may outperform behavior policies for specific individuals or subgroups (Keramati et al., 2021b), making it misleading to rely solely on aggregate policy value estimates. In practice however, these groups are often unknown, prompting the need for methods to learn interpretable characterizations of the circumstances where the evaluation policy benefits certain individuals over others.

In this paper, we propose performing OPE using interpretable concepts (Koh et al., 2020; Madeira et al., 2023) instead of relying solely on state and action information. We demonstrate that this

approach offers significant practical benefits for evaluation. These concepts can capture critical aspects in historical data, such as key transitions in a patient's treatment or features affecting short-term outcomes that serve as proxies for long-term results. By learning interpretable concepts from data, we introduce a new family of concept-based IS estimators that provide more accurate value estimates and stronger statistical guarantees. Additionally, these estimators allow us to identify which concepts contribute most to variance in evaluation. When the evaluation is unreliable, we can modify, intervene on, or remove these high-variance concepts to assess how the resulting evaluation improves (Marcinkevičs et al., 2024; Madeira et al., 2023).

A physician treating two patients infected with the same virus, with similar disease dynamics focuses on overall trends rather than precise viral load values when administering treatments. That is, if a drug lowers one patient's viral load below a threshold, it may also help the other, whereas it may be ineffective for a patient with a different disease trajectory. This distinction between a *concept* — a generalizable trend, such as viral load reduction — and a *state* — specific measurements at individual time points e.g. viral load — is key. Learning concepts that capture these trends, rather than isolated values, can better guide treatment decisions and evaluation. This idea is illustrated in Fig 1.

Our work makes the following contributions: i) We introduce a new family of IS estimators based on interpretable concepts; ii) We derive theoretical conditions ensuring lower variance compared to existing IS estimators; iii) We propose an end-to-end algorithm for optimizing parameterized concepts when concepts are unknown, using OPE characteristics like variance; iv) We show, through synthetic and real experiments, that our estimators for both known and unknown concepts outperform existing ones; v) We interpret the learned concepts to explain OPE characteristics and suggest intervention strategies to further improve OPE estimates.



Figure 1: Simple example of a state vs concept. In this scenario, the state is the viral load in a patient's blood, whereas the concept is defined as the viral load being above or below a certain threshold x. The concept divides patients into two groups, in which different treatments are administered, indicated by the frequency of syringes. We do evaluation based on these two concepts, rather than the unique values of the viral loads.

2 Preliminaries

Concept Bottleneck Models Conventional CBMs learn a mapping from some input features $x \in \mathbb{R}^d$ to targets y via some interpretable concepts $c \in \mathbb{R}^k$ based on training data of the form $\{x_n, c_n, y_n\}_{n=1}^N$. This mapping is a composition of a mapping from inputs to concepts, $f : \mathbb{R}^d \to \mathbb{R}^k$, and a mapping from concepts to targets, $g : \mathbb{R}^k \to \mathbb{R}$. These may be trained via independent, sequential or joint training (Marcinkevičs et al., 2024). Variations which consider learning concepts in a greedy fashion or in a semi-supervised way include Wu et al. (2022); Havasi et al. (2022).

Markov Decision Processes (MDP). An MDP is defined by a tuple $\mathcal{M} = (S, \mathcal{A}, P, R, \gamma, T)$. S and \mathcal{A} are the state and action spaces, $P : S \times \mathcal{A} \to \Delta(S)$ and $R : S \times \mathcal{A} \to \Delta(\mathbb{R})$ are the transition and reward functions, $\gamma \in [0, 1]$ is the discount factor, $T \in \mathbb{Z}^+$ is the fixed time horizon. A policy $\pi : S \to \Delta(\mathcal{A})$ is a mapping from each state to a probability distribution over actions in \mathcal{A} . A T-step trajectory following policy π is denoted by $\tau = [(s_t, a_t, r_t, s_{t+1})]_{t=1}^T$ where $s_1 \sim d_1, a_t \sim \pi(s_t), r_t \sim r(s_t, a_t), s_{t+1} \sim p(s_t, a_t)$. The value function of policy π , denoted by $V_{\pi} : S \to \mathbb{R}$, maps each state to the expected sum of rewards starting from that state following policy π . That is, $V_{\pi}(s) = \mathbb{E}_{\pi}[\sum_{t=1}^T \gamma^{t-1}r_t|s_1 = s]$.

Off-Policy Evaluation. In OPE, we have a dataset of *T*-step trajectories $\mathcal{D} = {\tau^{(n)}}_{n=1}^{N}$ independently generated by a *behaviour policy* π_b . Our goal is to estimate the value function of another

evaluation policy, π_e . We aim to use \mathcal{D} to produce an estimator, \hat{V}_{π_e} , that has low mean squared error, $MSE(V_{\pi_e}, \hat{V}_{\pi_e}) = \mathbb{E}_{\mathcal{D} \sim P_{\pi_b}^{\tau}}[(V_{\pi_e} - \hat{V}_{\pi_e})^2]$. Here, $P_{\pi_b}^{\tau}$ denotes the distribution of trajectories τ , under π_b , from which \mathcal{D} is sampled.

3 Concept-Based Off-Policy Evaluation

The goal of our work is to incorporate the notion of concepts into off-policy evaluation for improved interpretability and variance reduction. In this section, we formally define the notion of a concept, outline its desiderata and formally introduce a class of OPE estimators. In subsequent sections, we discuss how these estimators can be used when a) concepts are known from domain expertise (see Section 4), and b) concepts are unknown and must be learnt using a parametric representation (see Section 5).

3.1 Defining a Concept for OPE

Given a dataset $\mathcal{D} = \{\tau^{(n)}\}_{n=1}^N$ of N T-step trajectories, let $\phi : S \times \mathcal{A} \times R \times S \to \mathcal{C} \in \mathbb{R}^d$ denote a function that maps trajectory histories h_t to interpretable concepts in d-dimensional concept space \mathcal{C} . This mapping results in the concept vector $c_t = [c_t^1, c_t^2, ..., c_t^d]$ at time t, defined by $\phi(h_t)$. These concepts can capture various vital information in the history h_t , such as transition dynamics, short-term rewards, influential states, inter-dependencies in actions across time steps, etc. Without loss of generality, we consider concepts c_t as functions of current state s_t , however this could be extended to include historical information. This considers the scenario where concepts capture important information based on the criticality of the state. The concept function ϕ satisfies the following desiderata: explainability, conciseness, better trajectory coverage and diversity. A detailed description of these desiderata is provided in Supplement B.

3.2 Concept-Based Estimators for OPE.

We propose a new class of concept-based OPE estimators, adapting existing non-concept-based methods to integrate concepts into OPE. Here, we present the results specifically for per-decision IS and standard IS estimators, as these serve as the foundation for several other estimators. We also demonstrate in Supplement D how these methods can be extended to other estimators.

Definition 3.1 (Concept-Based Importance Sampling (CIS)).

$$\hat{V}_{\pi_e}^{CIS} = \frac{1}{N} \sum_{n=1}^{N} \rho_{0:T}^{(n)} \sum_{t=0}^{T} \gamma^t r_t^{(n)}; \quad \rho_{0:T}^{(n)} = \prod_{t'=0}^{T} \frac{\pi_e^c(a_{t'}^{(n)}|c_{t'}^{(n)})}{\pi_b^c(a_{t'}^{(n)}|c_{t'}^{(n)})}$$

Definition 3.2 (Concept-based Per-Decision Importance Sampling, (CPDIS)).

$$\hat{V}_{\pi_e}^{CPDIS} = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=0}^{T} \gamma^t \rho_{0:t}^{(n)} r_t^{(n)}; \quad \rho_{0:t}^{(n)} = \prod_{t'=0}^{t} \frac{\pi_e^c(a_{t'}^{(n)}|c_{t'}^{(n)})}{\pi_b^c(a_{t'}^{(n)}|c_{t'}^{(n)})}$$

Concept-based variants of IS replace the traditional IS ratio with one that leverages the concept c_t at time t instead of the state s_t . This enables customized evaluations for various concept types, such as: 1) subgroups with similar short-term outcomes, 2) cases with comparable state-visitation densities, and 3) subjects with high-variance transitions. Details on selecting concept types are provided in Supplement C.

4 Concept-based OPE under Known Concepts

We first consider the scenario where the concepts are known apriori using domain knowledge and human expertise. These concepts must satisfy the desiderata defined in Supplement B.

4.1 Theoretical Analysis of Known Concepts

In this subsection, we discuss the theoretical guarantees of OPE under known concepts. We make the completeness assumption where every action of a particular state has a non-zero probability of appearing in the batch data. When this assumption is satisfied, we obtain unbiasedness and lower variance when compared with traditional estimators. Proofs follow in Supplement E.

Assumption 4.1 (Completeness). $\forall s \in S, a \in A$, if $\pi_e(a|s), \pi_e^c(a|c) > 0$ then $\pi_b(a|s), \pi_b^c(a|c) > 0$.

This assumption states that if an action appears in the evaluation with some probability, it also has some probability of being in the batch data.

Assumption 4.2. $\forall s \in S, a \in A, |\pi_e^c(a|c) - \pi_e(a|s)| < \beta \text{ and } |\pi_b^c(a|c) - \pi_b(a|s)| < \beta.$ This assumption states that for all states s, the policies conditioned on concepts are allowed to differ from the state policies by atmost β , which is defined by the practitioner.

This assumption constrains concept-based policies to be close to state-based policies, with a maximum allowable difference of β , defined by the practitioner. This is to ensure that the evaluation policy π_e^c under concepts is reflective of the original policy π_e . If the practitioner is confident in the state representation, they may set a lower β to find concepts that align closely with state policies. Conversely, a higher β allows for more deviation between concept and state policies.

Theorem 4.3 (Bias). Under known-concepts, when assumption 4.1 holds, both $\hat{V}_{\pi_e}^{CIS}$ and $\hat{V}_{\pi_e}^{CPDIS}$ are unbiased estimators of the true value function V_{π_e} .

Theorem 4.4 (Variance comparison with traditional OPE estimators). When $Cov(\rho_{0:t}^c r_t, \rho_{0:k}^c r_k) \leq Cov(\rho_{0:t}r_t, \rho_{0:k}r_k)$, the variance of known concept-based IS estimators is lower than traditional estimators, i.e. $\mathbb{V}_{\pi_b}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_b}[\hat{V}^{IS}], \mathbb{V}_{\pi_b}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_b}[\hat{V}^{PDIS}].$

As noted in Jiang & Li (2016), the covariance assumption across timesteps is crucial yet challenging for OPE variance comparisons. Concepts being interpretable allows a user to design policies which align with this assumption, thereby reducing variance. We also compare concept-based estimators to the MIS estimator, the gold standard for minimizing variance via steady-state distribution ratios.

Theorem 4.5 (Variance comparison with MIS estimator). When $Cov(\rho_{0:t}^c r_t, \rho_{0:k}^c r_k) \leq Cov(\frac{d^{\pi_e}(s_t, a_t)}{d^{\pi_b}(s_t, a_t)}r_t, \frac{d^{\pi_e}(s_k, a_k)}{d^{\pi_b}(s_k, a_k)}r_k)$, the variance of known concept-based IS estimators is lower than the Variance of MIS estimator, i.e. $\mathbb{V}_{\pi_b}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_b}[\hat{V}^{MIS}]$, $\mathbb{V}_{\pi_b}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_b}[\hat{V}^{MIS}]$.

Finally, we evaluate the CR-bounds on the MSE and quantify the tightness achieved using concepts.

Theorem 4.6 (Confidence bounds for Concept-based estimators). The Cramer-Rao bound on the Mean-Square Error of CIS and CPDIS estimator under known-concepts is tightened by a factor of K^{2T} , where K is the ratio of the cardinality of the concept-space and state-space.

With limited samples, certain relevant states are underrepresented in the behavior policy, leading to a low $\pi_b(.|s)$ and corresponding high IS ratio. However, an alternative state s' in the data may closely resemble s (e.g., similar blood pressure values). Thus, even if s is missing, it can be characterized by s' through the concept function $\phi(s')$, as $\phi(s) \approx \phi(s')$. Consequently, while $\pi_b(.|s)$ is low, $\pi_b^c(.|s)$ is higher, as s is effectively represented via s'. This reduces the IS ratio in concept space and tightens the overall bounds.

4.2 Experimental Setup and Metrics

Environments: We consider a synthetic domain: WindyGridworld and the real world MIMIC-III dataset for acutely hypotensive ICU patients as our experiment domains for the rest of the paper.

WindyGridworld: The goal is to reach the top-right corner of the grid. The states are defined by the 2D co-ordinates, and actions are directions up, down, left, right. We (as human experts) define a concept $c_t = \phi$ (distance to target, wind) as a function of the distance to the target and the wind acting on the agent at a given state. This concept can take 25 unique values, ranging from 0 to 24. For example: $c_t = 0$ when distance to target $\in [15, 19] \times [15, 19]$ and wind = [0, 0]. The first and second co-ordinates represent the horizontal and vertical features respectively. Detailed description of known concepts in Supplement H. All experiments conducted over 20 seeds.

MIMIC: The goal is to treat and manage hypotensive patients. The state space consists of the physiological quantities of the patient while actions correspond to quantities of IV-fluids and vasopressors.

Concepts $c_t \in \mathbb{Z}^{15}$ are a function of 15 different vital signs (interpretable features) of a patient at a given timestep. The vital signs considered are: Creatinine, FiO₂, Lactate, Partial Pressure of Oxygen (PaO₂), Partial Pressure of CO₂, Urine Output, GCS score, and electrolytes such as Calcium, Chloride, Glucose, HCO₃, Magnesium, Potassium, Sodium, and SpO₂. Each vital sign is binned into 10 discrete levels, ranging from 0 (very low) to 9 (very high). Experiments conducted over 20 seeds.

For example, a patient with the concept representation [0, 2, 1, 1, 2, 0, 9, 5, 2, 0, 6, 2, 1, 5, 9] shows the following conditions: acute kidney injury-AKI (very low creatinine), severe hypoxemia (very low PaO₂), metabolic alkalosis (very high SpO₂), and critical electrolyte imbalances (low potassium and magnesium), along with severe hypoglycemia. The normal GCS score indicates preserved neurological function, but over-oxygenation and potential respiratory failure are likely. The combination of anuria, AKI, and hypoglycemia points strongly toward hypotension or shock as underlying causes.

Policy descriptions: In the case of WindyGridworld, we run a PPO Schulman et al. (2017) algorithm for 10k epochs and consider the evaluation policy π_e as the policy at epoch 10k, while the behavior policy π_b is taken as the policy at epoch 5k. For the MIMIC case, we generate the behavior policy π_b by running an Approximate Nearest Neighbors algorithm with 200 neighbors, using Manhattan distance as the distance metric. The evaluation policy π_e involves a more aggressive use of vasopressors (10% more) compared to the behavior policy. See Supplement G for further details.

Metrics: In the case of the synthetic domain, we measure bias, variance, MSE, and the effective sample size (ESS) to assess the quality of our concept-based OPE estimates. The ESS is defined as $N \times \frac{\mathbb{V}_{\pi_e}[\hat{V}_{\pi_e}^{on-policy}]}{\mathbb{V}_{\pi_b}[\hat{V}_{\pi_e}]}$, where N is the number of trajectories in the off-policy data, and $\hat{V}_{\pi_e}^{on-policy}$

and \hat{V}_{π_e} are the on-policy and OPE estimates of the value function, respectively. For MIMIC, where the true on-policy estimate is unknown due to the unknown transition dynamics and environment model, we only consider variance as the metric. Additionally, we compare the Inverse Propensity scores (IPS) under concepts and states to better underscore the reasons for variance reduction.

4.3 Results and Discussion

Good concept-based estimators demonstrate reduced variance, improved ESS, and lower MSE compared to traditional estimators, although they come with slightly higher bias. Fig 2 compares known-concept and traditional OPE estimators. We observe a consistent reduction in variance and an increase in ESS across all sample sizes for the concept-based estimators. Although our theoretical analysis suggests that known-concept estimators are unbiased in the asymptotic case, practical results indicate some bias due to finite sample size. While unbiased estimates are generally preferred, they can lead to higher errors when the behavior policy does not cover all states. This issue is especially pronounced in limited data settings, which are common in medical applications. Despite this bias-variance trade-off, the MSE for concept-based OPE estimators shows a 1-2 order of magnitude improvement over traditional estimators due to significant variance reduction. In the real-world MIMIC example, concept-based estimators exhibit a variance reduction of one order of magnitude compared to traditional OPE estimators. This shows that characterizing diverse states, such as varying grid world positions or patient vital signs, in terms of shared concepts based on common attributes, improves OPE characterization.

The frequency of higher IPS scores is reduced in good concepts compared to states. Fig 2c compares IPS scores in good concept and state estimators. We observe, the frequency of lower IPS scores is higher under concepts as opposed to states. This indicates the source of variance reduction in Concept-based OPE lies in the lowering of the IPS scores, which is also backed theoretically in Theorem 4.4 when the rewards r_t equal 1.

Imperfect concepts baseline: While known concepts display superior OPE performance, in real world scenarios, concepts are often poorly described. We thus, perform an experiment where the concepts known but poor and study the resulting OPE performance. For Windygridworld, we define concepts as functions solely of the horizontal distance to the target. This approach neglects critical information such as vertical distance to the target, wind effects, and region penalties. As a result, these concepts violate one of the primary desiderata: diversity. By capturing only one important



Figure 2: *Known Concepts.* (a) *Windy Gridworld:* Concept-based estimators with good concepts have lower variance and MSE, and higher ESS compared to traditional OPE estimators, with a higher bias. For poor concepts, we observe the OPE performance to be poor across all metrics compared to traditional estimators. (b) *MIMIC:* Good Concept-based estimators have lower variance compared to traditional OPE estimators, while poor concepts have higher variance. (c) Under good concepts, the frequency of high IPS scores is lower compared to traditional estimators, whereas for poor concepts, the frequency is higher across both domains. This indicates the source of variance reduction in good concept estimators lies in the lowered IPS scores.

concept dimension while disregarding others, these poor concepts fail to represent the full complexity of the environment.

Poor concepts exhibit inferior OPE characteristics across all metrics. From fig 2a and 2b, we observe that poor concepts exhibit higher bias, variance, and MSE, along with lower ESS, compared to traditional OPE estimators. Additionally, fig 2c shows an increased frequency of high IPS scores for poor concepts. This demonstrates that not all concept-based estimators improve performance; their quality is crucial and depends on the desiderata they satisfy. It also underscores the need for an algorithm that learns concepts with favorable OPE characteristics, particularly in complex domains or scenarios with imperfect experts. We explore this in the next section. Nonetheless, poor concepts still allow for interventions, as their impact on OPE metrics can be systematically analyzed and addressed.

5 Concept-based OPE under Unknown Concepts

While domain knowledge and predefined concepts can enhance OPE, in real-world situations concepts are typically unknown. In this section, we address cases where concepts are unknown and must be estimated. We learn a parametric representation of concepts via CBMs, which initially may not meet the required desiderata. This section introduces a methodology to optimize parameterized concepts to meet explicitly these desiderata, alongside improving OPE metrics like variance.

Learning concepts that characterize relevant trajectory information. Algorithm 1 outlines the training methodology. We split the batch of trajectories \mathcal{D} into training trajectories \mathcal{T}_{train} and evaluation trajectories \mathcal{T}_{OPE} , with the evaluation policy π_e , the behavior policy π_b , and an OPE estimator (e.g. CIS/CPDIS) known beforehand. We aim to learn our concepts using a CBM parameterized by θ . The

Algorithm 1 Unknown Concept-based Off Policy Evaluation							
Require: Trajectories { $\mathcal{T}_{train}, \mathcal{T}_{OPE}$ }, Policies { π_e, π_b }, OPE Estimator.							
Ensure: CBM θ , concept policies $\tilde{\pi}^c \{\theta_b, \theta_e\}$							
Loss terms: { $L_{\text{output}}, L_{\text{interpretability}}, L_{\text{diversity}}, L_{\text{OPE-metric}}, L_{\text{policy}}$ } = 0							
1: while Not Converged do							
2: for trajectory in $\mathcal{T}_{\text{train}}$ do							
3: for (s, a, r, s', o) in trajectory do \triangleright Choices for $o: s'$ (Next state) / r (Next reward							
4: $c', o' \leftarrow \text{CBM}(s)$ \triangleright CBM predicts concept c' and output label o'							
5: $L_{\text{output}} += C_{\text{output}}(o, o') $ \triangleright Eg: MSE/Cross-entropy between true next state and predicted next state							
6: $L_{\text{interpretability}} += C_{\text{interpretability}}(c')$ \triangleright Eg: L1-loss over weights							
7: $L_{\text{diversity}} += C_{\text{diversity}}(c')$ \triangleright Eg: Cosine distance between sub-concepts							
8: $L_{\text{policy}} += C_{\text{policy}}(c') \triangleright \text{Eg: MSE/Cross-entropy between predicted logits and true logits in Assn 4.2}$							
9: end for							
10: end for							
11: Returns \leftarrow Estimator($\mathcal{T}_{train}, \pi_e, \pi_b, \text{CBM}$) \triangleright Eg: CIS/CPDIS							
: $\text{Loss}(\theta, \theta_b, \theta_e) = L_{\text{output}} + L_{\text{interpret.}} + L_{\text{diversity}} + L_{\text{policy}} + C_{\text{OPE-metric}}(\text{Returns}) \triangleright \text{Eg: Variance}$							
13: Gradient Descent on $\{\theta, \theta_b, \theta_e\}$ using $\text{Loss}(\theta, \theta_b, \theta_e)$							
14: end while							
15: Return Concept OPE Returns \leftarrow OPE Estimator($\mathcal{T}_{OPE}, \pi_e, \pi_b, CBM$)							

CBM maps states to outputs through an intermediary concept layer. In this work, the output o is the next state, indicating that the bottleneck concepts capture transition dynamics. Other possible outputs could include short-term rewards, long-term returns, or any user-defined information of interest present in the batch data. In addition to learning concepts, we also learn parameterized concept policies $\tilde{\pi}^c$ which maps concepts to actions parameterized by θ_b , θ_e for π_b , π_e respectively.

Optimizing concepts for variance reduction in OPE. For each transition tuple (s, a, r, s'), the CBM computes a concept vector c' and an output o'. Since the concepts are initially unknown, they do not inherently satisfy the concept desiderata and must be learned through constraints. Lines 5-7 impose soft constraints on the concepts to meet these desiderata using loss functions. The losses are updated based on output, interpretability, and diversity, with MSE used for C_{output} , L1 loss for $C_{\text{interpretability}}$, and cosine distance for $C_{\text{diversity}}$. In Line 8, we constrain the difference between the concept policies and the original policies to satisfy Assumption 4.2. For our experiments, we take maximum allowable difference $\beta = 0$, however a user can choose a different value to allow for more deviation in the concept policies π^c and original policies π . In line 11, we evaluate the OPE estimator's returns based on the concepts at the current iteration with metrics like variance. The aggregate loss, $\text{Loss}(\theta)$, guides gradient descent on CBM parameters θ . Finally, the OPE estimator is applied to \mathcal{T}_{OPE} using learned concepts, yielding concept-based OPE returns. The $C_{OPE-metric}$ reflects an OPE criteria (e.g. Variance reduction) which we are trying to optimise directly via gradient descent. Integrating multiple competing loss components makes this problem complex, and, to our knowledge, this is the first approach that incorporates the OPE metric directly into the loss function.

5.1 Theoretical Analysis of Unknown Concepts

The theoretical implications mainly differ in the bias, consequently MSE and their Confidence bounds on moving from known to unknown concepts, as analyzed below. Proofs are listed in Supplement F.

Theorem 5.1 (Bias). Under Assumptions 4.1, 4.2, the unknown concept-based estimators are biased. The change of measure theorem from probability distributions π_b to π_b^c is not applicable on moving from known to unknown concepts, leading to bias. In the special case where $\pi_b^c(.|c_t) = \pi_b(.|s_t)$, the estimator is unbiased.

Theorem 5.2 (Variance comparison with traditional OPE estimators). Under Assumption 4.2, when $Cov(\rho_{0:t}^{c}r_{t},\rho_{0:k}^{c}r_{k}) \leq Cov(\rho_{0:t}r_{t},\rho_{0:k}r_{k})$, the variance of concept-based IS estimators is lower than the traditional estimators, i.e. $\mathbb{V}_{\pi_{b}}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{IS}], \mathbb{V}_{\pi_{b}}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{PDIS}].$





Figure 3: *Learned Concepts.* (a-b) *Windy Gridworld and MIMIC:* In both environments, we see a reduction in variance in learned concepts compared to both traditional and known concept estimators, at a cost of higher bias. (c) The frequency of high IPS scores is lower for learned concepts compared to traditional and known concept estimators. This indicates our proposed algorithm learns alternative concepts which further reduce variance.

Theorem 5.3 (Variance comparison with MIS estimator). Under Assumption 4.2, when $Cov(\rho_{0:t}^c r_t, \rho_{0:k}^c r_k) \leq Cov(\frac{d^{\pi_e}(s_t, a_t)}{d^{\pi_b}(s_t, a_t)}r_t, \frac{d^{\pi_e}(s_k, a_k)}{d^{\pi_b}(s_k, a_k)}r_k)$, like known concepts, the variance is lower than the Variance of MIS estimator, i.e. $\mathbb{V}_{\pi_b}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_b}[\hat{V}^{MIS}], \mathbb{V}_{\pi_b}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_b}[\hat{V}^{MIS}].$

Similar to known concepts, when the covariance assumption is satisfied, even unknown concept-based estimators can provide lower variances than traditional and MIS estimators. In known concepts however, this assumption has to be satisfied by the practitioner, whereas in unknown concepts, this assumption can be used as a loss function in our methodology to implicitly reduce variance. (Line 12)

Theorem 5.4 (Confidence bounds for Concept-based estimators). *The Cramer-Rao bound on the Mean-Square Error of CIS and CPDIS estimator loosen by* $\epsilon(|\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}]|^2)$, under unknown concepts over known-concepts. Here, $\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}]$ is the on-policy estimate of concept-based IS (PDIS) estimator.

The confidence bounds of unknown concepts mirror that of known-concepts, with the addition of the bias term whose maximum value is the true on-policy estimate of the estimator. This is typically unknown in real-world scenarios and requires additional domain knowledge to mitigate.

5.2 Experimental setup

Environments, Policy descriptions, Metrics: Same as those in known concepts section.

Concept representation: In both examples, we use a 4-dimensional concept $c_t \in \mathcal{R}^4$, where each sub-concept is a linear weighted function of human-interpretable features f, i.e., $c_t^i = w \cdot f(s_t)$, with w optimized as previously discussed. Detailed descriptions of the features and optimized concepts after CBM training are provided in Supplement I. For MIMIC, features f are normalized vital signs, as threshold information for discretization is unavailable. In brevity of space, we move the training and hyperparameter details to Supplement H.

5.3 Results and Discussion

Learned concepts using Algorithm 1 yield improvements across all metrics except bias compared to <u>traditional OPE estimators</u>. Significant improvements in variance, MSE, and ESS are observed for the Windy Gridworld and MIMIC datasets, with gains of 1-2 and 2-3 orders of magnitude, respectively. This improvement is due to our algorithm's ability to identify concepts that satisfy the desiderata, including achieving variance reduction as specified in line 12 of the algorithm. However, like known concepts, optimized concepts show a higher bias than traditional estimators. This is because, unlike variance, bias cannot be optimized in the loss function without the true on-policy estimate, which is typically unavailable in real-world settings. As a result, external information may be essential for further bias reduction.

Learned concepts yield improvements across all metrics besides bias over known concept estimators. From Fig 3a,3b, we observe that our methodology improves variance, MSE, and ESS by 1–2 orders of magnitude compared to known concepts. This suggests our algorithm can learn concepts that outperform human-defined ones in OPE metrics. Fig 3c further supports this, showing a lower frequency of high IPS scores for learned concepts than for known ones. This indicates our algorithm discovers novel concepts that satisfy concept desiderata in Section B while enhancing OPE characteristics, particularly variance. Such capability is valuable in domains with imperfect experts or complex real-world settings where perfect expertise is unattainable. However, these learned concepts introduce higher bias, as the training algorithm prioritizes variance reduction over bias minimization. This could be mitigated by regularizing variance during training.

Learned concepts are interpretable, show conciseness and diversity. We list the optimized concepts in Supplement I. These concepts exhibit sparse weights, enhancing their conciseness, with significant variation in weights across different dimensions of the concepts, reflecting diversity. This work focuses on linearly varying concepts, but more complex concepts, such as symbolic representations (Majumdar et al., 2023), could better model intricate environments.

6 Interventions on Concepts for Insights on Evaluation

Concepts provide interpretations, allowing practitioners to identify sources of variance—an advantage over traditional state abstractions like Pavse & Hanna (2022a). Concepts also clarify reasons behind OPE characteristics, such as high variance, enabling corrective interventions based on domain knowledge or human evaluation. We outline the details of performing interventions next.

6.1 Methodology

Given trajectory history h_t and concept c_t , we define c_t^{int} as the intervention (alternative) concept an expert proposes at time t. We define criteria $\kappa : (h_t, c_t) \to \{0, 1\}$ as a function constructed from domain expertise that takes in (h_t, c_t) as input and outputs a boolean value. This criteria function determines whether an intervention needs to be conducted over the current concept c_t or not. For e.g., if a practitioner has access to true on-policy values, he/she can estimate which concepts suffer from bias. If a concept doesn't suffer from bias, the criteria $\kappa(h_t, c_t) = 1$ is satisfied and the concept is not intervened upon, else $\kappa(h_t, c_t) = 0$ and the intervened concept c_t^{int} is used instead. The final concept \tilde{c}_t is then defined as: $\tilde{c}_t = \kappa(h_t, c_t) \cdot c_t + (1 - \kappa(h_t, c_t)) \cdot c_t^{int}$. Under the absence of true on-policy values, the practitioner may chose to intervene using a different criteria instead.

We define criteria κ for our experiments as follows. In gridworld, we assume access to oracle concepts, listed in Supp. Material H. When the learned concept c_t matches the true concept, $\kappa(h_t, c_t) = 1$, otherwise 0. In MIMIC, the interventions are based on a patient's urine output at a specific timestep with $\kappa(h_t, c_t) = 1$ when urine output > 30 ml/hr, and 0 otherwise. Performing interventions based on urine output enables us to assess the role of kidney function in hypotension management. In this work, we consider 3 possible intervention strategies either based on states or domain knowledge.

Interventions that replace concepts with state representations and state-based policies. We intervene on the concept with the state and use policies dependent on state to perform OPE, i.e $c_t^{\text{int}} = s_t$, $\pi_e^c(a_t | \tilde{c}_t) = \pi_e(a_t | s_t)$, $\pi_b^c(a_t | \tilde{c}_t) = \pi_b(a_t | s_t)$. This can be thought of as a comparative measure a practitioner can look for between the concept and the state representations.



Figure 4: Interpretations of learned concepts. *Windy Gridworld:* Fig 4a and 4b compare true oracle concepts with learned concepts derived from the proposed methodology. We observe a deviation between learned and oracle concepts (circled in red), identifying potential interventions. We compare our learnt concepts with a state abstraction baseline in Fig 4c obtained using K-means clustering. We observe the clusters to significantly differ from both oracle and optimized concepts, underscoring the meaningfulness of learned concepts. *MIMIC:* From Fig 4d, we observe patients with low urine output exhibit greater variance in learned concepts compared to high-output patients, revealing potential intervention targets.

Interventions that replace concepts with state representations and maximum likelihood estimator (MLE) of state-based policies. We replace the erroneous concept with the corresponding state and use the MLE of the state conditioned policy to perform OPE, i.e $c_t^{\text{int}} = s_t$, $\pi_e^c(a_t|\tilde{c}_t) = MLE(\pi_e(a_t|s_t))$, $\pi_b^c(a_t|\tilde{c}_t) = MLE(\pi_b(a_t|s_t))$. This can be thought of as a comparative measure a practitioner can look for between the concept and states, while priortizing over the most confident action.

Interventions using a qualitative concept while retaining concept-based policies. In this approach, a human expert replaces the concept using external domain knowledge. This is similar to Tang & Wiens (2023), where the authors augment the dataset with counterfactual annotations to improve sample efficiency in regions where the coverage is low. However, while Tang & Wiens focus on quantitative counterfactual annotations in the state representation, we employ human interventions to qualitatively edit concepts. In case of gridworld, we consider the oracle concepts as our qualitative concept, while for MIMIC, we consider the learnt CPDIS estimator as qualitative concept while intervening on CIS estimator.

6.2 Results and Interpretations from Interventions on learned concepts

We interpret the optimized concepts in Fig.4. In the gridworld environment, we compare the ground-truth concepts with the optimized ones and observe two additional concepts predicted in the bottom-right region. This likely stems from overfitting to reduce variance in the OPE loss, suggesting a need for inspection and possible intervention. Additionally, we compare our clusters with state-abstraction baseline (clustering in the state-space), and observe the clusters to be widely different from the learnt concepts. In MIMIC, prior studies indicate that patients with urine output above 30 ml/hr are less susceptible to hypotension than those with lower output (Kellum & Prowle, 2018; Singer et al., 2016; Vincent & De Backer, 2013). Using this, we analyze patient trajectories and find that lower urine output correlates with higher variance, while higher output corresponds to lower variance. This insight helps identify patients who may benefit from targeted interventions.

Interpretable concepts allow for targeted interventions that further enhance OPE estimates.

Using qualitative interventions, we observe a reduction in bias and, consequently, MSE in WindyGridworld. This occurs because replacing erroneous concepts with oracle concepts introduces previously missing on-policy information during the optimization of unknown concepts, while preserving the order of variance and ESS estimates. Similarly, in MIMIC, intervening on states with low urine output reduces variance by 1–2 orders of magnitude. This is further supported by the decreased frequency of high IPS scores after intervention, as shown in Fig 5c.



Figure 5: Interventions. (a-b) Qualitative interventions reduce bias in the learned concept estimator in Windy Gridworld and lower variance in MIMIC. In contrast, interventions based on traditional state-based policies $\pi(.|s)$ reduce bias and MSE compared to non-intervened concepts in Gridworld but are outperformed by qualitative interventions. (c) The frequency of high IPS scores decreases after applying qualitative interventions. Furthermore, $\pi(.|s)$ -based interventions exhibit a higher frequency of IPS scores before intervention, indicating that not all strategies are equally effective.

Not all interventions improve Concept OPE characteristics and should be used at the practitioner's discretion. We analyze OPE after applying interventions using traditional state-based policies, $\pi(.|s)$. In gridworld, state-based interventions increase bias and MSE compared to qualitative ones, while in MIMIC, they lead to higher variance. This occurs because traditional state policies (π_b and π_e) fail to compensate for the lack of on-policy information, undermining the advantages of concept-based policies (π_b^c and π_e^c). In contrast, qualitative interventions—oracle concepts in WindyGridworld or urine output thresholds in MIMIC—retain these benefits and effectively address domain-specific challenges. Additionally, as shown in Fig 5c, $\pi(.|s)$ -based interventions result in a higher frequency of both low and high IPS scores. However, since the effect of high IPS scores dominates, OPE variance increases compared to non-intervented OPE. Nevertheless, this framework allows practitioners to inspect and select among alternative interventions as needed.

7 Conclusions, Limitations and Future Work

We introduced a new family of concept-based OPE estimators, demonstrating that known-concept estimators can outperform traditional ones with greater accuracy and theoretical guarantees. For unknown concepts, we proposed an algorithm to learn interpretable concepts that improve OPE evaluations by identifying performance issues and enabling targeted interventions to reduce variance. These advancements benefit safety-critical fields like healthcare, education, and public policy by supporting reliable, interpretable policy evaluations. By reducing variance and providing policy insights, this approach enhances informed decision-making, facilitates personalized interventions, and refines policies before deployment for greater real-world effectiveness. A limitation of our work is trajectory distribution mismatch when learning unknown concepts, particularly in low-sample settings, which can lead to high-variance OPE. Targeted interventions help mitigate this issue. We also did not address hidden confounding variables or potential CBM concept leakage, focusing instead on evaluation. Future work will address these challenges and extend our approach to more general, partially observable environments.

References

- Panagiotis Anagnostou, Petros T. Barmbas, Aristidis G. Vrahatis, and Sotiris K. Tasoulis. Approximate knn classification for biomedical data, 2020. URL https://arxiv.org/abs/2012. 02149.
- David M. Bossens and Philip S. Thomas. Low variance off-policy evaluation with state-based importance sampling, 2024. URL https://arxiv.org/abs/2212.03932.
- Markus Böck, Julien Malle, Daniel Pasterk, Hrvoje Kukina, Ramin Hasani, and Clemens Heitzinger. Superhuman performance on sepsis mimic-iii data by distributional reinforcement learning. *PLOS ONE*, 17:e0275358, 11 2022. DOI: 10.1371/journal.pone.0275358.
- Yinlam Chow, Marek Petrik, and Mohammad Ghavamzadeh. Robust policy optimization with baseline guarantees. *arXiv preprint arXiv:1506.04514*, 2015.
- Shayan Doroudi, Philip S Thomas, and Emma Brunskill. Importance sampling for fair policy selection. *Grantee Submission*, 2017.
- Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. More robust doubly robust off-policy evaluation. In *International Conference on Machine Learning*, pp. 1447–1456. PMLR, 2018.
- Raphael Fonteneau, Susan A Murphy, Louis Wehenkel, and Damien Ernst. Batch mode reinforcement learning based on the synthesis of artificial trajectories. *Annals of operations research*, 208:383– 416, 2013.
- Scott Fujimoto, David Meger, and Doina Precup. A deep reinforcement learning approach to marginalized importance sampling with the successor representation, 2023.
- Ary L. Goldberger, Luis A. Nunes Amaral, L Glass, Jeffrey M. Hausdorff, Plamen Ch. Ivanov, Roger G. Mark, Joseph E. Mietus, George B. Moody, Chung-Kang Peng, and Harry Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *Circulation*, 101 23:E215–20, 2000. URL https://api. semanticscholar.org/CorpusID:642375.
- Omer Gottesman, Joseph Futoma, Yao Liu, Sonali Parbhoo, Leo Anthony Celi, Emma Brunskill, and Finale Doshi-Velez. Interpretable off-policy evaluation in reinforcement learning by highlighting influential transitions, 2020.
- Deepak Gupta, Russell Loane, Soumya Gayen, and Dina Demner-Fushman. Medical image retrieval via nearest neighbor search on pre-trained image features, 2022. URL https://arxiv.org/abs/2210.02401.
- Josiah Hanna, Peter Stone, and Scott Niekum. Bootstrapping with models: Confidence intervals for off-policy evaluation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Marton Havasi, Sonali Parbhoo, and Finale Doshi-Velez. Addressing leakage in concept bottleneck models. *Advances in Neural Information Processing Systems*, 35:23386–23397, 2022.
- Daniel G Horvitz and Donovan J Thompson. A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260):663–685, 1952.
- Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, STOC '98, pp. 604–613, New York, NY, USA, 1998. Association for Computing Machinery. ISBN 0897919629. DOI: 10.1145/276698.276876. URL https://doi.org/10.1145/276698. 276876.

- Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pp. 652–661. PMLR, 2016.
- Alistair E W Johnson, Tom J Pollard, Lu Shen, Li-Wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3:160035, May 2016. ISSN 2052-4463. DOI: 10.1038/sdata.2016.35. URL https://europepmc.org/articles/PMC4878278.
- Pulkit Katdare, Nan Jiang, and Katherine Rose Driggs-Campbell. Marginalized importance sampling for off-environment policy evaluation. In Jie Tan, Marc Toussaint, and Kourosh Darvish (eds.), *Proceedings of The 7th Conference on Robot Learning*, volume 229 of *Proceedings of Machine Learning Research*, pp. 3778–3788. PMLR, 06–09 Nov 2023. URL https://proceedings. mlr.press/v229/katdare23a.html.
- John A. Kellum and John R. Prowle. Acute kidney injury in the critically ill: Clinical epidemiology and outcomes. *Nature Reviews Nephrology*, 14(10):641–656, 2018. DOI: 10.1038/ s41581-018-0052-0.
- Ramtin Keramati, Omer Gottesman, Leo Anthony Celi, Finale Doshi-Velez, and Emma Brunskill. Identification of subgroups with similar benefits in off-policy policy evaluation. *CoRR*, abs/2111.14272, 2021a. URL https://arxiv.org/abs/2111.14272.
- Ramtin Keramati, Omer Gottesman, Leo Anthony Celi, Finale Doshi-Velez, and Emma Brunskill. Identification of subgroups with similar benefits in off-policy policy evaluation. *arXiv preprint arXiv:2111.14272*, 2021b.
- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5338–5348. PMLR, 13–18 Jul 2020. URL https://proceedings. mlr.press/v119/koh20a.html.
- Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature medicine*, 24(11):1716—1720, November 2018a. ISSN 1078-8956. DOI: 10.1038/s41591-018-0213-5. URL https://doi.org/10.1038/s41591-018-0213-5.
- Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature medicine*, 24(11):1716–1720, 2018b.
- Matthieu Komorowski, Leo Anthony Celi, Omar Badawi, Anthony C. Gordon, and A. Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature Medicine*, 24:1716–1720, 2018c. DOI: 10.1038/s41591-018-0213-5. URL https: //doi.org/10.1038/s41591-018-0213-5.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinitehorizon off-policy estimation. *Advances in neural information processing systems*, 31, 2018a.
- Yao Liu, Omer Gottesman, Aniruddh Raghu, Matthieu Komorowski, Aldo A Faisal, Finale Doshi-Velez, and Emma Brunskill. Representation balancing mdps for off-policy policy evaluation. Advances in Neural Information Processing Systems, 31, 2018b.
- Yao Liu, Pierre-Luc Bacon, and Emma Brunskill. Understanding the curse of horizon in off-policy evaluation via conditional importance sampling, 2020. URL https://arxiv.org/abs/1910.06508.
- Yao Liu, Yannis Flet-Berliac, and Emma Brunskill. Offline policy optimization with eligible actions, 2022. URL https://arxiv.org/abs/2207.00632.

- Pedro Madeira, André Carreiro, Alex Gaudio, Luís Rosado, Filipe Soares, and Asim Smailagic. Zebra: Explaining rare cases through outlying interpretable concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3781–3787, 2023.
- Anita Mahinpei, Justin Clark, Isaac Lage, Finale Doshi-Velez, and Weiwei Pan. Promises and pitfalls of black-box concept learning models. arXiv preprint arXiv:2106.13314, 2021.
- Ritam Majumdar, Vishal Jadhav, Anirudh Deodhar, Shirish Karande, Lovekesh Vig, and Venkataramana Runkana. Symbolic regression for pdes using pruned differentiable programs, 2023. URL https://arxiv.org/abs/2303.07009.
- Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. In *AAMAS*, volume 1077, 2014.
- Ričards Marcinkevičs, Sonia Laguna, Moritz Vandenhirtz, and Julia E Vogt. Beyond concept bottleneck models: How to make black boxes intervenable? *arXiv preprint arXiv:2401.13544*, 2024.
- Anton Matsson and Fredrik D. Johansson. Case-based off-policy policy evaluation using prototype learning, 2021. URL https://arxiv.org/abs/2111.11113.
- S A Murphy, M J van der Laan, J M Robins, and Conduct Problems Prevention Research Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 96 (456):1410–1423, 2001. DOI: 10.1198/016214501753382327. URL https://doi.org/10.1198/016214501753382327. PMID: 20019887.
- Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections, 2019.
- Cosmin Paduraru. Off-policy evaluation in Markov decision processes. PhD thesis, 2013.
- Konstantinos P Panousis, Dino Ienco, and Diego Marcos. Hierarchical concept discovery models: A concept pyramid scheme. arXiv preprint arXiv:2310.02116, 2023.
- Brahma S. Pavse and Josiah P. Hanna. Scaling marginalized importance sampling to high-dimensional state-spaces via state abstraction, 2022a.
- Brahma S. Pavse and Josiah P. Hanna. Scaling marginalized importance sampling to high-dimensional state-spaces via state abstraction, 2022b. URL https://arxiv.org/abs/2212.07486.
- Achim Peine, Andreas Hallawa, Jan Bickenbach, Peter Sidler, Andreas Markewitz, Alexandre Levesque, Jeremy Levesque, and Nils Haake. Development and validation of a reinforcement learning algorithm to dynamically optimize mechanical ventilation in critical care. *npj Digital Medicine*, 4(1):32, 2021. DOI: 10.1038/s41746-021-00388-6.
- Doina Precup, Richard S. Sutton, and Satinder P. Singh. Eligibility traces for off-policy policy evaluation. In *Proceedings of the Seventeenth International Conference on Machine Learning*, ICML '00, pp. 759–766, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc. ISBN 1558607072.
- Yoshihide Sawada and Keigo Nakamura. Concept bottleneck model with additional unsupervised concepts. *IEEE Access*, 10:41758–41765, 2022.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
- Simon P Shen, Yecheng Ma, Omer Gottesman, and Finale Doshi-Velez. State relevance for off-policy evaluation. In *International Conference on Machine Learning*, pp. 9537–9546. PMLR, 2021.

- Mervyn Singer, Clifford S. Deutschman, Christopher W. Seymour, et al. The third international consensus definitions for sepsis and septic shock (sepsis-3). *JAMA*, 315(8):801–810, 2016. DOI: 10.1001/jama.2016.0287.
- Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
- Shengpu Tang and Jenna Wiens. Counterfactual-augmented importance sampling for semi-offline policy evaluation, 2023. URL https://arxiv.org/abs/2310.17146.
- Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, pp. 2139–2148. PMLR, 2016.
- Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and q-function learning for off-policy evaluation, 2020.
- Jean-Louis Vincent and Daniel De Backer. Circulatory shock. *New England Journal of Medicine*, 369(18):1726–1734, 2013. DOI: 10.1056/NEJMra1208943.
- Carissa Wu, Sonali Parbhoo, Marton Havasi, and Finale Doshi-Velez. Learning optimal summaries of clinical time-series with concept bottleneck models. In *Machine Learning for Healthcare Conference*, pp. 648–672. PMLR, 2022.
- Tengyang Xie, Yifei Ma, and Yu-Xiang Wang. Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. *Advances in neural information processing systems*, 32, 2019.
- Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estimation of stationary values, 2020a.
- Shangtong Zhang, Bo Liu, and Shimon Whiteson. Gradientdice: Rethinking generalized offline estimation of stationary values, 2020b.

Supplementary Material

A Related Work

Off-Policy Evaluation. There is a long history of methods for performing OPE, broadly categorized into model-based or model-free (Sutton & Barto, 2018). Model-based methods, such as the Direct Method (DM), learn a model of the environment to simulate trajectories and estimate the policy value (Paduraru, 2013; Chow et al., 2015; Hanna et al., 2017; Fonteneau et al., 2013; Liu et al., 2018b). These methods often rely on strong assumptions about the parametric model for statistical guarantees. Model-free methods, like IS, correct sampling bias in off-policy data through reweighting to obtain unbiased estimates (e.g., Precup et al. (2000); Horvitz & Thompson (1952); Thomas & Brunskill (2016)). Doubly robust (DR) estimators (e.g., Jiang & Li (2016); Farajtabar et al. (2018)) combine model-based DM and model-free IS for OPE but may fail to reduce variance when both DM and IS have high variance. Various methods have been developed to refine estimation accuracy in IS, such as truncating importance weights and estimating weights from steady-state visitation distributions (Liu et al., 2018a; Xie et al., 2019; Doroudi et al., 2017; Bossens & Thomas, 2024).

Off-Policy Evaluation based on Subgroups. Keramati et al. (2021b) extend OPE to estimate treatment effects for subgroups and provide actionable insights on which subgroups may benefit from specific treatments, assuming subgroups are known or identified using regression trees. Unlike regression trees, which are limited in scalability, our approach learns interpretable concepts to characterize individuals, on the basis of which we introduce a new family of IS estimators. Similarly, Shen et al. (2021) propose reducing variance by omitting likelihood ratios for certain states. Our work complements this by summarizing relevant trajectory information using concepts, rather than explicitly omitting states irrelevant to the return. The advantage of using concepts as opposed to states is that we can easily interpret and intervene on these concepts unlike the states.

Marginalized Importance Sampling (MIS) estimators (Uehara et al., 2020; Liu et al., 2018a; Nachum et al., 2019; Zhang et al., 2020b;a) mitigate the high variance of traditional IS by reweighting data tuples using density ratios computed from the state visitation at each time step. These estimators enhance robustness by focusing on states with high visitation density ratios, thereby marginalizing out less visited states. However, MIS has its challenges: computing density ratios can introduce high variance, particularly in complex state spaces, and it obscures which aspects of the state space contribute directly to variance. Some studies, such as Katdare et al. (2023) and Fujimoto et al. (2023), improve MIS by decomposing density ratio estimation into components like large density ratio mismatch and transition probability mismatch. Our work differs from MIS by characterizing trends in a trajectory using interpretable concepts rather than solely relying on density ratios. This approach enables targeted interventions on specific concepts of interest, leading to more accurate return estimates and reduced variance in OPE. Unlike MIS, our method provides interpretability, which becomes increasingly important as problem complexity grows. Proposals for hybrid estimators, such as those in Payse & Hanna (2022a), suggest using low-dimensional abstractions of state spaces with MIS to manage high-dimensional spaces more effectively. Our work differs in the sense that we use concepts instead of state abstractions which can be easily plugged into the existing IS OPE definitions as elaborated in Sections 3.2 and Supplement D.

Concept Bottleneck Models. Concept Bottleneck Models (CBMs) (Koh et al., 2020) are a class of prediction models that first predict a set of human interpretable concepts, and subsequently use these concepts to predict a downstream label. Variations of these models include learning soft probabilistic concepts (Mahinpei et al., 2021), learning hierarchical concepts (Panousis et al., 2023) and learning concepts in a semi-supervised manner (Sawada & Nakamura, 2022). The key advantage of these models is they allow us to explicitly intervene on concepts and interpret what might happen to a downstream label if certain concepts were changed (Marcinkevičs et al., 2024). Unlike previous works, we leverage this idea to introduce a new class of estimators for OPE where we group trajectories based on interpretable concepts which are relevant for the downstream evaluation task.

B Concept Desiderata

Explainability: Explainability ensures that the concept function ϕ is composed of humaninterpretable functions f_1, f_2, \ldots, f_n . Each interpretable function f_i depends on the current state, past actions, rewards, and states, i.e., $s_t, a_{0:t-1}, r_{0:t-1}, s_{0:t-1}$. Mathematically:

$$c_t = \phi(s_t, a_{0:t-1}, r_{0:t-1}, s_{0:t-1}) = \psi(f_1(s_t, a_{0:t-1}, r_{0:t-1}, s_{0:t-1}), \dots, f_n(s_t, a_{0:t-1}, r_{0:t-1}, s_{0:t-1}))$$
(1)

Here, ψ maps the human-interpretable functions f_i to the concept c_t , and both ϕ and ψ share the same co-domain space C. In essence, ϕ can be defined using a single interpretable function or a combination of multiple interpretable functions.

As a running example in this paper (applicable across domains), the concept function $\phi(s_t)$ for diagnosing hypertension can be expressed using human-interpretable features:

$$\begin{split} c_t &= \phi(s_t) \\ &= \phi(\text{SBP}, \text{DBP}, \text{HR}, \text{Glucose levels}, \text{GCS}, \text{Age}, \text{Weight}) \\ &= \psi(f_1(\text{SBP}), f_2(\text{DBP}), f_3(\text{HR}), f_4(\text{Glucose levels}), f_5(\text{GCS}), f_6(\text{Age}, \text{Weight})) \end{split}$$

Where:

- $f_1(SBP)$ maps Systolic Blood Pressure to a category (e.g., Low, Normal, High).
- $f_2(DBP)$ maps Diastolic Blood Pressure to a category (e.g., Low, Normal, High).
- $f_3(HR)$ maps Heart Rate to a category (e.g., Low, Normal, High).
- f_4 (Glucose levels) maps blood glucose levels to a category (e.g., Low, Normal, High).
- $f_5(GCS)$ maps GCS scores to a category.
- f_6 (Age, Weight) maps age and weight to Body Mass Index (BMI).

This ensures that the concept $\phi(s_t)$ for diagnosing hypertension is built from human-interpretable features, making the diagnostic process explainable. Each function f_i translates raw medical data into intuitive categories that are meaningful to medical practitioners.

Conciseness: Conciseness ensures that the concept function ϕ represents the minimal mapping of interpretable functions f_1, f_2, \ldots, f_n to the concept c_t . If multiple mappings $\psi_1, \psi_2, \ldots, \psi_m$ satisfy ϕ , we choose the mapping ψ that provides the simplest composition of f_i to describe c_t .

E.g. Obesity can be represented by different combinations of human-interpretable functions. We select the least complex representation that remains interpretable. The two possible representations are:

$$c_t = \psi_1(f_1(\text{height}), f_2(\text{weight}), f_3(\text{SBP}), f_4(\text{DBP}))$$

$$c_t = \psi_2(f_5(\text{BMI}), f_3(\text{SBP}))$$

Since BMI encapsulates both height and weight, and either SBP or DBP accurately summarizes blood pressure pertinent to Obesity, the concept $c_t = \psi_2(f_5(BMI), f_3(SBP))$ is more concise.

Better Trajectory Coverage: Concept-based policies have a higher coverage than traditional state policies. Mathematically:

$$\sum_{\tau \in \mathcal{T}_1} \sum_{t=0}^T \pi^c(a_t | c_t) \ge \sum_{\tau \in \mathcal{T}_2} \sum_{t=0}^T \pi(a_t | s_t)$$
(2)

Here, π^c , π represent policies conditioned on concepts and states respectively, \mathcal{T}_1 , \mathcal{T}_2 is the set of all possible trajectories under π^c , π and T is the total number of timesteps.

Diversity: The diversity property ensures that each dimension of the concept at a given timestep captures distinct and independent aspects of the state space, minimizing overlap.

As an example, the concept function $\phi(s_t)$ for a comprehensive patient health assessment can be represented as:

$$\begin{split} \phi(s_t) &= [c_t^1, c_t^2, \dots, c_t^d] \\ &= [c_t^1(\text{Cardiovascular Health}), c_t^2(\text{Metabolic Health}), c_t^3(\text{Respiratory Health})] \\ &= [\psi_1(f_1(\text{blood pressure}), f_2(\text{cholesterol levels}), f_3(\text{heart rate variability})), \\ &\psi_2(f_1(\text{blood glucose levels}), f_2(\text{BMI}), f_3(\text{metabolic history})), \\ &\psi_3(f_1(\text{lung function}), f_2(\text{oxygen saturation}), f_3(\text{respiratory history}))] \end{split}$$

Each dimension of the concept c_t^i captures unique information, contributing to a holistic assessment of the patient's health without redundancy.

C Choice of Concept Types

Concepts capturing subgroups with short-term benefits. If ϕ maps state s_t and action a_t to immediate reward r_t , the resulting concepts can identify subgroups with similar short-term benefits, facilitating more personalized OPE, as seen in Keramati et al. (2021a). Unlike Keramati et al. (2021b), we do not limit ϕ to a regression tree.

Concepts capturing high-variance transitions. If ϕ highlights changes in state s_t and action a_t that cause significant shifts in value estimates, it can capture influential transitions or dynamics from historical data, similar to Gottesman et al. (2020).

Concepts capturing least influential states. If ϕ identifies the least (or most) influential states s_t , it can help focus more on critical states, reducing variance by only applying IS ratios to those states Bossens & Thomas (2024).

Concepts capturing state-density information. If ϕ extracts information from histories to predict state-action visitation counts, concept-based OPE with ϕ functions similarly to Marginalized OPE estimators, like Xie et al. (2019), which reweight trajectories based on state-visitation distributions. However, density-based concepts may be less interpretable and harder to intervene in the context of OPE.

D Generalized Concept-based OPE estimators

Building on the OPE estimators discussed in the main paper, we extend the integration of concepts into other popular OPE estimators. Without making any additional assumptions about the estimators' definitions, concepts can be seamlessly incorporated into the original formulations of these estimators.

Definition D.1 (Concept-based Weighted Importance Sampling, CWIS).

$$\hat{V}_{\pi_{e}}^{CWIS} = \frac{\sum_{n=1}^{N} \rho_{0:T}^{(n)} \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)}}{\sum_{n=1}^{N} \rho_{0:T}^{(n)}}; \quad \rho_{0:T}^{(n)} = \prod_{t'=0}^{T} \frac{\pi_{e}(a_{t'}^{(n)}|c_{t'}^{(n)})}{\pi_{b}(a_{t'}^{(n)}|c_{t'}^{(n)})};$$

Definition D.2 (Concept-based Per-Decision Weighted Importance Sampling, CPDWIS).

$$\hat{V}_{\pi_e}^{CPDWIS} = \frac{\sum_{n=1}^{N} \sum_{t=0}^{T} \rho_{0:t}^{(n)} \gamma^t r_t^{(n)}}{\sum_{n=1}^{N} \rho_{0:T}^{(n)}}; \quad \rho_{0:t}^{(n)} = \prod_{t'=0}^{t} \frac{\pi_e(a_{t'}^{(n)} | c_{t'}^{(n)})}{\pi_b(a_{t'}^{(n)} | c_{t'}^{(n)})}$$

Definition D.3 (Concept-based Doubly Robust Estimator, CDR).

$$\hat{V}_{CDR} = \frac{1}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \prod_{k=0}^{t} \frac{\pi_{e}(a_{k}^{(i)} \mid c_{k}^{(i)})}{\pi_{b}(a_{k}^{(i)} \mid c_{k}^{(i)})} \left(r_{t}^{(i)} - \hat{Q}(s_{t}^{(i)}, a_{t}^{(i)})\right) + \hat{V}(s_{t}^{(i)})$$

Assuming good model-based estimates $\hat{V}(s_t)$, $\hat{Q}(s_t^{(i)}, a_t^{(i)})$, all the advantages seen in the traditional DR estimator translate over to the concept-space representation. It's important to note, the concepts are only used to reweight the Importance Sampling ratios and are not incorporated in the model-based estimates. This allows concepts to have a general form and are not under any markovian assumption, thus satisfying the Bellman equation.

Definition D.4 (Concept-based Marginalized Importance Sampling Estimator, CMIS).

$$\hat{V}_{CMIS} = \sum_{n=1}^{N} \sum_{t=0}^{T} \frac{d_{\pi_{e}^{c}}(c_{t})}{d_{\pi_{b}^{c}}(c_{t})} \gamma^{t} r_{t}$$

Different algorithms from the DICE family attempt to estimate the state-distribution ratio $\frac{d_{\pi_e}(s_t)}{d_{\pi_b}(s_t)}$. MIS in the concept representation accounts for concept-visitation counts. These counts retain all the statistical guarantees of the state representation. However, a drawback is that concept-visitation counts are less intuitive than the original concept definition. This makes it harder to assess the quality of the OPE.

E Known Concept-based OPE Estimators: Theoretical Proofs

In this section, we provide the detailed proofs for the known concept scenario.

Theorem. For any arbitrary function f, $\mathbb{E}_{c \sim d_{\pi}c} f(c) = \mathbb{E}_{s \sim d_{\pi}} f(\phi(s))$

Proof: See Pavse & Hanna (2022b).

E.1 IS

E.1.1 Bias

$$Bias = |\mathbb{E}_{\pi_b^c}[\hat{V}_{\pi_e^c}^{CIS}] - \mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e^c}^{CIS}]|$$
(a)

$$= \left| \mathbb{E}_{\pi_{b}^{c}} \left[\rho_{0:T}^{(n)} \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)} \right] - \mathbb{E}_{\pi_{e}^{c}} [\hat{V}_{\pi_{e}}^{CIS}] \right|$$
(b)

$$= \left|\sum_{n=1}^{N} \left(\prod_{t=0}^{T} \pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)})\right) \rho_{0:T}^{(n)} \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)} - \mathbb{E}_{\pi_{e}}[\hat{V}_{\pi_{e}^{c}}^{CIS}]\right|$$
(c)

$$= \left|\sum_{n=1}^{N} \prod_{t=0}^{T} \left(\pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)}) \frac{\pi_{e}^{c}(a_{t}^{(n)}|c_{t}^{(n)})}{\pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)})} \right) \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}}^{CIS}] \right|$$
(d)

$$= \left|\sum_{n=1}^{N} \prod_{t=0}^{T} \pi_{e}^{c}(a_{t}^{(n)}|c_{t}^{(n)}) \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}^{c}}^{CIS}]\right| = 0$$
(e)

Explanation of steps:

- (a) We start by expressing the definition of Bias as the difference between expected values of the value function sampled under the behavior policy π_b^c and the concept-based evaluation policy $\pi_e^c(a|c)$.
- (b) We expand the respective definitions.
- (c) Each term is expanded to represent the probability of the trajectories, factoring in the importance sampling ratio.
- (d) Grouping similar terms. This change of measure is possible as the concepts are known and can be modify the trajectory probabilities.
- (e) The denominator of the IS term cancels with the probability of the trajectory under π_b^c . Using the definition of $\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e^c}^{CIS}] = \sum_{n=1}^N \prod_{t=0}^T \pi_e^c(a_t^{(n)}|c_t^{(n)}) \sum_{t=0}^T \gamma^t r_t^{(n)}$.

E.1.2 Variance

$$\mathbb{V}[\hat{V}_{\pi_{e}^{c}}^{CIS}] = \mathbb{E}_{\pi_{b}^{c}}[(\hat{V}_{\pi_{e}^{c}}^{CIS})^{2}] - (\mathbb{E}_{\pi_{b}^{c}}[\hat{V}_{\pi_{e}^{c}}^{CIS}])^{2}$$
(a)

We first evaluate the expectation of the square of the estimator:

$$\mathbb{E}_{\pi_{b}^{c}}[(\hat{V}_{\pi_{e}}^{CIS})^{2}] = \mathbb{E}_{\pi_{b}^{c}}\left[\left(\rho_{0:T}^{(n)}\sum_{t=0}^{T}\gamma^{t}r_{t}^{(n)}\right)^{2}\right]$$
(b)

$$= \mathbb{E}_{\pi_b^c} \left[\sum_{t=0}^T \sum_{t'=0}^T \rho_{0:T}^2 \gamma^{(t+t')} r_t^{(n)} r_{t'}^{(n')} \right]$$
(c)

$$= \sum_{n=1}^{N} \prod_{t=0}^{T} \frac{(\pi_e^c(a_t^{(n)}|c_t^{(n)}))^2}{\pi_b^c(a_t^{(n)}|c_t^{(n)})} \sum_{t=0}^{T} \sum_{t'=0}^{T} \gamma^{(t+t')} r_t r_{t'}$$
(d)

Evaluating the second term in the variance expression:

$$(\mathbb{E}_{\pi_{b}^{c}}[\hat{V}_{\pi_{c}^{c}}^{CIS}])^{2} = \left(\mathbb{E}_{\pi_{b}^{c}}[\rho_{0:T}^{(n)}\sum_{t=0}^{T}\gamma^{t}r_{t}^{(n)}]\right)^{2}$$
(e)

$$=\sum_{n=1}^{N}\prod_{t=0}^{T}\left(\pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)})(\frac{\pi_{e}^{c}(a_{t}^{(n)}|c_{t}^{(n)})}{\pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)})})\right)^{2}\sum_{t=0}^{T}\sum_{t'=0}^{T}\gamma^{(t+t')}r_{t}r_{t'}$$
(f)

$$=\sum_{n=1}^{N}\prod_{t=0}^{T}\left(\pi_{e}^{c}(a_{t}^{(n)}|c_{t}^{(n)})\right)^{2}\sum_{t=0}^{T}\sum_{t'=0}^{T}\gamma^{(t+t')}r_{t}r_{t'}$$
(g)

Subtracting the squared expectation from the expectation of the squared estimator:

$$\mathbb{V}[\hat{V}_{\pi_{e}^{c}}^{CIS}] = \sum_{n=1}^{N} \prod_{t=0}^{T} \left((\pi_{e}^{c}(a_{t}^{(n)}|c_{t}^{(n)})^{2}(\frac{1}{\pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)})} - 1) \right) \sum_{t=0}^{T} \sum_{t'=0}^{T} \gamma^{(t+t')} r_{t} r_{t'} \tag{h}$$

Explanation of steps:

- (a) We begin with the definition of variance for our estimator.
- (b) We evaluate the first term of the Variance.
- (c),(d) We expand the square of the estimator as the square of a sum of weighted returns.
 - (e) We calculate the square of the expectation of the estimator.
 - (f) We expand this squared expectation.
 - (g) The denominator of the IS ratio cancels with the probability of the trajectory.

E.1.3 Variance comparison between CIS ratios and IS ratios

Theorem.
$$\mathbb{V}[\prod_{t=0}^{T} \frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)}] \le \mathbb{V}[\prod_{t=0}^{T} \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)}]$$

Proof: The proof is similar to Pavse & Hanna (2022b), where we generalize to concepts from state abstractions. Using Lemma E and Assumption 4.1, we can say that:

$$\mathbb{E}_{c \sim d_{\pi^c}} \prod_{t=0}^T \frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)} = \mathbb{E}_{s \sim d_{\pi}} \prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} = 1$$
(a)

Denoting the difference between the two variances as D:

$$D = \mathbb{V}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}\right] - \mathbb{V}\left[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right]$$
(b)
$$= \mathbb{E}_{\pi_{b}}\left[\prod_{t=0}^{T} \left(\frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}\right)^{2}\right] - \left[\mathbb{E}_{\pi_{b}}\prod_{t=0}^{T} \left(\frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}\right)\right]^{2} - \mathbb{E}_{\pi_{b}^{c}}\left[\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)^{2}\right] + \left[\mathbb{E}_{\pi_{b}^{c}}\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)\right]^{2} - \mathbb{E}_{\pi_{b}^{c}}\left[\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)^{2}\right] + \left[\mathbb{E}_{\pi_{b}^{c}}\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)\right]^{2}$$
(b)

$$= \mathbb{E}_{\pi_b} \left[\prod_{t=0}^T \left(\frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} \right)^2 \right] - \mathbb{E}_{\pi_b^c} \left[\prod_{t=0}^T \left(\frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)} \right)^2 \right]$$
(d)

$$=\sum_{s}\prod_{t=0}^{T}\pi_{b}(a_{t}|s_{t})\left[\prod_{t=0}^{T}\left(\frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}\right)^{2}\right]-\sum_{c}\prod_{t=0}^{T}\pi_{b}^{c}(a_{t}|c_{t})\left[\prod_{t=0}^{T}\left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)^{2}\right]$$
(e)

$$=\sum_{c} \left(\sum_{s} \prod_{t=0}^{T} \pi_{b}(a_{t}|s_{t}) [\prod_{t=0}^{T} \left(\frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \right)^{2}] - \prod_{t=0}^{T} \pi_{b}^{c}(a_{t}|c_{t}) [\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})} \right)^{2}] \right)$$
(f)

$$= \sum_{c} \left(\sum_{s} \left[\prod_{t=0}^{T} \left(\frac{\pi_{e}(a_{t}|s_{t})^{2}}{\pi_{b}(a_{t}|s_{t})} \right) \right] - \left[\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|c_{t})^{2}}{\pi_{b}^{c}(a_{t}|c_{t})} \right) \right] \right)$$
(g)

We will analyse the difference of variance for 1 fixed concept and denote it as D':

$$D' = \left[\prod_{t=0}^{T} \left(\frac{\pi_e^c(a_t|c_t)^2}{\pi_b^c(a_t|c_t)}\right)\right] - \left(\sum_s \left[\prod_{t=0}^{T} \left(\frac{\pi_e(a_t|s_t)^2}{\pi_b(a_t|s_t)}\right)\right]\right)$$
(h)

Now, if we can show $D' \ge 0$ for |c|, where |c| is the cardinality of the concept representation, then the difference will always be positive, thus completing our proof. We will use induction to prove $D' \ge 0$ on the total number of concepts from 1 to |c| = n < |S|. Now, our induction statement T(n)to prove is, $D' \ge 0$ where n = |c'|. For n = 1, the statement is trivially true where every concept can be represented as the traditional representation of the state.Our inductive hypothesis states that

$$D' = \left(\left[\prod_{t=0}^{T} \left(\frac{\pi_e^c(a_t|c_t)^2}{\pi_b^c(a_t|c_t)} \right) \right] - \sum_s \left[\prod_{t=0}^{T} \left(\frac{\pi_e(a_t|s_t)^2}{\pi_b(a_t|s_t)} \right) \right] \right) \ge 0$$
(i)

Now, we define $S = \sum_{s} [\prod_{t=0}^{T} \left(\frac{\pi_e(a_t|s_t)^2}{\pi_b(a_t|s_t)} \right)], C = \prod_{t=0}^{T} \pi_e^c(a_t|c_t)^2, C' = \prod_{t=0}^{T} \pi_b^c(a_t|c_t).$ After making the substitutions, we obtain

$$C^2 \le SC' \tag{j}$$

This result holds true for |c| = n as per the induction. Now, we add a new state s_{n+1} to the concept as part of the induction, and obtain the following difference:

$$D' = S \times \frac{\pi_e(a|s_{n+1})^2}{\pi_b(a|s_{n+1})} - \frac{C}{C'} \times \frac{\pi_e(a|s_{n+1})^2}{\pi_b(a|s_{n+1})}$$
(k)

Let $\pi_e(a|s_{n+1}) = X$ and $\pi_b(a|s_{n+1}) = Y$. Substituting, we get:

$$D' = S\frac{X^2}{Y} - \frac{C}{C'}\frac{X^2}{Y} = \frac{(SC' - C)X^2}{C'Y}$$
(1)

D' is minimum when C is maximized, hence we substitute $C \le \sqrt{SC'}$ from the induction hypothesis in the expression

$$D' \le \frac{(SC' - \sqrt{SC'})X^2}{C'Y} \tag{m}$$

As $SC' \ge 0$, the term $SC' - \sqrt{SC'}$ is never negative, leading to $D' \le 0$, since the remaining quantities are always positive. Thus, the induction hypothesis holds, and that concludes the proof.

E.1.4 Variance comparison between CIS and IS estimators

 $\begin{array}{ll} \textbf{Theorem.} \qquad When \qquad Cov(\prod_{t=0}^{t} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t},\prod_{t=0}^{k} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) \qquad \leq \\ Cov(\prod_{t=0}^{k} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t},\prod_{t=0}^{k} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{k}), \ the \ variance \ of \ known \ concept-based \ IS \ estimators \ is \\ lower \ than \ traditional \ estimators, \ i.e. \ \mathbb{V}_{\pi_{b}}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{IS}], \ \mathbb{V}_{\pi_{b}}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{PDIS}]. \end{array}$

Proof: Using Lemma E and Assumption 4.1, we can say that:

$$\mathbb{E}_{c \sim d_{\pi^c}} \left[\sum_{t=0}^T \prod_{t=0}^T \frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)} r_t(c_t, a_t) \right] = \mathbb{E}_{s \sim d_{\pi}} \left[\sum_{t=0}^T \prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_t(s_t, a_t) \right]$$
(a)

The Variance for a single example of a CIS estimator is given by

$$\mathbb{V}[\hat{V}^{CIS}] = \frac{1}{T^2} \left(\sum_{t=0}^T \mathbb{V}[\prod_{t=0}^T \frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)} r_t] + 2\sum_{t< k} Cov(\prod_{t=0}^T \frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)} r_t, \prod_{t=0}^T \frac{\pi_e^c(a_t|c_t)}{\pi_b^c(a_t|c_t)} r_k) \right) \quad (b)$$

The Variance for a single example of a IS estimator is given by

$$\mathbb{V}[\hat{V}^{IS}] = \frac{1}{T^2} \left(\sum_{t=0}^T \mathbb{V}[\prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_t] + 2 \sum_{t< k} Cov(\prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_t, \prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_k) \right) \quad (c)$$

We take the difference between the variances, and note the difference of the covariances is not positive as per the assumption. Hence, if we show the differences of variances per timestep is negative, we complete our proof.

$$\begin{split} D &= \mathbb{V}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}] - \mathbb{V}[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}] \qquad (d) \\ &= \mathbb{E}_{\pi_{b}^{c}}[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}] - [\mathbb{E}_{\pi_{b}^{c}}\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|s_{t})}r_{t}\right)]^{2} - \mathbb{E}_{\pi_{b}}[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)^{2}] + [\mathbb{E}_{\pi_{b}}\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)]^{2} \\ &= \mathbb{E}_{\pi_{b}^{c}}[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}] - \mathbb{E}_{\pi_{b}}[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)^{2}] \qquad (e) \\ &= \sum_{c} \prod_{t=0}^{T} (\pi_{b}^{c}(a_{t}|c_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}\right] - \sum_{s} \prod_{t=0}^{T} (\pi_{b}(a_{t}|s_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}\right] - \left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}\right] \\ &= \sum_{c} \sum_{s \in \phi^{-1}(c)} \prod_{t=0}^{T} (\pi_{b}^{c}(a_{t}|c_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}\right] - \sum_{s} \prod_{t=0}^{T} (\pi_{b}(a_{t}|s_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)^{2}\right] \\ &= \sum_{c} \sum_{s \in \phi^{-1}(c)} \prod_{t=0}^{T} (\pi_{b}^{c}(a_{t}|c_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}\right)^{2}\right] - \sum_{s} \prod_{t=0}^{T} (\pi_{b}(a_{t}|s_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)^{2}\right] \\ &\qquad (g) \\ &\leq R_{max}^{2} \left(\sum_{s \in \phi^{-1}(c)} \prod_{t=0}^{T} (\pi_{b}^{c}(a_{t}|c_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)^{2}\right] - \sum_{s} \prod_{t=0}^{T} (\pi_{b}(a_{t}|s_{t})) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)^{2}\right] \\ &\qquad (h) \end{aligned}$$

The rest of the proof is identical to the previous subsection, wherein we perform induction on the cardinality of the concept for and the term inside the bracket is never positive, thus completing the proof.

E.1.5 Upper Bound on the Variance

$$\mathbb{V}[\hat{V}_{\pi_{b}^{c}}^{CIS}] = \mathbb{E}_{\pi_{b}^{c}}\left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{T} \frac{\pi_{e}(a_{t'}|c_{t'})}{\pi_{b}(a_{t'}|c_{t'})}\right)^{2}\right) - \mathbb{E}_{\pi_{b}^{c}}\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{T} \frac{\pi_{e}(a_{t'}|c_{t'})}{\pi_{b}(a_{t'}|c_{t'})}\right)^{2}$$
(a)

$$\leq \mathbb{E}_{\pi_{b}^{c}}\left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{T} \frac{\pi_{e}(a_{t'}|c_{t'})}{\pi_{b}(a_{t'}|c_{t'})}\right)^{2}\right)$$
(b)

$$\leq \frac{1}{N} \sum_{n=1}^{N} \left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{T} \frac{\pi_{e}(a_{t'}|c_{t'})}{\pi_{b}(a_{t'}|c_{t'})} \right)^{2} \right) + \frac{7T^{2}R_{max}^{2}U_{c}^{2T}ln(\frac{2}{\delta})}{3(N-1)} + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3}-N^{2}}} \sum_{i< j}^{N} (X_{i}^{2}-X_{j}^{2})^{2}}$$
(c)

$$\leq T^2 R_{max}^2 U_c^{2T} (\frac{1}{N} + \frac{\ln \frac{2}{\delta}}{3(N-1)}) + \sqrt{\frac{\ln(\frac{2}{\delta})}{N^3 - N^2}} \sum_{i < j}^N (X_i^2 - X_j^2)^2 \tag{d}$$

Explanation of steps:

- (a) We begin with the definition of variance.
- (b) The second term is always greater than 0
- (c) Applying Bernstein inequality with probability 1- δ . X_i refers to the CIS estimate for 1 sample.
- (d) Grouping terms 1 and 2 together, where $U_c = max \frac{\pi_e^c(a|c)}{\pi_e^c(a|c)}$.

The first term of the variance dominates the second with increase in number of samples. Thus, Variance is of the complexity $\mathcal{O}(\frac{T^2 R_{max}^2 U_c^{2T}}{N})$

E.1.6 Upper Bound on the MSE

$$MSE = Bias^{2} + Variance = Variance \sim \mathcal{O}(\frac{T^{2}R_{max}^{2}U_{c}^{2T}}{N})$$
(a)

The Upper Bound on the MSE of Concept-based IS estimator is of the same form as the Cramer-Rao bounds of the traditional IS estimator as stated in Jiang & Li (2016). We investigate when the MSE bounds can be tightened in the concept representation. We first say,

$$U_{c} = max \frac{\pi_{e}^{c}(a|c)}{\pi_{b}^{c}(a|c)} = U_{s} \frac{K_{1}}{K_{2}}$$
 (b)

Here, $U_s = max \frac{\pi_e(a|s)}{\pi_b(a|s)}$, K_1 is the cardinality of the states which have the same concept c under evaluation policy π_e , while K_2 refers to the same quantity under the behavior policy π_b . Typically, the maximum value of the IS ratio occurs when $\pi_e(a|s) >> \pi_b(a|s)$, i.e. the action taken is very likely under the evaluation policy π_e while it's unlikely under the behavior policy π_b . This typically happens when that particular state has less coverage, or doesn't appear in the data generated by the behavior policy π_b . Under concepts however, similar states are visited and categorized, which improves the information on the state s through c, leading to $K_2 > 1$. On the other hand, as both $\pi_e^c(a|s)$ and $\pi^e(a|s)$ are close to 1, $K_1 = 1$. Thus, $K = \frac{K_1}{K_2} < 1$ and Hence,

$$\mathcal{O}(\frac{T^2 R_{max}^2 U_c^{2T}}{N}) \sim \mathcal{O}(\frac{T^2 R_{max}^2 (U_s K)^{2T}}{N}) \sim \mathcal{O}(\frac{T^2 R_{max}^2 U_s^{2T}}{N}) K^{2T}$$
(3)

Thus, the Concept-based MSE bounds are tightened by a factor of K^{2T} .

E.1.7 Variance comparison with MIS estimator

Theorem. Let ρ be the product of the Importance Sampling ratio in the state space, and d^{π_e}, d^{π_b} be the stationary density ratios. Then,

$$\mathbb{E}(\rho_{0:T}|s_t, a_t) = \frac{d^{\pi_e}(s_t, a_t)}{d^{\pi_b}(s_t, a_t)}$$

Proof: See Liu et al. (2020)

Theorem. Let X_t and Y_t be two sequences of random variables. Then

$$\mathbb{V}(\sum_{t} Y_t) - \mathbb{V}(\sum_{t} \mathbb{E}[Y_t | X_t]) \ge 2 \sum_{t < k} \mathbb{E}[Y_t Y_k] - 2 \sum_{t < k} \mathbb{E}[\mathbb{E}[Y_t | X_t] \mathbb{E}[Y_k | X_k]]$$

Proof: See Liu et al. (2020)

Theorem. When $Cov(\prod_{t=0}^{t} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) \leq Cov(\frac{d^{\pi_{e}}(s_{t},a_{t})}{d^{\pi_{b}}(s_{t},a_{t})}r_{t}, \frac{d^{\pi_{e}}(s_{k},a_{k})}{d^{\pi_{b}}(s_{k},a_{k})}r_{k}), the variance of known CIS estimators is lower than the Variance of MIS estimator, i.e. <math>\mathbb{V}_{\pi_{b}}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{MIS}].$

Proof: We start from the assumption:

$$Cov(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}, \prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) = \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})} \prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}r_{k}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}]\mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}]\mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}]\mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}]\mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|c_{t})}r_{t}$$

$$= \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \prod_{t=0}^{T} \frac{\pi_{e}(a_{k}|s_{k})}{\pi_{b}(a_{k}|s_{k})} K^{2} r_{t} r_{k}\right] - \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} K r_{t}\right] \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} K r_{k}\right] \quad (b)$$

$$\leq \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \prod_{t=0}^{T} \frac{\pi_{e}(a_{k}|s_{k})}{\pi_{b}(a_{k}|s_{k})} r_{t}r_{k}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} r_{t}] \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} r_{k}]$$
(c)

$$\leq \mathbb{E}[(\frac{d^{\pi_{e}}(s_{t},a_{t})}{d^{\pi_{b}}(s_{t},a_{t})})(\frac{d^{\pi_{e}}(s_{k},a_{k})}{d^{\pi_{b}}(s_{k},a_{k})})r_{t}r_{k}] - \mathbb{E}[\frac{d^{\pi_{e}}(s_{t},a_{t})}{d^{\pi_{b}}(s_{t},a_{t})})r_{t}]\mathbb{E}[\frac{d^{\pi_{e}}(s_{k},a_{k})}{d^{\pi_{b}}(s_{k},a_{k})})r_{k}]$$
(d)

Explanation of steps:

(a) We begin with the definition of covariance.

(b),(c) Using the definition of π^c , with K (the ratio of state-space distribution ratio) < 1.

(d) Applying Lemma E.1.7 to both the terms.

Finally, using Lemma E.1.7, substituting $Y_t = \prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_t$ and $X_t = s_t, a_t, r_t$ completes our proof.

E.2 PDIS

E.2.1 Bias

$$Bias = |\mathbb{E}_{\pi_b^c}[\hat{V}_{\pi_e^c}^{CPDIS}] - \mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e^c}^{CPDIS}]|$$
(a)

$$= \left| \mathbb{E}_{\pi_{b}^{c}} \left[\sum_{t=0}^{I} \gamma^{t} \rho_{0:t}^{(n)} r_{t}^{(n)} \right] - \mathbb{E}_{\pi_{e}^{c}} [\hat{V}_{\pi_{e}^{c}}^{CPDIS}] \right|$$
(b)

$$= \left|\sum_{n=1}^{N} \left(\prod_{t=0}^{T} \pi_{b}^{c}(a_{t}^{(n)}|c_{t}^{(n)})\right) \sum_{t=0}^{T} \gamma^{t} \rho_{0:t}^{(n)} r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}^{c}}^{CPDIS}]\right|$$
(c)

$$= \left|\sum_{n=1}^{N}\sum_{t=0}^{T}\gamma^{t}\left(\prod_{t'=0}^{t}\pi_{b}^{c}(a_{t'}^{(n)}|c_{t'}^{(n)})(\frac{\pi_{e}^{c}(a_{t'}^{(n)}|c_{t'}^{(n)})}{\pi_{b}^{c}(a_{t'}^{(n)}|c_{t'}^{(n)})})\right)r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}}^{CPDIS}]\right| = 0 \quad (\mathbf{d})$$

Explanation of steps: Similar to CIS.

E.2.2 Variance

Following the process similar to CIS estimator:

$$\mathbb{V}[\hat{V}_{\pi_{b}^{c}}^{CPDIS}] = \mathbb{E}_{\pi_{b}^{c}}[(\hat{V}_{\pi_{b}^{c}}^{CPDIS})^{2}] - (\mathbb{E}_{\pi_{b}^{c}}[\hat{V}_{\pi_{b}^{c}}^{CPDIS}])^{2}$$
(a)

We first evaluate the expectation of the square of the estimator:

$$\mathbb{E}_{\pi_b^c}[(\hat{V}_{\pi_b^c}^{CPDIS})^2] = \mathbb{E}_{\pi_b^c}\left[\left(\sum_{t=0}^T \gamma^t \rho_{0:t} r_t\right)^2\right] \tag{b}$$

$$= \mathbb{E}_{\pi_b^c} \left[\sum_{t=0}^T \sum_{t'=0}^T \rho_{0:t} \rho_{0:t'} \gamma^{(t+t')} r_t r_{t'} \right]$$
(c)

$$=\sum_{n=1}^{N}\sum_{t=0}^{T}\sum_{t'=0}^{T}\left(\prod_{t''=0}^{t}\pi_{b}^{c}(a_{t}|c_{t})(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})})\right)\left(\prod_{t'''=0}^{t'}(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})})\right)\gamma^{(t+t')}r_{t}r_{t'} \tag{d}$$

$$=\sum_{n=1}^{N}\sum_{t=0}^{T}\sum_{t'=0}^{T}\left(\prod_{t''=0}^{t}\pi_{e}^{c}(a_{t}|c_{t})\right)\left(\prod_{t'''=0}^{t'}\left(\frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}\right)\right)\gamma^{(t+t')}r_{t}r_{t'}$$
(e)

Evaluating the second term in the variance expression:

$$(\mathbb{E}_{\pi_{b}^{c}}[\hat{V}_{\pi_{b}^{c}}^{CPDIS}])^{2} = \left(\sum_{n=1}^{N}\sum_{t=0}^{T}\mathbb{E}_{\pi_{b}^{c}}[\gamma^{t}\rho_{0:t}r_{t}]\right)^{2}$$
(f)
$$= \sum_{n=1}^{N}\sum_{t=0}^{T}\sum_{t'=0}^{T}\left(\prod_{t''=0}^{t}\pi_{b}^{c}(a_{t''}|c_{t''})(\frac{\pi_{e}^{c}(a_{t''}|c_{t''})}{\pi_{b}^{c}(a_{t''}|c_{t''})})\right) \left(\prod_{t'''=0}^{t'}\pi_{b}^{c}(a_{t'''}|s_{t'''})(\frac{\pi_{e}^{c}(a_{t'''}|c_{t'''})}{\pi_{b}^{c}(a_{t'''}|c_{t'''})})\right) \gamma^{(t+t')}r_{t}r_{t'}$$
(g)
$$= \sum_{n=1}^{N}\sum_{t=0}^{T}\sum_{t'=0}^{T}\left(\prod_{t''=0}^{t}\pi_{e}^{c}(a_{t''}|c_{t''})\right) \left(\prod_{t'''=0}^{t'}\pi_{e}^{c}(a_{t'''}|s_{t'''})\right) \gamma^{(t+t')}r_{t}r_{t'}$$
(h)

Subtracting the squared expectation from the expectation of the squared estimator:

$$\mathbb{V}[\hat{V}_{\pi_{b}^{c}}^{CPDIS}] = \sum_{n=1}^{N} \sum_{t=0}^{T} \sum_{t'=0}^{T} \left(\prod_{t'''=0}^{t} \pi_{e}^{c}(a_{t'''}|c_{t'''}) \right) \left(\prod_{t''=0}^{t'} \pi_{e}^{c}(a_{t''}|c_{t''})(\frac{1}{\pi_{b}^{c}(a_{t''}|c_{t''})} - 1) \right) \gamma^{(t+t')} r_{t} r_{t'}$$
(i)

Explanation of steps: Similar to CIS.

E.2.3 Variance comparison between CPDIS ratios and PDIS ratios

Theorem. $\mathbb{V}[\sum_{t=0}^{T} \prod_{t'=0}^{t} \frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})}] \leq \mathbb{V}[\sum_{t=0}^{T} \prod_{t'=0}^{t} \frac{\pi_{e}(a_{t'}|s_{t'})}{\pi_{b}(a_{t'}|s_{t'})}]$

Proof: Similar to CIS estimator.

E.2.4 Variance comparison between CPDIS and PDIS estimators

Theorem. If for any fixed $0 \le t \le k < T$, if

$$Cov(\prod_{t'=0}^{t} \frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})}r_{t}, \prod_{t'=0}^{k} \frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})}r_{k}) \le Cov(\prod_{t'=0}^{t} \frac{\pi_{e}(a_{t'}|s_{t'})}{\pi_{b}(a_{t'}|s_{t'})}r_{t}, \prod_{t=0}^{T} \frac{\pi_{e}(a_{t'}|s_{t'})}{\pi_{b}(a_{t'}|s_{t'})}r_{k})$$

then $\mathbb{V}[\hat{V}^{CPDIS}] \leq \mathbb{V}[\hat{V}^{PDIS}].$

Proof: Similar to CIS estimator.

E.2.5 Upper Bound on the Variance

$$\mathbb{V}[\hat{V}_{\pi_{b}^{c}}^{CPDIS}] = \mathbb{E}_{\pi_{b}^{c}}\left(\left(\sum_{t=0}^{T}\gamma^{t}r_{t}\prod_{t'=0}^{t}\frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})}\right)^{2}\right) - \mathbb{E}_{\pi_{b}^{c}}\left(\sum_{t=0}^{T}\gamma^{t}r_{t}\prod_{t'=0}^{t}\frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})}\right)^{2} \qquad (a)$$
$$\leq \mathbb{E}_{\pi_{b}^{c}}\left(\left(\sum_{t=0}^{T}\gamma^{t}r_{t}\prod_{t'=0}^{t}\frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})}\right)^{2}\right) \qquad (b)$$

$$\leq \frac{1}{N} \sum_{n=1}^{N} \left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{t} \frac{\pi_{e}^{c}(a_{t'}|c_{t'})}{\pi_{b}^{c}(a_{t'}|c_{t'})} \right)^{2} \right) + \frac{7T^{2}R_{max}^{2}U_{c}^{2T}ln(\frac{2}{\delta})}{3(N-1)} + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3}-N^{2}}} \sum_{i
(c)$$

$$\leq T^2 R_{max}^2 U_c^{2T} \left(\frac{1}{N} + \frac{\ln \frac{2}{\delta}}{3(N-1)}\right) + \sqrt{\frac{\ln(\frac{2}{\delta})}{N^3 - N^2}} \sum_{i < j}^N (X_i^2 - X_j^2)^2 \tag{d}$$

Explanation of steps: Similar to CIS.

E.2.6 Upper Bound on the MSE

$$MSE = Bias^{2} + Variance = Variance \sim \mathcal{O}(\frac{T^{2}R_{max}^{2}U_{c}^{2T}}{N}) \sim \mathcal{O}(\frac{T^{2}R_{max}^{2}U_{s}^{2T}}{N})K^{2T}$$
(4)

Proof: Similar to CIS estimator.

E.2.7 Variance comparison with MIS estimator

Theorem. When $Cov(\prod_{t=0}^{t} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) \leq Cov(\frac{d^{\pi_{e}}(s_{t},a_{t})}{d^{\pi_{b}}(s_{t},a_{t})}r_{t}, \frac{d^{\pi_{e}}(s_{k},a_{k})}{d^{\pi_{b}}(s_{k},a_{k})}r_{k}),$ the variance of known CPDIS estimators is lower than the Variance of MIS estimator, i.e. $\mathbb{V}_{\pi_{b}}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{MIS}].$

Proof: Similar to CIS estimator.

F Unknown Concept-based OPE Estimators: Theoretical Proofs

In this section, we provide the theoretical proofs of the unknown concept scenarios.

F.1 IS

F.1.1 Bias

We begin by stating the expression for the expected value of the CIS estimator under π_b :

$$Bias = |\mathbb{E}_{\pi_b}[\hat{V}_{\pi_e}^{CIS}] - \mathbb{E}_{\pi_e}[\hat{V}_{\pi_e^c}^{CIS}]| \tag{a}$$

$$= |\mathbb{E}_{\pi_b} \left[\rho_{0:T}^{(n)} \sum_{t=0}^T \gamma^t r_t^{(n)} \right] - \mathbb{E}_{\pi_e} [\hat{V}_{\pi_e^c}^{CIS}]|$$
(b)

$$= \left|\sum_{n=1}^{N} \left(\prod_{t=0}^{T} \pi_b(a_t^{(n)} | s_t^{(n)})\right) \rho_{0:T}^{(n)} \sum_{t=0}^{T} \gamma^t r_t^{(n)} - \mathbb{E}_{\pi_e}[\hat{V}_{\pi_e^c}^{CIS}]\right|$$
(c)

$$= \left|\sum_{n=1}^{N} \prod_{t=0}^{T} \left(\pi_b(a_t^{(n)} | s_t^{(n)}) \frac{\pi_e^c(a_t^{(n)} | \tilde{c}_t^{(n)})}{\pi_b^c(a_t^{(n)} | \tilde{c}_t^{(n)})} \right) \sum_{t=0}^{T} \gamma^t r_t^{(n)} - \mathbb{E}_{\pi_e}[\hat{V}_{\pi_e^c}^{CIS}] \right|$$
(d)

Explanation of steps:

- (a) We start by expressing the definition of Bias as the difference between expected values of the value function sampled under the behavior policy π_b and the concept-based evaluation policy $\pi_e(a|c)$.
- (b) We expand the respective definitions.
- (c) Each term is expanded to represent the probability of the trajectories, factoring in the importance sampling ratio.
- (d) Similar terms are grouped together to concisely represent the impact of the importance sampling ratios.

The bias of the CIS estimator is minimum when the concepts \tilde{c}_t equals the traditional state representations s_t , thus, implying imperfect concept-based sampling induces bias. As the concepts are unknown, the reparameterization of the probabilities of the behavior trajectories isn't possible, thus leading to a finite bias as opposed to Known-concept representations.

F.1.2 Variance

We start with the definition of variance for the CIS estimator:

$$\mathbb{V}[\hat{V}_{\pi_{e}}^{CIS}] = \mathbb{E}_{\pi_{b}}[(\hat{V}_{\pi_{e}}^{CIS})^{2}] - (\mathbb{E}_{\pi_{b}}[\hat{V}_{\pi_{e}}^{CIS}])^{2}$$
(a)

We first evaluate the expectation of the square of the estimator:

$$\mathbb{E}_{\pi_{b}}[(\hat{V}_{\pi_{e}}^{CIS})^{2}] = \mathbb{E}_{\pi_{b}}\left[\left(\rho_{0:T}^{(n)}\sum_{t=0}^{T}\gamma^{t}r_{t}^{(n)}\right)^{2}\right]$$
(b)

$$= \mathbb{E}_{\pi_b} \left[\sum_{t=0}^T \sum_{t'=0}^T \rho_{0:T}^2 \gamma^{(t+t')} r_t^{(n)} r_{t'}^{(n')} \right]$$
(c)

$$=\sum_{n=1}^{N}\prod_{t=0}^{T}\left(\pi_{b}(a_{t}^{(n)}|s_{t}^{(n)})(\frac{\pi_{e}^{c}(a_{t}^{(n)}|\tilde{c}_{t}^{(n)})}{\pi_{b}^{c}(a_{t}^{(n)}|\tilde{c}_{t}^{(n)})})^{2}\right)\sum_{t=0}^{T}\sum_{t'=0}^{T}\gamma^{(t+t')}r_{t}r_{t'}$$
(d)

Evaluating the second term in the variance expression:

$$(\mathbb{E}_{\pi_b}[\hat{V}_{\pi_e}^{CIS}])^2 = \left(\mathbb{E}_{\pi_b}[\rho_{0:T}^{(n)}\sum_{t=0}^T \gamma^t r_t^{(n)}]\right)^2 \tag{e}$$

$$=\sum_{n=1}^{N}\prod_{t=0}^{T}\left(\pi_{b}(a_{t}^{(n)}|s_{t}^{(n)})(\frac{\pi_{e}^{c}(a_{t}^{(n)}|\tilde{c}_{t}^{(n)})}{\pi_{b}^{c}(a_{t}^{(n)}|\tilde{c}_{t}^{(n)})})\right)^{2}\sum_{t=0}^{T}\sum_{t'=0}^{T}\gamma^{(t+t')}r_{t}r_{t'}$$
(f)

Subtracting the squared expectation from the expectation of the squared estimator:

$$\mathbb{V}[\hat{V}_{\pi_{e}}^{CIS}] = \sum_{n=1}^{N} \prod_{t=0}^{T} \left((\pi_{b}(a_{t}^{(n)}|s_{t}^{(n)}) - \pi_{b}(a_{t}^{(n)}|s_{t}^{(n)})^{2}) (\frac{\pi_{e}^{c}(a_{t}^{(n)}|\tilde{c}_{t}^{(n)})}{\pi_{b}^{c}(a_{t}^{(n)}|\tilde{c}_{t}^{(n)})})^{2} \right) \sum_{t=0}^{T} \sum_{t'=0}^{T} \gamma^{(t+t')} r_{t} r_{t'}$$
(g)

Explanation of steps:

- (a) We begin with the definition of variance for our estimator.
- (b) We expand the square of the estimator as the square of a sum of weighted returns.
- (c),(d) We further expand the expected value of this squared sum and evaluate the expected values under the assumption that trajectories are sampled independently.
 - (e) We calculate the square of the expectation of the estimator.
 - (f) We expand this squared expectation.
 - (g) We obtain the final expression for variance by subtracting the squared expectation from the expectation of the squared estimator, simplifying to consider the covariance terms.

F.1.3 Variance comparison between Concept IS ratios and Traditional IS ratios

Theorem.
$$\mathbb{V}[\prod_{t=0}^{T} \frac{\pi_e^c(a_t|\tilde{c}_t)}{\pi_b^c(a_t|\tilde{c}_t)}] \le \mathbb{V}[\prod_{t=0}^{T} \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)}]$$

Proof: The proof is similar to Pavse & Hanna (2022b) and the ones we used in known concepts, where we generalize to parameterized concepts from state abstractions. The proof remains intact because we make no assumptions on how the concepts are derived, as long as they satisfy the desiderata. Using Lemma E and Assumption 4.1, we can say that:

$$\mathbb{E}_{c \sim d_{\pi^c}} \prod_{t=0}^T \frac{\pi_e^c(a_t | \tilde{c}_t)}{\pi_b^c(a_t | \tilde{c}_t)} = \mathbb{E}_{s \sim d_{\pi}} \prod_{t=0}^T \frac{\pi_e(a_t | s_t)}{\pi_b(a_t | s_t)} = 1$$
(a)

Denoting the difference between the two variances as D:

$$D = Var[\prod_{t=0}^{T} \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)}] - Var[\prod_{t=0}^{T} \frac{\pi_e^c(a_t|\tilde{c}_t)}{\pi_b^c(a_t|\tilde{c}_t)}]$$
(b)

$$= \mathbb{E}_{\pi_b} \left[\prod_{t=0}^T \left(\frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} \right)^2 \right] - \mathbb{E}_{\pi_b} \left[\prod_{t=0}^T \left(\frac{\pi_e^c(a_t|\tilde{c}_t)}{\pi_b^c(a_t|\tilde{c}_t)} \right)^2 \right]$$
(c)

$$=\sum_{s}\prod_{t=0}^{T}\pi_{b}(a_{t}|s_{t})\left[\prod_{t=0}^{T}\left(\frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}\right)^{2}\right]-\sum_{c}\prod_{t=0}^{T}\pi_{b}^{c}(a_{t}|\tilde{c}_{t})\left[\prod_{t=0}^{T}\left(\frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}\right)^{2}\right]$$
(d)

$$= \sum_{c} \left(\sum_{s} \prod_{t=0}^{T} \pi_{b}(a_{t}|s_{t}) [\prod_{t=0}^{T} \left(\frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \right)^{2}] - \prod_{t=0}^{T} \pi_{b}^{c}(a_{t}|\tilde{c}_{t}) [\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})} \right)^{2}] \right) \quad (e)$$

$$= \sum_{c} \left(\sum_{s} \left[\prod_{t=0}^{T} \left(\frac{\pi_{e}(a_{t}|s_{t})^{2}}{\pi_{b}(a_{t}|s_{t})} \right) \right] - \left[\prod_{t=0}^{T} \left(\frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})^{2}}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})} \right) \right] \right)$$
(f)

We will analyse the difference of variance for 1 fixed concept and denote it as D':

$$D' = \left[\prod_{t=0}^{T} \left(\frac{\pi_e^c(a_t|\tilde{c}_t)^2}{\pi_b^c(a_t|\tilde{c}_t)}\right)\right] - \left(\sum_s \left[\prod_{t=0}^{T} \left(\frac{\pi_e(a_t|s_t)^2}{\pi_b(a_t|s_t)}\right)\right]\right)$$
(g)

Now, if we can show $D' \ge 0$ for |c|, where |c| is the cardinality of concept representation, then the difference will always be positive, thus completing our proof. We will use induction to prove $D' \ge 0$

on the total number of concepts from 1 to |c| = n < |S|. Now, our induction statement T(n) to prove is, $D' \ge 0$ where n = |c'|. For n = 1, the statement is trivially true where every concept can be represented as the traditional representation of the state. Our inductive hypothesis states that

$$D' = \left(\left[\prod_{t=0}^{T} \left(\frac{\pi_e^c(a_t | \tilde{c}_t)^2}{\pi_b^c(a_t | \tilde{c}_t)} \right) \right] - \sum_s \left[\prod_{t=0}^{T} \left(\frac{\pi_e(a_t | s_t)^2}{\pi_b(a_t | s_t)} \right) \right] \right) \ge 0$$
 (h)

Now, we define $S = \sum_{s} [\prod_{t=0}^{T} \left(\frac{\pi_e(a_t|s_t)^2}{\pi_b(a_t|s_t)} \right)]$, $C = \prod_{t=0}^{T} \pi_e^c(a_t|\tilde{c}_t)^2$, $C' = \prod_{t=0}^{T} \pi_b^c(a_t|\tilde{c}_t)$. After making the substitutions, we obtain

$$C^2 \le SC' \tag{i}$$

This result holds true for |c| = n as per the induction. Now, we add a new state s_{n+1} to the concept as part of the induction, and obtain the following difference:

$$D' = S \times \frac{\pi_e(a|s_{n+1})^2}{\pi_b(a|s_{n+1})} - \frac{C}{C'} \times \frac{\pi_e(a|s_{n+1})^2}{\pi_b(a|s_{n+1})}$$
(j)

Let $\pi_e(a|s_{n+1}) = X$ and $\pi_b(a|s_{n+1}) = Y$. Substituting, we get:

$$D' = S\frac{X^2}{Y} - \frac{C}{C'}\frac{X^2}{Y} = \frac{(SC' - C)X^2}{C'Y}$$
(k)

D' is minimum when C is maximized, hence we substitute $C \le \sqrt{SC'}$ from the induction hypothesis in the expression

$$D' \le \frac{(SC' - \sqrt{SC'})X^2}{C'Y} \tag{1}$$

As $SC' \ge 0$, the term $SC' - \sqrt{SC'}$ is never negative, leading to $D' \le 0$, since the remaining quantities are always positive. Thus, the induction hypothesis holds, and that concludes the proof.

F.1.4 Variance comparison between unknown CIS and IS estimators

Theorem. When $Cov(\prod_{t=0}^{t} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|s_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}) \leq Cov(\prod_{t=0}^{k} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}), the variance of unknown CIS estimator is lower than IS estimator, i.e. <math>\mathbb{V}_{\pi_{b}}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{IS}].$

Proof: Using Lemma E and Assumption 4.1, we can say that:

$$\mathbb{E}_{c \sim d_{\pi^c}} \left[\sum_{t=0}^T \prod_{t=0}^T \frac{\pi_e^c(a_t | \tilde{c}_t)}{\pi_b^c(a_t | \tilde{c}_t)} r_t(c_t, a_t) \right] = \mathbb{E}_{s \sim d_{\pi}} \left[\sum_{t=0}^T \prod_{t=0}^T \frac{\pi_e(a_t | s_t)}{\pi_b(a_t | s_t)} r_t(s_t, a_t) \right]$$
(a)

The Variance for a single example of a CIS estimator is given by

$$\mathbb{V}[\hat{V}^{CIS}] = \frac{1}{T^2} \left(\sum_{t=0}^T \mathbb{V}[\prod_{t=0}^T \frac{\pi_e^c(a_t | \tilde{c}_t)}{\pi_b^c(a_t | \tilde{c}_t)} r_t] + 2 \sum_{t < k} Cov(\prod_{t=0}^T \frac{\pi_e^c(a_t | \tilde{c}_t)}{\pi_b^c(a_t | \tilde{c}_t)} r_t, \prod_{t=0}^T \frac{\pi_e^c(a_t | \tilde{c}_t)}{\pi_b^c(a_t | \tilde{c}_t)} r_k) \right) \quad (b)$$

The Variance for a single example of a IS estimator is given by

$$\mathbb{V}[\hat{V}^{IS}] = \frac{1}{T^2} \left(\sum_{t=0}^T \mathbb{V}[\prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_t] + 2\sum_{t< k} Cov(\prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_t, \prod_{t=0}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} r_k) \right) \quad (c)$$

We take the difference between the variances, and note the difference of the covariances is not positive as per the assumption. Hence, if we show the differences of variances per timestep is negative, we complete our proof.

$$D = \mathbb{V}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}] - \mathbb{V}[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}]$$
(d)

$$= \mathbb{E}_{\pi_{b}^{c}} \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})} r_{t} \right)^{2} \right] - \mathbb{E}_{\pi_{b}} \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} r_{t} \right)^{2} \right]$$
(e)

$$=\sum_{c}\prod_{t=0}^{T} \left(\pi_{b}^{c}(a_{t}|\tilde{c}_{t})\right) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}\right)^{2}\right] - \sum_{s}\prod_{t=0}^{T} \left(\pi_{b}(a_{t}|s_{t})\right) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)^{2}\right]$$
(f)
$$=\sum_{c}\sum_{s\in\phi^{-1}(c)}\prod_{t=0}^{T} \left(\pi_{b}^{c}(a_{t}|\tilde{c}_{t})\right) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}\right)^{2}\right] - \sum_{s}\prod_{t=0}^{T} \left(\pi_{b}(a_{t}|s_{t})\right) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}\right)^{2}\right]$$
(g)

$$\leq R_{max}^{2} \left(\sum_{c} \sum_{s \in \phi^{-1}(c)} \prod_{t=0}^{T} \left(\pi_{b}^{c}(a_{t}|\tilde{c}_{t}) \right) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})} \right)^{2} \right] - \sum_{s} \prod_{t=0}^{T} \left(\pi_{b}(a_{t}|s_{t}) \right) \left[\left(\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \right)^{2} \right] \right)$$
(h)

The rest of the proof is identical to the previous subsection of known concepts, wherein we apply induction over the cardinality of the concepts and show the term inside the bracket is never positive, thus completing the proof.

F.1.5 Upper Bound on the Bias

Unlike known concepts, there exists a finite bias in case of unknown concepts, and the finite bounds need to be analyzed.

$$Bias = \left|\sum_{n=1}^{N} \prod_{t=0}^{T} \left(\pi_b(a_t^{(n)} | s_t^{(n)}) \frac{\pi_e^c(a_t^{(n)} | \tilde{c}_t^{(n)})}{\pi_b^c(a_t^{(n)} | \tilde{c}_t^{(n)})} \right) \sum_{t=0}^{T} \gamma^t r_t^{(n)} - \mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e^c}^{CIS}] \right|$$
(a)

$$\leq |\sum_{n=1}^{N} \prod_{t=0}^{T} \left(\pi_{e}^{c}(a_{t}^{(n)} | \tilde{c}_{t}^{(n)}) (\frac{\pi_{b}(a_{t}^{(n)} | s_{t}^{(n)})}{\pi_{b}(a_{t}^{(n)} | \tilde{c}_{t}^{(n)})}) \right) \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)} | + |\mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}^{c}}^{CIS}]|$$
(b)

$$\leq \frac{1}{N} \left| \sum_{n=1}^{N} \prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}^{(n)} | \tilde{c}_{t}^{(n)})}{\pi_{b}^{c}(a_{t}^{(n)} | \tilde{c}_{t}^{(n)})} \sum_{t=0}^{T} \gamma^{t} r_{t}^{(n)} \right| + \frac{7TR_{max} U_{c}^{T} ln(\frac{2}{\delta})}{3(N-1)} + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3} - N^{2}}} \sum_{i
(c)$$

$$\leq TR_{max}U_{c}^{T}(\frac{1}{N} + \frac{ln_{\delta}^{2}}{3(N-1)}) + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3} - N^{2}}\sum_{i < j}^{N}(X_{i} - X_{j})^{2}} + |\mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}}^{CIS}]| \tag{d}$$

Explanation of steps:

- (a) We begin with the evaluated Bias expression.
- (b) Applying triangle inequality.
- (c) Applying Bernstein inequality with probability 1- δ . X_i refers to the CIS estimate for 1 sample.
- (d) Grouping terms 1 and 2 together, where $U_c = max \frac{\pi_c^e(a|\tilde{c})}{\pi_c^e(a|\tilde{c})}$.

The first term of the bias dominates the second in terms of the number of samples, with the true expectation of the CIS estimator being unknown in general cases. Generally, the maximum possible

reward is known, which leads to the first term dominating the Bias expression. Thus, Bias is of the complexity $\mathcal{O}(\frac{TR_{max}U_c^T}{N})$

F.1.6 Upper Bound on the Variance

$$\mathbb{V}[\hat{V}_{\pi_{b}}^{CPDIS}] = \mathbb{E}_{\pi_{b}}\left(\left(\sum_{t=0}^{T}\gamma^{t}r_{t}\prod_{t'=0}^{T}\frac{\pi_{e}^{c}(a_{t'}|\tilde{c}_{t'})}{\pi_{b}^{c}(a_{t'}|\tilde{c}_{t'})}\right)^{2}\right) - \mathbb{E}_{\pi_{b}}\left(\sum_{t=0}^{T}\gamma^{t}r_{t}\prod_{t'=0}^{T}\frac{\pi_{e}^{c}(a_{t'}|\tilde{c}_{t'})}{\pi_{b}^{c}(a_{t'}|\tilde{c}_{t'})}\right)^{2}$$
(a)

$$\leq \mathbb{E}_{\pi_b}\left(\left(\sum_{t=0} \gamma^t r_t \prod_{t'=0} \frac{\pi_e^{\circ}(a_{t'}|c_{t'})}{\pi_b^{\circ}(a_{t'}|\tilde{c}_{t'})}\right)\right)$$
(b)

$$\leq \frac{1}{N} \sum_{n=1}^{N} \left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{T} \frac{\pi_{e}^{c}(a_{t'} | \tilde{c}_{t'})}{\pi_{b}^{c}(a_{t'} | \tilde{c}_{t'})} \right)^{2} \right) + \frac{7T^{2} R_{max}^{2} U_{c}^{2T} ln(\frac{2}{\delta})}{3(N-1)} + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3} - N^{2}}} \sum_{i < j}^{N} (X_{i}^{2} - X_{j}^{2})^{2}}$$
(c)

$$\leq T^2 R_{max}^2 U_c^{2T} \left(\frac{1}{N} + \frac{\ln \frac{2}{\delta}}{3(N-1)}\right) + \sqrt{\frac{\ln(\frac{2}{\delta})}{N^3 - N^2}} \sum_{i < j}^N (X_i^2 - X_j^2)^2 \tag{d}$$

Explanation of steps:

- (a) We begin with the definition of variance.
- (b) The second term is always greater than 0
- (c) Applying Bernstein inequality with probability 1- δ . X_i refers to the CIS estimate for 1 sample.
- (d) Grouping terms 1 and 2 together, where $U_c = max \frac{\pi_e(a|\tilde{c})}{\pi_b(a|\tilde{c})}$.

The first term of the variance dominates the second with increase in number of samples. Thus, Variance is of the complexity $\mathcal{O}(\frac{T^2 R_{max}^2 U_c^{2T}}{N})$

F.1.7 Upper Bound on the MSE

$$MSE = Bias^2 + Variance \tag{a}$$

$$\sim \mathcal{O}(\frac{TR_{max}U_{c}^{T}}{N})^{2} + \epsilon(|\mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}}^{CIS}]|^{2}) + \mathcal{O}(\frac{T^{2}R_{max}^{2}U_{c}^{2T}}{N})$$
(b)

$$\sim \mathcal{O}\left(\frac{T^2 R_{max}^2 U_c^{2T}}{N}\right) + \epsilon\left(|\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CIS}]|^2\right) \tag{c}$$

$$\sim \mathcal{O}(\frac{T^2 R_{max}^2 U_s^{2T}}{N}) K^{2T} + \epsilon (|\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CIS}]|^2)$$
(d)

The arguments are similar to the known-concept bounds of the MSE, with the difference being the expressions for U_c, U_s, K are over approximations of concepts instead of true concepts and an irreducible error over $\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CIS}]$ as the distribution is sampled in the concept representations instead of state representations.

F.1.8 Variance comparison with MIS estimator

Theorem. When $Cov(\prod_{t=0}^{t} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) \leq Cov(\frac{d^{\pi_{e}}(s_{t},a_{t})}{d^{\pi_{b}}(s_{t},a_{t})}r_{t}, \frac{d^{\pi_{e}}(s_{k},a_{k})}{d^{\pi_{b}}(s_{k},a_{k})}r_{k}), the variance is lower than the Variance of MIS estimator, i.e. <math>\mathbb{V}_{\pi_{b}}[\hat{V}^{CIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{MIS}].$

Proof: We start from the assumption:

$$Cov(\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}, \prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{k}) = \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})} \prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}r_{k}] - \mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}]\mathbb{E}[\prod_{t=0}^{T} \frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})}r_{t}]$$
(a)

$$= \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} K^{2} r_{t} r_{k}\right] - \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} K r_{t}\right] \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} K r_{k}\right]$$
(b)

$$\leq \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} \prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} r_{t}r_{k}\right] - \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} r_{t}\right] \mathbb{E}\left[\prod_{t=0}^{T} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})} r_{k}\right]$$
(c)

$$\leq \mathbb{E}[(\frac{d^{\pi_{e}}(s_{t}, a_{t})}{d^{\pi_{b}}(s_{t}, a_{t})})(\frac{d^{\pi_{e}}(s_{k}, a_{k})}{d^{\pi_{b}}(s_{k}, a_{k})})r_{t}r_{k}] - \mathbb{E}[\frac{d^{\pi_{e}}(s_{t}, a_{t})}{d^{\pi_{b}}(s_{t}, a_{t})})r_{t}]\mathbb{E}[\frac{d^{\pi_{e}}(s_{k}, a_{k})}{d^{\pi_{b}}(s_{k}, a_{k})})r_{k}]$$
(d)

Explanation of steps:

- (a) We begin with the definition of covariance.
- (b),(c) Using the definition of π^c , with K (the ratio of state-space distribution ratio) < 1.
 - (c) Applying Lemma E.1.7 to both the terms.

Finally, using Lemma E.1.7, substituting $Y_t = \prod_{t'=0}^T \left(\frac{\pi_e(a_{t'}|s_{t'})}{\pi_b(a_{t'}|s_{t'})}\right) r_t$ and $X_t = s_t, a_t, r_t$ completes our proof.

F.2 PDIS

F.2.1 Bias

$$Bias = |\mathbb{E}_{\pi_b}[\hat{V}_{\pi_e}^{CPDIS}] - \mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e^c}^{CPDIS}]|$$
(a)

$$= \left| \mathbb{E}_{\pi_b} \left[\sum_{t=0}^T \gamma^t \rho_{0:t}^{(n)} r_t^{(n)} \right] - \mathbb{E}_{\pi_e^c} [\hat{V}_{\pi_e^c}^{CPDIS}] \right|$$
(b)

$$= \left|\sum_{n=1}^{N} \left(\prod_{t=0}^{T} \pi_{b}(a_{t}^{(n)}|s_{t}^{(n)})\right) \sum_{t=0}^{T} \gamma^{t} \rho_{0:t}^{(n)} r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}^{c}}^{CPDIS}]\right|$$
(c)

$$= \left|\sum_{n=1}^{N}\sum_{t=0}^{T}\gamma^{t}\left(\prod_{t'=0}^{t}\pi_{e}^{c}(a_{t'}^{(n)}|\tilde{c}_{t'}^{(n)})(\frac{\pi_{b}(a_{t'}^{(n)}|s_{t'}^{(n)})}{\pi_{b}(a_{t'}^{(n)}|\tilde{c}_{t'}^{(n)})})\right)r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}^{c}}^{CPDIS}]\right|$$
(d)

Explanation of steps: Similar to CIS.

F.2.2 Variance

Following the process similar to CIS estimator:

$$\mathbb{V}[\hat{V}_{\pi_b}^{CPDIS}] = \mathbb{E}_{\pi_b}[(\hat{V}_{\pi_b}^{CPDIS})^2] - (\mathbb{E}_{\pi_b}[\hat{V}_{\pi_b}^{CPDIS}])^2 \tag{a}$$

We first evaluate the expectation of the square of the estimator:

$$\mathbb{E}_{\pi_b}[(\hat{V}_{\pi_b}^{CPDIS})^2] = \mathbb{E}_{\pi_b} \left[\left(\sum_{t=0}^T \gamma^t \rho_{0:t} r_t \right)^2 \right]$$
(b)

$$= \mathbb{E}_{\pi_b} \left[\sum_{t=0}^T \sum_{t'=0}^T \rho_{0:t} \rho_{0:t'} \gamma^{(t+t')} r_t r_{t'} \right]$$
(c)

$$=\sum_{n=1}^{N}\sum_{t=0}^{T}\sum_{t'=0}^{T}\left(\prod_{t''=0}^{t}\pi_{b}(a_{t}|s_{t})(\frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})})\right)\left(\prod_{t'''=0}^{t'}(\frac{\pi_{e}^{c}(a_{t}|\tilde{c}_{t})}{\pi_{b}^{c}(a_{t}|\tilde{c}_{t})})\right)\gamma^{(t+t')}r_{t}r_{t'}$$
(d)

Evaluating the second term in the variance expression:

$$(\mathbb{E}_{\pi_{b}}[\hat{V}_{\pi_{b}}^{CPDIS}])^{2} = \left(\sum_{n=1}^{N}\sum_{t=0}^{T}\mathbb{E}_{\pi_{b}}[\gamma^{t}\rho_{0:t}r_{t}]\right)^{2}$$
(e)
$$= \sum_{n=1}^{N}\sum_{t=0}^{T}\sum_{t'=0}^{T}\left(\prod_{t''=0}^{t}\pi_{b}(a_{t''}|s_{t''})(\frac{\pi_{e}^{c}(a_{t''}|\tilde{c}_{t''})}{\pi_{b}^{c}(a_{t''}|\tilde{c}_{t''})})\right) \left(\prod_{t'''=0}^{t'}\pi_{b}(a_{t'''}|s_{t'''})(\frac{\pi_{e}^{c}(a_{t'''}|\tilde{c}_{t'''})}{\pi_{b}^{c}(a_{t'''}|\tilde{c}_{t'''})})\right) \gamma^{(t+t')}r_{t}r_{t}$$
(f)

Subtracting the squared expectation from the expectation of the squared estimator:

$$\mathbb{V}[\hat{V}_{\pi_{b}}^{CPDIS}] = \sum_{n=1}^{N} \sum_{t=0}^{T} \sum_{t'=0}^{T} \left(\prod_{t'''=0}^{t} \pi_{b}(a_{t'''}|s_{t'''}) (\frac{\pi_{e}^{c}(a_{t'''}|\tilde{c}_{t'''})}{\pi_{b}^{c}(a_{t'''}|\tilde{c}_{t'''})}) \right) \left(\prod_{t''=0}^{t'} (1 - \pi_{b}(a_{t''}|s_{t''})) (\frac{\pi_{e}^{c}(a_{t''}|\tilde{c}_{t''})}{\pi_{b}^{c}(a_{t''}|\tilde{c}_{t'''})}) \right) \gamma^{(t+t')} r_{t} r_{t'}$$

Explanation of steps: Similar to CIS

F.2.3 Variance comparison between unknown CPDIS and PDIS estimators

 $\begin{array}{l|lll} \textbf{Theorem} & \textbf{F.1.} & \textbf{When} & Cov(\prod_{t=0}^{t} \frac{\pi_{c}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{c}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) & \leq \\ Cov(\prod_{t=0}^{k} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}(a_{t}|s_{t})}{\pi_{b}(a_{t}|s_{t})}r_{k}), & \text{the variance of parameterized CPDIS estimators} \\ is lower than PDIS estimator, i.e. & \mathbb{V}_{\pi_{b}}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{PDIS}]. \end{array}$

Proof: Similar to CIS estimator.

F.2.4 Upper Bound on the Bias

Unlike known concepts, there exists a finite bias in case of unknown concepts, and the bounds need to be analyzed.

$$Bias = \left|\sum_{n=1}^{N}\sum_{t=0}^{T}\gamma^{t}\left(\prod_{t'=0}^{t}\pi_{e}^{c}(a_{t'}^{(n)}|\tilde{c}_{t'}^{(n)})(\frac{\pi_{b}(a_{t'}^{(n)}|s_{t'}^{(n)})}{\pi_{b}^{c}(a_{t'}^{(n)}|\tilde{c}_{t'}^{(n)})})\right)r_{t}^{(n)} - \mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}}^{CPDIS}]\right|$$
(a)

$$\leq |\sum_{n=1}^{N} \sum_{t=0}^{T} \gamma^{t} \left(\prod_{t'=0}^{t} \pi_{e}^{c}(a_{t'}^{(n)} | \tilde{c}_{t'}^{(n)}) (\frac{\pi_{b}(a_{t'}^{(n)} | s_{t'}^{(n)})}{\pi_{b}^{c}(a_{t'}^{(n)} | \tilde{c}_{t'}^{(n)})}) \right) r_{t}^{(n)} + |\mathbb{E}_{\pi_{e}^{c}} [\hat{V}_{\pi_{e}}^{CPDIS}]|$$
(b)

$$\leq \frac{1}{N} \sum_{n=1}^{N} \sum_{t=0}^{T} \gamma^{t} \prod_{t'=0}^{t} (\frac{\pi_{e}^{c}(a_{t'}^{(n)} | \tilde{c}_{t'}^{(n)})}{\pi_{b}^{c}(a_{t'}^{(n)} | \tilde{c}_{t'}^{(n)})}) r_{t}^{(n)} + \frac{7TR_{max}U_{c}^{T}ln(\frac{2}{\delta})}{3(N-1)} + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3}-N^{2}}} \sum_{i$$

$$\leq TR_{max}U_{c}^{T}(\frac{1}{N} + \frac{ln_{\overline{\delta}}^{2}}{3(N-1)}) + \sqrt{\frac{ln(\frac{2}{\overline{\delta}})}{N^{3} - N^{2}}\sum_{i < j}^{N}(X_{i} - X_{j})^{2}} + |\mathbb{E}_{\pi_{e}^{c}}[\hat{V}_{\pi_{e}}^{CPDIS}]|$$
(d)

Explanation of steps: Similar to CIS.

F.2.5 Upper Bound on the Variance

$$\mathbb{V}[\hat{V}_{\pi_{b}}^{CPDIS}] = \mathbb{E}_{\pi_{b}}\left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{t} \frac{\pi_{e}^{c}(a_{t'}|\tilde{c}_{t'})}{\pi_{b}^{c}(a_{t'}|\tilde{c}_{t'})}\right)^{2}\right) - \mathbb{E}_{\pi_{b}}\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{t} \frac{\pi_{e}^{c}(a_{t'}|\tilde{c}_{t'})}{\pi_{b}^{c}(a_{t'}|\tilde{c}_{t'})}\right)^{2}$$
(a)

$$\leq \mathbb{E}_{\pi_b}\left(\left(\sum_{t=0}^{r} \gamma^t r_t \prod_{t'=0}^{r} \frac{\pi_e^c(a_{t'}|\tilde{c}_{t'})}{\pi_b^c(a_{t'}|\tilde{c}_{t'})}\right) \right)$$
(b)

$$\leq \frac{1}{N} \sum_{n=1}^{N} \left(\left(\sum_{t=0}^{T} \gamma^{t} r_{t} \prod_{t'=0}^{t} \frac{\pi_{e}^{c}(a_{t'}|\tilde{c}_{t'})}{\pi_{b}^{c}(a_{t'}|\tilde{c}_{t'})} \right)^{2} \right) + \frac{7T^{2}R_{max}^{2}U_{c}^{2T}ln(\frac{2}{\delta})}{3(N-1)} + \sqrt{\frac{ln(\frac{2}{\delta})}{N^{3}-N^{2}}} \sum_{i

$$(c)$$$$

$$\leq T^2 R_{max}^2 U_c^{2T} (\frac{1}{N} + \frac{\ln \frac{2}{\delta}}{3(N-1)}) + \sqrt{\frac{\ln(\frac{2}{\delta})}{N^3 - N^2}} \sum_{i < j}^N (X_i^2 - X_j^2)^2 \tag{d}$$

Explanation of steps: Similar to CIS.

F.2.6 Upper Bound on the MSE

$$MSE = Bias^2 + Variance \tag{a}$$

$$\sim \mathcal{O}(\frac{TR_{max}U_c^T}{N})^2 + \epsilon(|\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CPDIS}]|^2) + \mathcal{O}(\frac{T^2R_{max}^2U_c^{2T}}{N})$$
(b)

$$\sim \mathcal{O}(\frac{T^2 R_{max}^2 U_c^{2T}}{N}) + \epsilon(|\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CPDIS}]|^2) \tag{c}$$

$$\sim \mathcal{O}(\frac{T^2 R_{max}^2 U_s^{2T}}{N}) K^{2T} + \epsilon (|\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CPDIS}]|^2) \tag{d}$$

The arguments are similar to the known-concept bounds of the MSE, with the difference being the expressions for U_c, U_s, K are over approximations of concepts instead of true concepts and an irreducible error over $\mathbb{E}_{\pi_e^c}[\hat{V}_{\pi_e}^{CPDIS}]$ as the distribution is sampled in the concept representations instead of state representations.

F.2.7 Variance comparison with MIS estimator

Theorem. When $Cov(\prod_{t=0}^{t} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{t}, \prod_{t=0}^{k} \frac{\pi_{e}^{c}(a_{t}|c_{t})}{\pi_{b}^{c}(a_{t}|c_{t})}r_{k}) \leq Cov(\frac{d^{\pi_{e}}(s_{t},a_{t})}{d^{\pi_{b}}(s_{t},a_{t})}r_{t}, \frac{d^{\pi_{e}}(s_{k},a_{k})}{d^{\pi_{b}}(s_{k},a_{k})}r_{k}),$ the variance is lower than the Variance of MIS estimator, i.e. $\mathbb{V}_{\pi_{b}}[\hat{V}^{CPDIS}] \leq \mathbb{V}_{\pi_{b}}[\hat{V}^{MIS}].$

Explanation of Steps: Similar to CIS

G Environments

	20 -					
	16	Wind: ← (-1) Penalty: -1	Wind: ← (-1) Penalty: -1	Wind: ←↑ (-1,+1) Penalty: -2	Penalty: 0	Penalty: 0
Y	10 -	Wind: ← (-1) Penalty: -1	Wind: ←↑ (-1,+1) Penalty: -2	Penalty: 0	Penalty: 0	Wind: →↓ (+1,-1) Penalty: -2
	Ω.	Wind: ←↑ (-1,+1) Penalty: -2	Penalty: 0	Penalty: 0	Wind: →↓ (+1,-1) Penalty: -2	Wind: ↓ (+0,-1) Penalty: -1
	4	Penalty: 0	Penalty: 0	Wind: →↓ (+1,-1) Penalty: -2	Wind: ↓ (+0,-1) Penalty: -1	Wind: ↓ (+0,-1) Penalty: -1
	4	Penalty: 0	Wind: →↓ (+1,-1) Penalty: -2	Wind: ↓ (+0,-1) Penalty: -1	Wind: ↓ (+0,-1) Penalty: -1	Wind: ↓ (+0,-1) Penalty: -1
	0 -	0 4	4 8	3 1 X	2 1	6 20

Figure 6: Schematic of windy-gridworld environment. The top-right corner refers to the goal target of the agent. The wind direction and reward penalty is indicated in each region.

WindyGridworld Fig 6 illustrates the Windy Gridworld environment, a 20x20 grid divided into regions with varying wind directions and penalties. The agent's goal is to navigate from a randomly chosen starting point to a fixed goal in the top-right corner. Off-diagonal winds increase in strength near non-windy regions, affecting the agent's movement. Each of the four available actions moves the agent four steps in the chosen direction. Reaching the goal earns a +5 reward while moving away results in a -0.2 penalty. Additional negative rewards are based on regional penalties within the grid. Each episode ends after 200 steps.

The grid is split into 25 blocks, each measuring 4x4 units with each region having a penalty based on the wind-strength. Blocks affected by wind display the direction and strength (e.g., ' \leftarrow ↑ (-2,+2)' indicates northward and westward winds with a strength of 2 units each). This setup encourages the agent to navigate through non-penalty areas for optimal rewards.

MIMIC-III We use the publicly available MIMIC-III database (Johnson et al., 2016) from PhysioNet (Goldberger et al., 2000), which records the treatment and progression of ICU patients at the Beth Israel Deaconess Medical Center in Boston, Massachusetts. We focus on the task of managing acutely hypotensive patients in the ICU. Our preprocessing follows the original MIMIC-III steps detailed in Komorowski et al. (2018c) and used in subsequent works (Keramati et al., 2021b; Matsson & Johansson, 2021). After processing the data in Excel, we group patients by 'ICU-stayID' to form distinct trajectories.

The state space includes 15 features: Creatinine, FiO_2 , Lactate, Partial Pressure of Oxygen, Partial Pressure of CO₂, Urine Output, GCS score, and electrolytes such as Calcium, Chloride, Glucose, HCO₃, Magnesium, Potassium, Sodium, and SpO₂. Each feature is binned into 10 levels from 0 (very low) to 9 (very high).

Treatments for hypotension include IV fluid bolus administration and vasopressor initiation, with doses categorized into four levels: "none," "low," "medium," and "high," forming a total action space of 16 discrete actions. The reward function depends on the next mean arterial pressure (MAP) and ranges from -1 to 0, linearly distributed between 20 and 65. A MAP above 65 indicates that the patient is not experiencing hypotension.

H Additional Experimental Details

H.1 Known Concepts Policy Extraction

PPO exoerimental details: We run a PPO algorithm over our windygridworld environment to derive our evaluation and behavior policies. The hyperparameters are as follows: 10k training episodes, batch size of 64, discount factor of 0.99, learning rate of 1e-4, rate of exploration of 0.2, 5 epochs between actor and critic update, update timestep interval of 2000. We take the evaluation policy as the policy learnt at epoch 10k while the behavior policy as the policy learnt at epoch 5k.

H.2 Unknown Concepts Experimental setup

Environments, Policy descriptions, Metrics: Same as those in known concepts section.

Training and Hyperparameter Details: We use 400 training, 50 validation, and 50 test trajectories sampled from the behavior policy to train the CBMs, which predict the next state transitions from the current state. The model architecture includes an input layer, a bottleneck, two 256-neuron layers, and an output layer, all with ReLU activations. Training is performed using the Adam Optimizer with a learning rate of 1e-3 on an Nvidia P100 GPU (16 GB) within the Pytorch framework.

Training targets multiple loss components: the OPE metric, interpretability, diversity, and CBM output. The non-convex nature of the loss landscape can lead to issues such as non-convergence and NaN values. To address this, we employ a three-stage training strategy.

In the first stage, we optimize all losses except the OPE metric to stabilize the initial training process. In the second stage, the OPE metric is gradually incorporated into the optimization until convergence is achieved. Finally, in the third stage, we freeze the CBM weights to refine the remaining losses while controlling variations in the OPE metric.

This strategy balances the learning of critical on-policy features with maintaining relevant OPE metrics, thereby enhancing concept learning and policy generalization. Despite these efforts, managing the complexity of the loss landscape remains a significant challenge, particularly in dynamic environments, and represents an important direction for future research.

H.3 Known, Oracle and Intervened Concepts

Table 1 refers to the list of concepts used for Windygridworld experiments. The table consists of Oracle concepts, good known concepts, poor known concepts, and learned concepts. In case of

interventions, we intervene on misclassified learned concepts using good known concepts. In case of MIMIC, oracle concepts are unknown as it's a real world practical example. In case of good known concepts, the concepts are obtained by applying Approximate nearest neighbors algorithm over 200 nearest neighbors for all features discretized over 10 levels. For poor MIMIC concepts, we discretize the features into just 2 levels, thus leading to ambiguous features and thereby lack of diversity.

v	V	Known (Good)	Known (Poor)	Oracle	Learned
	I	Concept	Concept	Concept	Concept
(0,4)	(0,4)	0	0	0	0
(4,8)	(0,4)	1	1	1	1
(8,12)	(0,4)	2	2	1	1
(12,16)	(0,4)	3	3	1	0
(16,20)	(0,4)	4	4	1	1
(0,4)	(4,8)	5	0	2	2
(4,8)	(4,8)	6	1	0	0
(8,12)	(4,8)	7	2	1	1
(12,16)	(4,8)	8	3	1	1
(16,20)	(4,8)	9	4	1	0
(0,4)	(8,12)	10	0	2	2
(4,8)	(8,12)	11	1	2	2
(8,12)	(8,12)	12	2	0	0
(12,16)	(8,12)	13	3	1	1
(16,20)	(8,12)	14	4	1	1
(0,4)	(12,16)	15	0	2	2
(4,8)	(12,16)	16	1	2	2
(8,12)	(12,16)	17	2	2	2
(12,16)	(12,16)	18	3	3	3
(16,20)	(12,16)	19	4	3	3
(0,4)	(16,20)	20	0	2	2
(4,8)	(16,20)	21	1	2	2
(8,12)	(16,20)	22	2	2	2
(12,16)	(16,20)	23	3	2	2
(16,20)	(16,20)	24	4	3	3

 Table 1: WindyGridworld Concept Information

H.4 Additional description on polices for MIMIC-III

For the MIMIC-III dataset, it is common to generate behavior trajectories using K-nearest neighbors (KNN) as the true on-policy trajectories are unavailable. Examples of works that generate behavior trajectories or policies using KNNs include (Gottesman et al., 2020; Böck et al., 2022; Liu et al., 2022; Keramati et al., 2021b; Komorowski et al., 2018b; Peine et al., 2021). In this paper, we employ a popular variant of KNN, known as approximate nearest neighbors (ANN) search.

The advantages of ANN over traditional KNN include scalability, reduced computational cost, efficient indexing, and support for dynamic data. These benefits allow us to generate behavior and evaluation policies with a larger number of neighbors (200 in this study, which is double that used in prior works employing KNN) while achieving faster inference times. Examples of papers that use approximate nearest neighbors in medical applications include (Anagnostou et al., 2020; Gupta et al., 2022). For readers interested in the foundational work outlining the benefits of ANN over KNN, we refer to the seminal paper (Indyk & Motwani, 1998).

I Optimized Parameterized Concepts

Feature		С	IS		CPDIS			
	c_1	c_2	c_3	c_4	c_1	c_2	c_3	c_4
f_1 : X-coordinate	0.15	-0.07	0.05	0.19	-0.23	0.33	-0.03	0.03
f_2 : Y-coordinate	-0.02	-0.23	0.07	-0.12	-0.22	0.25	0.02	-0.06
f_3 : Horizontal distance from target	-0.02	0.07	-0.10	0.00	-0.15	-0.30	0.02	-0.11
f_4 : Vertical distance from target	0.06	-0.26	-0.09	0.06	-0.11	0.10	-0.04	-0.21
f_5 : Horizontal Wind	0.05	0.12	-0.12	0.00	-0.15	0.20	0.29	-0.14
f_6 : Vertical Wind	0.26	0.01	-0.02	0.00	-0.18	0.06	-0.17	0.19
f_7 : Region penalty	0.24	0.18	-0.25	0.15	0.23	0.01	-0.11	0.22
f_8 : Distance to left wall	-0.14	-0.25	0.01	0.05	-0.13	0.24	0.16	0.14
f_9 : Distance to right wall	0.02	0.00	0.01	0.19	-0.12	-0.28	0.06	0.16
f_{10} : Distance to top wall	-0.01	-0.20	-0.21	0.07	-0.33	-0.05	-0.04	-0.01
f_{11} : Distance to bottom wall	-0.16	0.07	0.22	-0.22	0.06	-0.13	0.13	-0.22
f_{12} : Penalty of left subregion	-0.06	0.08	-0.08	-0.22	-0.07	-0.01	0.03	-0.16
f_{13} : Penalty of right subregion	-0.03	0.02	-0.20	-0.20	-0.07	-0.18	-0.34	-0.21
f_{14} : Penalty of top subregion	0.16	0.19	-0.08	-0.17	0.00	0.04	-0.07	0.21
f_{15} : Penalty of bottom subregion	0.08	0.24	0.05	-0.19	0.17	-0.07	-0.12	0.21
f_{16} : Distance to left subregion	-0.11	0.05	0.00	0.26	0.10	-0.07	0.22	0.04
f_{17} : Distance to right subregion	0.00	-0.17	0.04	0.13	0.05	-0.13	0.06	0.11
f_{18} : Distance to top subregion	0.07	-0.03	0.13	0.08	-0.12	0.01	0.06	0.00
f_{19} : Distance to bottom subregion	-0.06	-0.09	-0.06	-0.01	-0.19	-0.01	0.06	0.13
Constant	-0.06	-0.16	0.14	-0.01	-0.08	-0.01	0.00	0.13

Table 2: WindyGridworld: Coefficients of the human interpretable features learnt while optimizing parameterized concepts. Here, the concept c_t is a 4-dimensional vector $[c_1, c_2, c_3, c_4]$, where $c_i = w_i^T f_i$, with f_i being the human interpretable features.

 Table 3: MIMIC: Coefficients of the human interpretable features learnt while optimizing parameterized concepts.

Feature		C	IS		CPDIS				
	c_1	c_2	c_3	c_4	c_1	c_2	c_3	c_4	
f_1 : Creatinine	-0.08	-0.24	0.19	-0.18	-0.08	-0.24	0.19	-0.18	
f_2 : FiO ₂	-0.13	0.00	0.04	-0.06	-0.13	0.00	0.04	-0.06	
f_3 : Lactate	-0.24	-0.02	-0.23	0.21	-0.24	-0.02	-0.23	0.21	
f_4 : Partial Pressure of O_2	0.09	-0.07	-0.06	-0.12	0.09	-0.07	-0.06	-0.12	
f_5 : Partial Pressure of CO_2	-0.21	0.16	0.19	-0.03	-0.21	0.16	0.19	-0.03	
f_6 : Urine Output	0.06	0.07	0.06	0.22	0.06	0.07	0.06	0.22	
f_7 : GCS Score	0.11	-0.05	-0.01	0.15	0.11	-0.05	-0.01	0.15	
f_8 : Calcium	0.16	-0.20	0.06	0.16	0.16	-0.20	0.06	0.16	
f_9 : Chloride	0.02	-0.11	-0.04	0.14	0.02	-0.11	-0.04	0.14	
f_{10} : Glucose	0.06	-0.10	-0.10	-0.08	0.05	-0.10	-0.10	-0.08	
f_{11} : HCO ₂	0.21	0.14	-0.20	-0.22	0.20	0.14	-0.20	-0.22	
f_{12} : Magnesium	-0.15	-0.02	-0.20	0.01	-0.15	-0.02	-0.20	0.01	
f_{13} : Potassium	0.04	0.08	0.15	-0.26	0.04	0.08	0.15	-0.26	
f_{14} : Sodium	0.00	-0.02	0.24	0.19	0.00	-0.02	0.24	0.19	
f_{15} : SpO ₂	-0.17	-0.20	-0.06	-0.23	-0.17	-0.20	-0.06	-0.23	