OPTIMAL ROBUST SUBSIDY POLICIES FOR IRRA-TIONAL AGENT IN PRINCIPAL-AGENT MDPS

Anonymous authors

Paper under double-blind review

ABSTRACT

We investigate a principal-agent problem modeled within a Markov Decision Process, where the principal and the agent have their own rewards. The principal can provide subsidies to influence the agent's action choices, and the agent's resulting action policy determines the rewards accrued to the principal. Our focus is on designing a robust subsidy scheme that maximizes the principal's cumulative expected return, even when the agent displays bounded rationality and may deviate from the optimal action policy after receiving subsidies.

As a baseline, we first analyze the case of a perfectly rational agent and show that the principal's optimal subsidy coincides with the policy that maximizes social welfare, the sum of the utilities of both the principal and the agent. We then introduce a bounded-rationality model: the globally ϵ -incentive-compatible agent, who accepts any policy whose expected cumulative utility lies within ϵ of the personal optimum. In this setting, we prove that the optimal robust subsidy scheme problem simplifies to a one-dimensional concave optimization. This reduction not only yields a clean analytical solution but also highlights a key structural insight: optimal subsidies are concentrated along the social-welfare-maximizing trajectories. We further characterize the loss in social welfare—the degradation under the robust subsidy scheme compared to the maximum achievable—and provide an upper bound on this loss. Finally, we investigate a finer-grained, state-wise ϵ -incentive-compatible model. In this setting, we show that under two natural definitions of state-wise incentive-compatibility, the problem becomes intractable: one definition results in a non-Markovian agent action policy, while the other renders the search for an optimal subsidy scheme NP-hard.

1 Introduction

The principal–agent problem (often modeled as a Stackelberg game) has long been central to the study of strategic interactions where one party acts on behalf of another, yet with potentially misaligned incentives. This setting arises frequently in economics and governance: for example, governments design taxes, subsidies, and public investments to guide individual behavior toward socially beneficial outcomes. However, in decentralized markets, each participant ultimately pursues their own utility, and centralized guidance can only partially influence outcomes. A similar dynamic appears in machine learning, where reinforcement learning with human feedback (RLHF) is employed to align large language models (LLMs) with societal values such as ethics and legal compliance. In both cases, the principal faces the fundamental challenge of shaping an agent's behavior without direct control, while respecting both parties' interests.

In this paper, we investigate the principal–agent problem within the framework of a Markov Decision Process (MDP), where the principal can provide subsidies to influence the agent's action choices. More specifically, in our setting, each action under each state yields two distinct rewards: one for the principal and one for the agent. The principal may also assign non-negative subsidies to actions. The agent selects an action policy based on its own reward combined with subsidies offered by the principal. The principal, in turn, strategically designs these subsidies to influence the agent's choices, aiming to maximize the principal's overall payoff, which equals the total principal's reward associated with the agent's chosen action minus the subsidies provided.

A natural assumption in such models is that the agent always behaves rationally, selecting the trajectory that maximizes the sum of the agent's own reward and the subsidies provided by the principal. Yet in practice, this assumption is often violated: agents may deviate from perfect rationality due to bounded cognition, incomplete information, or limited computational power. For example, in economics, individuals may fail to optimize utility precisely because of uncertainty or behavioral biases. Similarly, in reinforcement learning, approximate training algorithms may yield suboptimal policies due to limited exploration or finite computation.

Motivated by these considerations, we ask:

How should the principal design subsidies when the agent may behave irrationally?

Our goal is to identify a **robust subsidy scheme** that guarantees the principal the best possible expected cumulative return in the worst-case scenario.

Our Contributions We introduce a theoretical framework based on Markov Decision Processes (MDPs) to model the principal-agent problem and formulate the design of an optimal robust subsidy scheme as a minimax optimization problem. Within this framework, we systematically analyze three agent models: the perfectly rational agent, the globally ϵ -incentive-compatible (IC) agent, and the state-wise ϵ -IC agent. For each model, we provide structural insights and algorithmic solutions.

We first study a *perfectly rational agent* as a baseline, who always selects actions that maximize its own reward. In Theorem 3.1, we characterize the optimal subsidy scheme and show in Proposition 3.2 that it suffices to subsidize only actions that maximize social welfare, defined as the sum of the principal's and agent's utilities. Under this scheme, the agent's best-response policy aligns with the social welfare-maximizing policy, establishing a clear benchmark for incentive alignment.

Next, we consider *globally* ϵ -*IC agents*, who tolerate at most an ϵ loss relative to their optimal reward under a given subsidy scheme. Unlike perfectly rational agents, these agents may adopt stochastic policies, making the principal's optimization a nontrivial bi-level problem. Theorem 4.1 shows that this problem can be equivalently reduced to maximizing a one-dimensional concave function over a bounded interval, allowing efficient solution via standard first-order methods. Structurally, in Proposition 4.2, we show the optimal subsidy mirrors the perfectly rational case by exclusively rewarding actions that align with maximizing social welfare; and, in the worst-case response, the agent's policy will assign positive probability to the socially optimal actions, though it may also mix with other actions. We further provide a quantitative analysis of the gap between the total payoff achieved under this robust scheme and the maximum possible social welfare, as shown in Proposition 4.3.

Finally, in Section 5, we examine *state-wise* ϵ -*IC agents*, for which the ϵ -tolerance must hold at each individual state. Two natural formalizations arise, each presenting distinct challenges. In the first formalization, the agent's worst-case response may necessitate a non-Markovian policy, thereby violating the foundational assumptions of the MDP framework and introducing history dependence that makes the problem computationally intractable. In the second formalization, while the agent's worst-case response remains polynomial-time computable, Theorem 5.1 demonstrates that the principal's problem becomes NP-hard. These findings illustrate that, although state-wise constraints are conceptually appealing, they introduce significant computational and modeling complexities that limit practical applicability.

Related work The principal–agent problem, a central concept in economics (Ross, 1973; Grossman & Hart, 1992), arises when a principal delegates tasks to an agent whose actions may be guided by self-interest. This framework underpins both contract theory (Laffont & Maskin, 1981; Guruganesh et al., 2021) and mechanism design (Myerson, 1982; Kadan et al., 2017).

Recent work has examined this problem in the setting of Markov Decision Processes (MDPs). Research in this area falls into two broad directions. The first, information design, seeks to influence the agent's beliefs, as in Bayesian persuasion (Gan et al., 2022; Wu et al., 2022; Bernasconi et al., 2023). The second, more closely aligned with our work, focuses on shaping the agent's incentives through policy teaching (Zhang & Parkes, 2008; Banihashem et al., 2022) or environment/model design (Thoma et al., 2024; Yu & Ho, 2022). A comprehensive survey is provided by Dütting et al. (2024). Among these, two approaches are most closely related to our study:

Contract-based models. This line of research integrates contract theory with MDPs, assuming the principal observes only states and offers state-dependent payments. Prior studies analyze subgame perfect equilibrium (Wu et al., 2024; Ivanov et al., 2024), showing that history-dependent contracts are necessary for farsighted agents (Bollini et al., 2024). These works typically assume perfectly rational agents and establish that the optimal contract design problem is NP-hard.

Reward shaping. In Reward shaping, the principal modifies the agent's incentives via additional rewards for specific state—action pairs, subject to a fixed budget (Ben-Porat et al., 2024), with the design problem remaining NP-hard. Extensions address behavioral uncertainty through robust reward design (Wu et al., 2025). In contrast, we incorporate incentive costs directly into the principal's objective, treating them as part of payoff optimization rather than an external constraint.

2 Problem Formulation

The Principal-Agent MDP Model We consider a principal-agent problem modeled as a time-inhomogeneous, finite-horizon Markov Decision Process (MDP). In this setting, the principal aims to achieve a goal by influencing an agent's actions. The principal can offer subsidies to incentivize the agent to follow a policy that benefits the principal.

Formally, we define the problem instance using the tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{H}, \mathcal{P}, r_P, r_A, \hat{s}, \Pi \rangle$, where:

- S is the set of the finite states and A is the set of actions. We assume that both states and actions are *discrete*.
- $\mathcal{H} = \{0, 1, \dots, H-1\}$ is the set of time steps, with H representing the time horizon.
- $P: \mathcal{S} \times \mathcal{A} \times \mathcal{H} \to \Delta(\mathcal{S})$ is the transition kernel , where P(s'|s,a,h) indicates the probability of transferring to state $s' \in \mathcal{S}$ after executing action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$ at timestep $h \in \mathcal{H}$.
- $r_P, r_A : \mathcal{S} \times \mathcal{A} \times \mathcal{H} \to \mathbb{R}$ are the reward functions of the principal and the agent, respectively, where $r_P(s, a, h)$ (resp. $r_A(s, a, h)$) denotes the reward obtained by the principal (resp. agent) when the agent executes action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$ at timestep $h \in \mathcal{H}$.
- Without loss of generality, \hat{s} is the fixed starting state for the agent.

Subsidy Scheme and Action Policy The principal commits to a subsidy scheme $\Delta r : \mathcal{S} \times \mathcal{A} \times \mathcal{H} \to \mathbb{R}_{\geq 0}$. Here, $\Delta r(s, a, h)$ is a non-negative payment from the principal to the agent for taking action a in state s at timestep h. We denote the set of all feasible subsidy policies as \mathcal{R}_{Δ} .

Given a subsidy Δr on action a in state s at timestep h, the effective rewards for the principal and agent become:

$$r_P^{\Delta r}(s,a,h) = r_P(s,a,h) - \Delta r(s,a,h) \quad \text{and} \quad r_A^{\Delta r}(s,a,h) = r_A(s,a,h) + \Delta r(s,a,h)$$

The agent observes the subsidy scheme and then chooses a Markovian **action policy** $\pi: \mathcal{S} \times \mathcal{H} \to \Delta(\mathcal{A})$. Based on the agent's (ir)rationality, for any given Δr , the agent will choose a policy from a specific set of feasible policies, which we denote by $\Pi(\Delta r)$.

Value Functions For any player $i \in \{P, A\}$, subsidy scheme Δr , and agent policy π , we define the standard state-value and action-value functions via the Bellman expectation equations:

$$\begin{split} V_i^{\pi,\Delta r}(s,h) &= \sum_{a \in \mathcal{A}} \pi(a|s,h) Q_i^{\pi,\Delta r}(s,a,h) \\ Q_i^{\pi,\Delta r}(s,a,h) &= r_i^{\Delta r}(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h) V_i^{\pi,\Delta r}(s',h+1) \end{split}$$

with the terminal condition $V_i^{\pi,\Delta r}(s,H)=0$. Furthermore, we use $\overline{V}_A^{\Delta r}(s,h)$ and $\overline{Q}_A^{\Delta r}(s,a,h)$ to denote the optimal state-value and action-value functions attainable by the agent,

$$\overline{V}_A^{\Delta r}(s,h) = \max_{a} \overline{Q}_A^{\Delta r}(s,a,h)$$
$$\overline{Q}_A^{\Delta r}(s,a,h) = r_A(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h) \overline{V}_A^{\Delta r}(s',h+1)$$

Additionally, $V_i^{\pi,\Delta r=0}(s,h)$, $Q_i^{\pi,\Delta r=0}(s,a,h)$, $\overline{V}_A^{\Delta r=0}(s,h)$ and $\overline{Q}_A^{\Delta r=0}(s,a,h)$ denote the corresponding value in the absence of subsidies.

Social Welfare We define social welfare as the aggregate reward of both the principal and the agent: $r_{sw}(s, a, h) \triangleq r_P(s, a, h) + r_A(s, a, h)$, which remains unaffected by the subsidy term Δr .

The social welfare value functions, $V_{\rm sw}^{\pi}$ and $Q_{\rm sw}^{\pi}$, characterize the expected social welfare under an agent policy π :

$$\begin{split} V_{\text{sw}}^{\pi}(s,h) &= \sum_{a \in \mathcal{A}} \pi(a|s,h) \, Q_{\text{sw}}^{\pi}(s,a,h), \\ Q_{\text{sw}}^{\pi}(s,a,h) &= r_{\text{sw}}(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h) \, V_{\text{sw}}^{\pi}(s',h+1). \end{split}$$

Analogously, the optimal social welfare value functions, $V_{\rm sw}^*$ and $Q_{\rm sw}^*$, are defined as:

$$\begin{split} V_{\text{sw}}^*(s,h) &= \max_{a \in \mathcal{A}} Q_{\text{sw}}^*(s,a,h), \\ Q_{\text{sw}}^*(s,a,h) &= r_{\text{sw}}(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h) \, V_{\text{sw}}^*(s',h+1). \end{split}$$

An action a is said to be **social-welfare-maximizing** in state s at timestep h if it is greedy with respect to the optimal Q-value, i.e., $a \in \arg\max_{a' \in \mathcal{A}} Q_{sw}^*(s, a', h)$.

Optimization Objective We consider a robust formulation where the principal seeks a subsidy scheme that performs best against the agent's worst-case response. The agent's adversarial action policy to a subsidy Δr is an agent policy $\pi_{\Delta r}$ that minimizes the principal's expected return within the feasible set $\Pi(\Delta r)$:

$$\pi_{\Delta r} \in \operatorname*{arg\,min}_{\pi \in \Pi(\Delta r)} V_P^{\pi,\Delta r}(\hat{s}, h = 0)$$

The principal's objective is to find the optimal subsidy scheme Δr^* that maximizes this worst-case outcome. The optimal value for the principal is therefore:

$$OPT \triangleq \max_{\Delta r \in \mathcal{R}_{\Delta}} \min_{\pi \in \Pi(\Delta r)} V_P^{\pi, \Delta r}(\hat{s}, h = 0)$$
(2.1)

3 WARM-UP: THE PERFECTLY RATIONAL AGENT

We begin with the simplest setting of a perfectly rational agent, defined as an agent that seeks to maximize its cumulative reward. Although this scenario is conceptually straightforward, it provides a crucial foundation for the subsequent analysis of more complex, irrational agents. We formalize this concept as follows.

Definition 3.1 (Perfectly Rational Agent). Given a subsidy scheme Δr , the action policy $\pi \in \Pi_0(\Delta r)$ of a perfectly rational agent satisfies the constraint

$$V_A^{\pi,\Delta r}(\hat{s}, h=0) \ge \overline{V}_A^{\Delta r}(\hat{s}, h=0).$$

Tie-breaking Rule A tie-breaking rule dictates the agent's choice when multiple actions yield identical rewards. In this setting with a perfectly rational agent, we assume that when two options provide the same personal reward, the agent selects the more cooperative action—that is, the one that benefits the principal more. For example, consider a single state with two actions. Both give the agent a reward of 0, but the principal receives 2 for the first action and 0 for the second. Even a negligible subsidy on the first action makes it strictly preferred. As the subsidy approaches zero, the agent's choice remains the action with a higher principal value. Thus, tie-breaking systematically favors actions that increase the principal's payoff. This assumption allows for a tractable proof of optimality in this section, but it is important to note that we will not rely on this rule in the more general frameworks developed later in the paper.

3.1 OPTIMAL SUBSIDY SCHEME

 With the definition of perfect rationality, we now address the problem of determining the optimal subsidy scheme Δr^* . The following theorem characterizes the principal's optimal payoff and the optimal subsidy scheme. Detailed proof is deferred to Appendix A.3.

Theorem 3.1 (Optimal Subsidy Scheme). For a perfectly rational agent, the principal's optimal payoff is given by

$$V_{sw}^*(\hat{s}, h = 0) - \overline{V}_A^{\Delta r = 0}(\hat{s}, h = 0),$$

that is, the maximum attainable social welfare (over all action policies) minus the maximum reward the agent can obtain in the absence of subsidies. Furthermore, there exists an optimal subsidy scheme Δr^* such that, for every state-action-timestep triple (s, a, h),

$$\Delta r^*(s, a, h) = \overline{V}_A^{\Delta r = 0}(s, h) - \overline{Q}_A^{\Delta r = 0}(s, a, h).$$
 (3.1)

Proof Sketch. The principal's optimal payoff is bounded above by $V_{\rm sw}^*(\hat s,h=0) - \overline{V}_A^{\Delta r=0}(\hat s,h=0)$, since the total value of the principal and agent cannot exceed the maximum possible social welfare, and the agent will not accept less than their stand-alone value without subsidies. This upper bound is achieved under the subsidy scheme Δr^* defined in equation (3.1). Under this scheme, the agent's adjusted Q-values are equalized across all actions: $Q_A^{\Delta r^*}(s,a,h) = \overline{V}_A^{\Delta r=0}(s,h)$ for all (s,a,h). Thus, the agent is indifferent among all actions. Our provisional tie-breaking rule then ensures the agent selects actions that maximize the principal's reward, allowing the principal's payoff to exactly reach the upper bound.

Although Theorem 3.1 identifies an optimal subsidy scheme that provides transfers on nearly all actions, the following proposition shows that, to achieve optimal rewards, the principal needs to subsidize only the social-welfare-maximizing actions. The detailed proof is deferred to Appendix A.4.

Proposition 3.2 (Social Welfare). There exists an optimal subsidy scheme Δr_{sw} that assigns positive transfers exclusively to social-welfare-maximizing actions. Under Δr_{sw} , the agent implements social-welfare-maximizing agent policy π_{sw} , allowing the principal to attain the maximum achievable social welfare.

4 Optimal Policies for Globally ϵ -IC Agents

When an agent is no longer perfectly rational, the optimality of its response ceases to be the sole factor guiding its decisions. To model such bounded rationality, a natural approach is to assume that the agent can tolerate a maximum reward loss of ϵ , in line with the classical notion of ϵ -incentive compatibility (IC). However, since we are dealing with sequential decision-making, several interpretations of ϵ -IC are possible. Here, we focus on the so-called *globally* ϵ -IC agent, which constrains only the cumulative reward loss over the entire decision horizon.

Definition 4.1. An agent is a globally ϵ -IC agent if and only if, given a subsidy scheme Δr , the action policy $\pi \in \Pi^g_{\epsilon}(\Delta r)$ satisfies

$$V_A^{\pi,\Delta r}(\hat{s}, h=0) \ge \overline{V}_A^{\Delta r}(\hat{s}, h=0) - \epsilon.$$

4.1 OPTIMAL SUBSIDY SCHEME

We now consider the problem of determining the optimal subsidy scheme Δr^* . Unlike the perfectly rational case, the agent's best-response policy may be stochastic.

To handle this, we reformulate the objective (2.1) using occupancy measures. Specifically, let $\mu(s,a,h)$ denote the probability that the agent takes action a in state s at timestep h. Replacing the policy π with its corresponding occupancy measure μ , the optimization problem becomes

$$\max_{\Delta r \in \mathcal{R}_{\Delta}} \min_{\mu \in M(\Delta r)} \sum_{s,a,h} \mu(s,a,h) \Big(r_P(s,a,h) - \Delta r(s,a,h) \Big), \tag{4.1}$$

where $M(\Delta r)$ is the set of occupancy measures satisfying the following constraints:

Initial state:
$$\sum_a \mu(\hat{s},a,h=0) = 1, \quad \sum_a \mu(s,a,h=0) = 0 \quad \forall s \neq \hat{s}, \tag{4.2a}$$

Transition:
$$\sum_{a} \mu(s, a, h) = \sum_{s', a'} \mu(s', a', h - 1) P(s|s', a', h - 1),$$
 (4.2b)

Non-negativity:
$$\mu(s, a, h) \ge 0$$
, (4.2c)

Global
$$\epsilon$$
-IC:
$$\sum_{s,a,h} \mu(s,a,h) \left(r_A(s,a,h) + \Delta r(s,a,h) \right) \ge \overline{V}_A^{\Delta r}(\hat{s},h=0) - \epsilon. \tag{4.2d}$$

Directly solving this program is challenging for two main reasons. First, the feasible set of μ is not fixed but depends on the choice of Δr , creating a coupling between the inner and outer variables that distinguishes our setting from standard minimax formulations. Second, defining $f(\Delta r) = \min_{\mu \in M(\Delta r)} \sum_{s,a,h} \mu(s,a,h) \left(r_P(s,a,h) - \Delta r(s,a,h) \right)$ shows that $f(\Delta r)$ is not concave in Δr (see Appendix A.2.1 for example). Consequently, the outer problem $\max_{\Delta r} f(\Delta r)$ is not a concave maximization , which rules out standard convex optimization methods.

In our main theorem, we show the problem can be reformulated to a one-dimensional concave optimization (Theorem 4.1). The approach leverages the dual of the inner optimization problem and swaps the order of optimization between the subsidy scheme Δr and the dual variables (α, V) . The optimal subsidy scheme can then be expressed as the difference between the V-function and Q-function, analogous to the perfectly rational case.

Theorem 4.1. The optimization problem (4.1) is equivalent to maximizing a concave function F(x), formulated as

$$\max_{x \in [0,1)} F(x) = xV_{sw}^*(\hat{s}, h = 0) - V_x^*(\hat{s}, h = 0) - \frac{x}{1 - x}\epsilon,$$

where, for each state s and timestep h, $V_x^*(s,h) \triangleq \max_{\pi} \Big\{ x V_{sw}^{\pi}(s,h) - V_P^{\pi,\Delta r=0}(s,h) \Big\}$.

Furthermore, for an optimal x^* , there exists an optimal subsidy scheme Δr^* such that

$$\Delta r^*(s, a, h) = V_{x^*}^*(s, h) - Q_{x^*}^*(s, a, h)$$
(4.3)

where
$$Q_{x^*}^*(s, a, h) \triangleq x^* r_{sw}(s, a, h) - r_P(s, a, h) + \sum_{s' \in S} P(s'|s, a, h) V_{x^*}^*(s', h+1).$$

Proof. We begin by considering the inner program over the state-action occupancy measure μ for a fixed subsidy scheme Δr . This program is a linear program. By introducing dual variables $\alpha \in \mathbb{R}_+$ for the globally ϵ -IC constraint (4.2d) and $V \in \mathbb{R}^{|\mathcal{S}|(H+1)}$ for the transition (4.2a) and initial state (4.2b) constraints, we can express the problem in its dual form. Combining this with the outer maximization over Δr , α , and V yields the following optimization problem:

$$\max_{\alpha \geq 0, V} V(\hat{s}, h = 0) - \alpha \epsilon + \alpha \max_{\Delta r} \overline{V}_A^{\Delta r}(\hat{s}, h = 0)$$

such that $V(s,h) \leq r_P(s,a,h) - \alpha r_A(s,a,h) - (1+\alpha)\Delta r(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h)V(s',h+1)$ for any $s \in \mathcal{S}$, $a \in \mathcal{A}$, and $h \in \mathcal{H}$; and with the terminal condition V(s,H) = 0 for any state $s \in \mathcal{S}$.

Next, we exchange $\max_{\Delta r}$ and $\max_{\alpha \geq 0, V}$ and analyze maximization over Δr for a fixed V and α . Notice that the objective is non-decreasing with respect to Δr , since $\overline{V}_A^{\Delta r}(\hat{s}, h = 0)$ represents the maximum value attainable by the agent under the subsidy Δr . Additionally, the constraints impose an upper bound on each $\Delta r(s, a, h)$:

$$\Delta r(s, a, h) \le \frac{1}{1 + \alpha} \Big(-V(s, h) + \sum_{s' \in S} P(s'|s, a, h)V(s', h+1) + r_P(s, a, h) - \alpha r_A(s, a, h) \Big).$$

Thus, the optimal choice for Δr is to take this upper bound, making the inequality hold with equality. Given α and V, substituting the optimal choice of Δr , the RHS of the above inequality, into

$$\begin{split} \overline{V}_A^{\Delta r}(\hat{s},h=0) &= \max_{\pi} \mathbb{E}_{\pi} \Big[\sum_{t=0}^{H-1} r_A(s_t,a_t,t) + \Delta r(s_t,a_t,t) \Big] \text{ gives} \\ \overline{V}_A^{\Delta r}(\hat{s},h=0) &= \max_{\pi} \frac{1}{1+\alpha} \mathbb{E}_{\pi} \Big[\sum_{t=0}^{H-1} \left(r_P(s_t,a_t,t) + r_A(s_t,a_t,t) \right) \\ &+ \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t,a_t,t) V(s_{t+1},t+1) - V(s_t,t) \Big] \\ &= \frac{1}{1+\alpha} \Big(V_{\text{sw}}^*(\hat{s},h=0) - V(\hat{s},h=0) \Big). \end{split}$$

Substituting this back, the problem reduces to

$$\begin{split} \max_{\alpha \geq 0} \max_{V} \frac{1}{1+\alpha} V(\hat{s},h=0) + \frac{\alpha}{1+\alpha} V_{\text{sw}}^*(\hat{s},h=0) - \alpha \epsilon \\ \text{s.t.} \quad V(s,h) \leq r_P(s,a,h) - \alpha r_A(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h) V(s',h+1), \\ V(s,H) < 0. \end{split}$$

Observing the inner optimization over V(s,h) coincides with form of minimizing cumulative reward in an MDP with modified reward $r_P - \alpha r_A$. By letting $x = \frac{\alpha}{1+\alpha}$ and introducing $V_x^*(s,h)$ equals $= -\frac{1}{1+\alpha}$ times the optimal value of V(s,h), the formulation equals

$$\begin{split} \max_{x \in (0,1]} \quad x \cdot V_{\text{sw}}^*(\hat{s}, h = 0) - V_x^*(\hat{s}, h = 0) - \frac{x}{1 - x} \epsilon \\ \text{where} \quad V_x^*(\hat{s}, h = 0) &\triangleq -(1 - x) \cdot \min_{\pi} \left\{ V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0) - \frac{x}{1 - x} V_A^{\pi, \Delta r = 0}(\hat{s}, h = 0) \right\} \\ &= \max_{\pi} \big\{ x V_{\text{sw}}^{\pi}(\hat{s}, h = 0) - V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0) \big\}. \end{split}$$

Restricting π to deterministic action policies does not change the value of $V_x^*(\hat{s}, h = 0)$, and under this restriction, $V_x^*(\hat{s}, h = 0)$ is the maximum of finitely many linear functions in x, so the objective function is concave over the interval [0, 1).

Markovian vs. Non-Markovian A process is called **Markovian** if it depends solely on its current state, independent of its past trajectory. Conversely, a process is **non-Markovian** if it can depend on historical states, i.e., it possesses "memory."

In our framework, both the principal and the agent may adopt non-Markovian strategies. For example, the principal might determine subsidies based not only on the agent's current action but also on past actions. Similarly, in equation (4.1), the agent could adopt a non-Markovian globally ϵ -IC policy to reduce the principal's reward. Nevertheless, the following two key observations establish that it suffices to restrict attention to Markovian strategies.

First observation: Given a Markovian subsidy scheme of the principal, there always exists a Markovian globally ϵ -IC policy for the agent that minimizes the principal's reward. This follows from the fact that the inner optimization problem in equation (4.1) is a linear program. Any non-Markovian ϵ -IC policy can be represented by an occupancy measure $\mu(s,a,h)$, which specifies the probability of taking action a in state s at timestep s. Such an occupancy measure can always be replicated by a Markovian policy, ensuring identical rewards for both the principal and the agent.

Second observation: Among all possible subsidy schemes—Markovian or non-Markovian—the Markovian scheme specified in equation (4.3) is optimal. A non-Markovian scheme can be transformed into a Markovian one by augmenting the state space to encode the relevant history. By Theorem 4.1, for each state—action pair in this augmented representation, the scheme in equation (4.3) coincides exactly with its Markovian counterpart.

Remark We briefly examine the boundary cases of x^* and ϵ in Theorem 4.1. When $\epsilon=0$, as $x^*\to 1$, the principal's value approaches $V^*_{\rm sw}(\hat s,h=0)-V^{\Delta r=0}_A(\hat s,h=0)$, consistent with the tie-breaking rule in the perfectly rational case. This shows that the globally ϵ -IC agent naturally generalizes the perfectly rational agent.

4.2 ACTION POLICY

According to Theorem 4.1, the optimal subsidy scheme Δr^* takes a form similar to that in the perfectly rational case. The following proposition shows that the principal can still allocate positive transfers exclusively to the social-welfare-maximizing actions. Furthermore, the agent is still willing to cooperate with the principal to a certain extent by choosing one social-welfare-maximizing agent policy $\pi_{\rm sw}$ with probability x^* , the optimal solution in Theorem 4.1. The detailed proof of the following proposition is deferred to Appendix A.5.

Proposition 4.2 (Optimal subsidy scheme and action policy). There exists an optimal subsidy scheme Δr_{sw} that assigns positive reward transfers solely to social-welfare-maximizing actions. Meanwhile, there exists a globally ϵ -IC action policy $\pi_{\Delta r_{sw}}$ minimizing the principal's reward, which is the mixture of a social-welfare-maximizing agent policy π_{sw} and one other action policy, placing a weight of at least x^* on π_{sw} .

Proof Sketch. The proof relies on two key insights. First, under the optimal subsidy scheme Δr^* , the policy $\pi_{\rm sw}$ achieves the maximum agent expected cumulative reward, $\overline{V}_A^{\Delta r^*}(\hat{s},h=0)$. This implies that it is sufficient to provide subsidies only along the trajectories induced by $\pi_{\rm sw}$, without affecting the optimal value for the principal. Second, there exists an action policy $\hat{\pi}$ whose agent value falls below $\overline{V}_A^{\Delta r^*}(\hat{s},h=0) - \epsilon$, which can be combined with $\pi_{\rm sw}$ to form the globally ϵ -IC policy $\pi_{\Delta r_{\rm sw}}$, such that the dual of the global ϵ -incentive compatibility constraint is tight.

4.3 SOCIAL WELFARE

We define the social welfare gap $\delta_{\rm sw}$ as the difference between the maximum attainable welfare and the welfare achieved under the optimal subsidy scheme Δr^* . When $\epsilon \to +\infty$, the agent can effectively bypass the global ϵ -IC constraint and freely select any action policy. In this limit, the welfare gap becomes $\delta_{\rm sw} = V_{\rm sw}^*(\hat{s},h=0) - \min_\pi V_{\rm sw}^\pi(\hat{s},h=0)$. Our objective is to characterize the upper bound on $\delta_{\rm sw}$ and the rate at which social welfare declines as a function of ϵ , particularly in the regime where ϵ remains small. We first establish the following upper bound on $\delta_{\rm sw}$.

Proposition 4.3. Given ϵ and the corresponding optimal solution $x^* \in (0,1)$, the social welfare gap is $\delta_{sw} = \frac{\epsilon}{1-x^*}$ and it is upper bounded by $O(\sqrt{\epsilon})$.

This $O(\sqrt{\epsilon})$ bound can be achieved in certain specific cases (see Appendix A.2.2 for an example). However, in most cases, the social welfare gap $\delta_{\rm sw}$ exhibits two different growth rates— $O(\sqrt{\epsilon})$ or $O(\epsilon)$ —depending on whether V_x^* is differentiable at x^* . A concrete example is provided below, while detailed discussions are deferred to Appendix A.6.1.

Example Consider a single-period scenario with three actions and $\epsilon=1$. For the first action, the principal's reward is 7 and the agent's reward is 3. For the second action, the principal's reward is 1 and the agent's reward is 2. For the third action, the principal's reward is 1 and the agent's reward is 0. Figure 3a shows that x can grow at rates of $O(\epsilon)$ and $O(\sqrt{\epsilon})$, corresponding to the cases in Figure 3b where x remains constant or grows at $O(\sqrt{\epsilon})$. Figure 3c depicts the piecewise-linear relationship between $V_x^*(\hat{s},h=0)$ and x, where the constant-x value in Figure 3b coincides with the break point of $V_x^*(\hat{s},h=0)$, a non-differentiable point of the objective function.

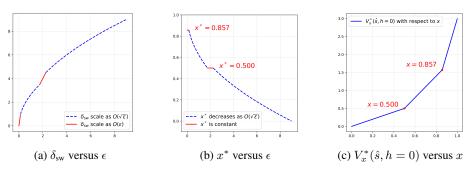


Figure 1: Curves of δ_{sw} and x^* versus ϵ when $V_x^*(\hat{s}, h = 0)$ is non-differentiable.

5 STATE-WISE ϵ -IC AGENT

In this section, we examine the *state-wise* ϵ -*IC agent*, which differs from the globally ϵ -*IC* agent in that incentive compatibility is enforced locally at each state and decision step. Intuitively, such an agent ensures that its chosen action remains within ϵ of the best immediate value available at that decision point. While the idea is simple, constructing a mathematically consistent and tractable formalization is more subtle. We provide two definitions below.

Value-Consistent State-Wise ϵ -**IC Agent** We first define the *value-consistent state-wise* ϵ -**IC** *agent*, where the agent's action at each state must approximate the optimal reward within ϵ .

Definition 5.1. An agent is a value-consistent state-wise ϵ -IC agent if, under a subsidy scheme Δr , the induced policy $\pi \in \Pi^v_{\epsilon}(\Delta r)$ satisfies $V^{\pi,\Delta r}_A(s,h) \geq \overline{V}^{\Delta r}_A(s,h) - \epsilon$ for all $s \in \mathcal{S}$ and $h \in \mathcal{H}$.

A key challenge with this formulation is that the agent's policy minimizing the principal's reward under a given subsidy scheme may be **non-Markovian**. In such cases, the agent's policy cannot be represented within polynomial size.

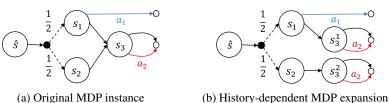


Figure 2: Illustration of value-consistent state-wise ϵ -IC agents.

To illustrate, consider the post-subsidy MDP in Figure 2a, where (i) for action a_1 at s_1 : principal reward 100, agent reward 3; (ii) for action a_2 at s_3 : principal reward 2, agent reward 2; and (iii) for all other actions: reward 0. Under a **Markovian policy**, the value-consistent state-wise ϵ -IC agent minimizes the principal's reward by selecting a_2 at s_3 , and steering toward s_3 from s_1 . This yields a principal reward of 2. However, under a **non-Markovian policy**, we can duplicate s_3 into two history-dependent states, s_3^1 and s_3^2 . At s_3^1 , the agent always selects a_2 , while at s_3^2 , the agent mixes between two actions with equal probability. This reduces the principal's expected reward to 1.5.

Greedy State-Wise ϵ -IC Agent To avoid non-Markovian behavior, we introduce the *greedy state-wise* ϵ -IC agent, which replaces recursive value computations with greedy look-ahead. Once the subsidy scheme is fixed, $\overline{V}_A^{\Delta r}$ becomes deterministic, and the agent greedily minimizes the principal's value through local decisions.

Definition 5.2. An agent is a greedy state-wise ϵ -IC agent if, under subsidy scheme Δr , the induced policy $\pi \in \Pi^s_{\epsilon}(\Delta r)$ satisfies, for all $s \in S$, $h \in \mathcal{H}$:

$$\sum_{a \in \mathcal{A}} \pi(a|s,h) \Big(r_A^{\Delta r}(s,a,h) + \sum_{s' \in \mathcal{S}} P(s'|s,a,h) \, \overline{V}_A^{\Delta r}(s',h+1) \Big) \; \geq \; \overline{V}_A^{\Delta r}(s,h) - \epsilon.$$

However, even in this simplified greedy setting, designing the principal's optimal subsidy scheme remains computationally intractable. The complete proof is deferred to Appendix A.7.

Theorem 5.1. Given a greedy state-wise ϵ -IC agent, computing the principal's optimal subsidy scheme is NP-hard.

6 Conclusion

In this paper, we study a principal-agent problem with the aim of designing a robust subsidy scheme that maximizes the cumulative expected return in the presence of an irrational agent. We demonstrate that, under the globally ϵ -IC assumption, the optimal subsidy scheme can be effectively determined, representing a natural extension of the perfectly rational case. We further show that formulating the state-wise ϵ -IC follower is computationally challenging. As future work, it would be interesting to consider scenarios in which the principal does not have prior knowledge of the agent's reward function or the value of ϵ , such as in a learning-based setting.

REFERENCES

- Kiarash Banihashem, Adish Singla, Jiarui Gan, and Goran Radanovic. Admissible policy teaching through reward design. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6037–6045, 2022.
- Omer Ben-Porat, Yishay Mansour, Michal Moshkovitz, and Boaz Taitler. Principal-agent reward shaping in mdps. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 9502–9510, 2024.
- Martino Bernasconi, Matteo Castiglioni, Alberto Marchesi, and Mirco Mutti. Persuading farsighted receivers in mdps: the power of honesty. *Advances in Neural Information Processing Systems*, 36:14987–15014, 2023.
- Matteo Bollini, Francesco Bacchiocchi, Matteo Castiglioni, Alberto Marchesi, and Nicola Gatti. Contracting with a reinforcement learning agent by playing trick or treat. *arXiv preprint arXiv:2410.13520*, 2024.
- Paul Dütting, Michal Feldman, and Inbal Talgam-Cohen. Algorithmic contract theory: A survey. *Foundations and Trends® in Theoretical Computer Science*, 16(3-4):211–412, 2024.
- Jiarui Gan, Rupak Majumdar, Goran Radanovic, and Adish Singla. Bayesian persuasion in sequential decision-making. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 5025–5033, 2022.
- Sanford J Grossman and Oliver D Hart. An analysis of the principal-agent problem. In *Foundations of insurance economics: Readings in economics and finance*, pp. 302–340. Springer, 1992.
- Guru Guruganesh, Jon Schneider, and Joshua R Wang. Contracts under moral hazard and adverse selection. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pp. 563–582, 2021.
- Dima Ivanov, Paul Dütting, Inbal Talgam-Cohen, Tonghan Wang, and David C Parkes. Principal-agent reinforcement learning: Orchestrating ai agents with contracts. *arXiv preprint arXiv:2407.18074*, 2024.
- Ohad Kadan, Philip J Reny, and Jeroen M Swinkels. Existence of optimal mechanisms in principal-agent problems. *Econometrica*, 85(3):769–823, 2017.
- Jean-Jacques Laffont and Eric Maskin. *The theory of incentives: An overview*. Université des sciences sociales, Faculté des sciences économiques, 1981.
- Paul Milgrom and Ilya Segal. Envelope theorems for arbitrary choice sets. *Econometrica*, 70(2): 583–601, 2002.
- Roger B Myerson. Optimal coordination mechanisms in generalized principal–agent problems. *Journal of mathematical economics*, 10(1):67–81, 1982.
- Stephen A Ross. The economic theory of agency: The principal's problem. *The American economic review*, 63(2):134–139, 1973.
- Vinzenz Thoma, Barna Pásztor, Andreas Krause, Giorgia Ramponi, and Yifan Hu. Contextual bilevel reinforcement learning for incentive alignment. *Advances in Neural Information Processing Systems*, 37:127369–127435, 2024.
- Jibang Wu, Zixuan Zhang, Zhe Feng, Zhaoran Wang, Zhuoran Yang, Michael I Jordan, and Haifeng Xu. Sequential information design: Markov persuasion process and its efficient reinforcement learning. *arXiv* preprint arXiv:2202.10678, 2022.
- Jibang Wu, Siyu Chen, Mengdi Wang, Huazheng Wang, and Haifeng Xu. Contractual reinforcement learning: Pulling arms with invisible hands. *arXiv preprint arXiv:2407.01458*, 2024.
- Shuo Wu, Haoxiang Ma, Jie Fu, and Shuo Han. Robust reward design for markov decision processes. *Journal of Artificial Intelligence Research*, 84, 2025.

Guanghui Yu and Chien-Ju Ho. Environment design for biased decision makers. In IJCAI, pp. 592–598, 2022.

Haoqi Zhang and David C Parkes. Value-based policy teaching with active indirect elicitation. In AAAI, volume 8, pp. 208–214, 2008.

APPENDIX

USAGE OF LLM

We employed the large language model (LLM) to assist in refining the language and enhancing the clarity of this manuscript. The LLM was **not** used for generating research ideas, identifying related work, performing analyses, or contributing to the substantive scientific content of this paper.

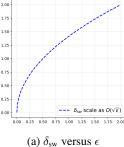
MISSING EXAMPLES

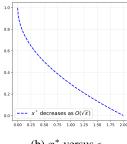
A.2.1COUNTEREXAMPLE ON CONVEXITY

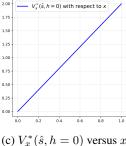
Consider a single-period scenario with three actions and $\epsilon = 2$. The principal receives 0 for the first action and 5 for both the second and third actions, while the agent's reward is 0 for the first action and 1 for other actions. Let Δr_1 and Δr_2 transfer 2 units to the second and third actions, respectively; both yield a principal value of 1. A convex combination, Δr_3 , transferring 1 unit to both actions, results in a leader value of 0, directly violating concavity. This illustrates that the outer optimization cannot be assumed concave.

A.2.2 Example of δ_{SW} Scaling as $O(\sqrt{\epsilon})$

Consider a single-period scenario with two actions and $\epsilon = 1$. The principal and agent values for the first action are 4 and 0, respectively, and for the second action, they are 0 and 2. In this setting, $1-x^*$ always scales as $O(\sqrt{\epsilon})$ and matches the upper bound. The core idea behind is in such instance, the function $V_r^*(\hat{s}, h = 0)$ is a complete linear function in interval [0, 1]. Figure 3 illustrates relationship between $\delta_{\rm sw}$, x^* , and ϵ , along with the behavior of $V_x^*(\hat{s}, h = 0)$ as a function of x.







(b) x^* versus ϵ Figure 3: Curves of $\delta_{\rm sw}$ and x^* versus ϵ when $V_x^*(\hat{s}, h = 0)$ is differentiable.

A.3 Proof of Theorem 3.1

By definition, for any subsidy scheme Δr with induced action policy $\pi_{\Delta r} \in \Pi_0(\Delta r)$, we have

$$V_P^{\pi_{\Delta r}, \Delta r}(\hat{s}, h = 0) + V_A^{\pi_{\Delta r}, \Delta r}(\hat{s}, h = 0) = V_{\text{sw}}^{\pi_{\Delta r}}(\hat{s}, h = 0) \le V_{\text{sw}}^*(\hat{s}, h = 0). \tag{A.1}$$

Moreover, since any subsidy scheme provides the agent with non-negative reward transfers, backward induction gives

$$V_A^{\pi,\Delta r}(\hat{s}, h = 0) = \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{H-1} r_A^{\Delta r}(s_t, a_t, t) \right]$$

$$\geq \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{H-1} r_A(s_t, a_t, t) \right]$$

$$= \overline{V}_A^{\Delta r = 0}(\hat{s}, h = 0).$$
(A.2)

Combining this with inequality (A.1), the optimal principal value is upper bounded by

OPT
$$\leq V_{\text{sw}}^*(\hat{s}, h = 0) - \overline{V}_A^{\Delta r = 0}(\hat{s}, h = 0).$$

It remains to show that

$$\Delta r^*(s,a,h) = \overline{V}_A^{\Delta r=0}(s,h) - \overline{Q}_A^{\Delta r=0}(s,a,h)$$

achieves this bound. For any policy π , substituting Δr^* into $V_A^{\pi,\Delta r^*}(\hat{s},h=0)$ and applying backward induction establishes that the agent's value for every action equals $V_A^{\Delta r=0}(\hat{s},h=0)$, which makes inequality (A.2) tight. In addition, since the social-welfare-maximizing policy $\pi_{\rm sw}$ renders inequality (A.1) exact, the principal's value under Δr^* is

$$V_P^{\pi_{\text{sw}},\Delta r^*}(\hat{s}, h = 0) = V_{\text{sw}}^*(\hat{s}, h = 0) - \overline{V}_A^{\Delta r = 0}(\hat{s}, h = 0),$$

which coincides with the upper bound. Consequently, under the cooperative tie-breaking rule, the agent selects π_{sw} , thereby achieving

OPT =
$$V_{sw}^*(\hat{s}, h = 0) - \overline{V}_A^{\Delta r = 0}(\hat{s}, h = 0).$$

A.4 Proof for Proposition 3.2

Recall that in Theorem 3.1, we defined the optimal reward transfer as

$$\Delta r^*(s, a, h) = \overline{V}_A^{\Delta r = 0}(s, h) - \overline{Q}_A^{\Delta r = 0}(s, a, h).$$

In fact, it suffices to retain the reward transfer only along the social-welfare-maximizing actions. In particular, we define subsidy scheme $\Delta r_{\rm sw}^*$ as

$$\Delta r_{\rm sw}(s,a,h) = \begin{cases} \Delta r^*(s,a,h), & \text{if } \pi_{\rm sw}(a|s,h) > 0, \\ 0, & \text{otherwise.} \end{cases}$$

so that the agent value of social-welfare-maximizing action policy $\pi_{\rm sw}$ under $\Delta r_{\rm sw}$ is

$$V_A^{\pi_{\text{sw}}, \Delta r_{\text{sw}}}(\hat{s}, h = 0) = \overline{V}_A^{\Delta r^*}(\hat{s}, h = 0) = \overline{V}_A^{\Delta r = 0}(\hat{s}, h = 0),$$

which yields a principal value of $V_{\rm sw}^*(\hat s,h=0) - \overline V_A^{\Delta r=0}(\hat s,h=0).$

A.5 PROOF OF PROPOSITION 4.2

Optimal subsidy scheme Recall from (4.3) that

$$\begin{split} \Delta r^*(s,a,h) &= V_{x^*}^*(s,h) - Q_{x^*}^*(s,a,h) \\ &= V_{x^*}^*(s,h) - x^* r_{\text{sw}}(s,a,h) + r_P(s,a,h) - \sum_{s' \in \mathcal{S}} P(s'|s,a,h) V_{x^*}^*(s',h+1) \\ &= (1-x^*) r_P(s,a,h) - x^* r_A(s,a,h) + V_{x^*}^*(s,h) - \sum_{s' \in \mathcal{S}} P(s'|s,a,h) V_{x^*}^*(s',h+1). \end{split}$$

Let π_{sw} be a deterministic social-welfare-maximizing agent policy. We now define a subsidy scheme Δr_{sw} that is restricted to π_{sw} :

$$\Delta r_{\mathrm{sw}}(s, a, h) \triangleq \begin{cases} \Delta r^*(s, a, h), & \text{if } \pi_{\mathrm{sw}}(a|s, h) > 0, \\ 0, & \text{otherwise.} \end{cases}$$

First, we claim that under either subsidy scheme, $\Delta r_{\rm sw}$ or Δr^* , the agent's optimal value is identical:

$$\overline{V}_A^{\Delta r^*}(\hat{s}, h = 0) = \overline{V}_A^{\Delta r_{\rm sw}}(\hat{s}, h = 0).$$

The central argument of the proof is that, under both Δr^* and $\Delta r_{\rm sw}$, action policy $\pi_{\rm sw}$ can achieve the agent's maximal value, and the values thus attained coincide. Specifically,

$$\overline{V}_A^{\Delta r^*}(\hat{s},h=0) \stackrel{(a)}{=} V_A^{\pi_{\mathrm{sw}},\Delta r^*}(\hat{s},h=0) \stackrel{(b)}{=} V_A^{\pi_{\mathrm{sw}},\Delta r_{\mathrm{sw}}}(\hat{s},h=0) \stackrel{(c)}{=} \overline{V}_A^{\Delta r_{\mathrm{sw}}}(\hat{s},h=0).$$

In what follows, we establish the validity of each equality (a)-(c) sequentially.

We first show (a):

$$\overline{V}_A^{\Delta r^*}(\hat{s},h=0) = V_A^{\pi_{\mathrm{sw}},\Delta r^*}(\hat{s},h=0).$$

To see this, consider any action policy π under Δr^* ,

$$V_A^{\pi,\Delta r^*}(\hat{s}, h = 0) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{H-1} r_A(s_t, a_t, t) + \Delta r^*(s_t, a_t, t) \right]$$

$$= \mathbb{E}_{\pi} \left[\sum_{t=0}^{H-1} r_A(s_t, a_t, t) + (1 - x^*) r_P(s_t, a_t, t) - x^* r_A(s_t, a_t, t) \right]$$

$$+ V_{x^*}^*(s_t, t) - \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t, t) V_{x^*}^*(s_{t+1}, t+1) \right]$$

$$= (1 - x^*) V_{\text{ew}}^{\pi}(\hat{s}, h = 0) + V_{x^*}^*(\hat{s}, h = 0).$$
(A.3)

Subtracting the agent value under any policy π from that under π_{sw} gives

$$\begin{split} V_A^{\pi,\Delta r^*}(\hat{s},h=0) - V_A^{\pi_{\text{sw}},\Delta r^*}(\hat{s},h=0) \\ &= (1-x^*)(V_{\text{sw}}^{\pi}(\hat{s},h=0) - V_{\text{sw}}^{\pi_{\text{sw}}}(\hat{s},h=0)) \\ &\leq 0, \end{split} \tag{A.4}$$

which implies $\pi_{\rm sw}$ can achieves the maximum agent value under Δr^* :

$$V_A^{\pi_{\text{sw}}, \Delta r^*}(\hat{s}, h = 0) = \overline{V}_A^{\Delta r^*}(\hat{s}, h = 0).$$
 (A.5)

Then, we show (b):

$$V_A^{\pi_{\text{sw}},\Delta r^*}(\hat{s}, h = 0) = V_A^{\pi_{\text{sw}},\Delta r_{\text{sw}}}(\hat{s}, h = 0).$$

Based on the definition of $\Delta r_{\rm sw}$, we can express $\Delta r_{\rm sw}$ for any triple (s, a, h) as

$$\Delta r_{\rm sw}(s, a, h) = \mathbf{1}_{\{x>0\}}(\pi_{\rm sw}(a|s, h))\Delta r(s, a, h),$$

where $\mathbf{1}_{\{x>0\}}(x)$ is the indicator function:

$$\mathbf{1}_{\{x>0\}} = \begin{cases} 1, & x > 0, \\ 0, & x \le 0. \end{cases}$$

By introducing the occupancy measure μ_{sw} of policy π_{sw} , which satisfies $(\sum_{a'} \mu_{sw}(s, a', h)) \cdot \pi_{sw}(a|s, h) = \mu_{sw}(s, a, h)$ for any (s, a, h), we have $\mathbf{1}_{\{x>0\}}(\pi_{sw}(a|s, h)) \geq \mathbf{1}_{\{x>0\}}(\mu_{sw}(s, a, h))$.

Hence,

$$\begin{split} V_A^{\pi_{\text{sw}},\Delta r^*}(\hat{s},h &= 0) = \sum_{s,a,h} \mu_{\text{sw}}(s,a,h) \big(r_A(s,a,h) + \Delta r^*(s,a,h) \big) \\ &= \sum_{s,a,h} \Big(\mu_{\text{sw}}(s,a,h) r_A(s,a,h) + \mathbf{1}_{\{x>0\}} (\mu_{\text{sw}}(s,a,h)) \cdot \mu_{\text{sw}}(s,a,h) \Delta r^*(s,a,h) \Big) \\ &\leq \sum_{s,a,h} \Big(\mu_{\text{sw}}(s,a,h) r_A(s,a,h) + \mathbf{1}_{\{x>0\}} (\pi_{\text{sw}}(a|s,h)) \cdot \mu_{\text{sw}}(s,a,h) \Delta r^*(s,a,h) \Big) \\ &= \sum_{s,a,h} \Big(\mu_{\text{sw}}(s,a,h) r_A(s,a,h) + \mu_{\text{sw}}(s,a,h) \Delta r_{\text{sw}}(s,a,h) \Big) \\ &= V_A^{\pi_{\text{sw}},\Delta r_{\text{sw}}}(\hat{s},h = 0). \end{split}$$

Meanwhile, since $\Delta r^*(s, a, h) \ge \Delta r_{sw}(s, a, h)$ for any (s, a, h), by construction, it follows that

$$V_A^{\pi_{\text{sw}}, \Delta r^*}(\hat{s}, h = 0) \ge V_A^{\pi_{\text{sw}}, \Delta r_{\text{sw}}}(\hat{s}, h = 0).$$

Combining the above results, we can conclude that

$$V_A^{\pi_{\text{sw}},\Delta r^*}(\hat{s}, h=0) = V_A^{\pi_{\text{sw}},\Delta r_{\text{sw}}}(\hat{s}, h=0).$$

Finally, we establish (c):

$$V_A^{\pi_{\text{sw}},\Delta r_{\text{sw}}}(\hat{s}, h = 0) = \overline{V}_A^{\Delta r_{\text{sw}}}(\hat{s}, h = 0).$$

To prove this, it suffices to show that for any policy π , $V_A^{\pi,\Delta r_{\rm sw}}(\hat{s},h=0) \leq V_A^{\pi_{\rm sw},\Delta r_{\rm sw}}(\hat{s},h=0)$, which directly implies that $\pi_{\rm sw}$ attains the maximum agent value under $\Delta r_{\rm sw}$, i.e., $\overline{V}_A^{\Delta r_{\rm sw}}(\hat{s},h=0) = V_A^{\pi_{\rm sw},\Delta r_{\rm sw}}(\hat{s},h=0)$.

The inequality, $V_A^{\pi,\Delta r_{\rm sw}}(\hat{s},h=0) \leq V_A^{\pi_{\rm sw},\Delta r_{\rm sw}}(\hat{s},h=0)$, follows from the chain

$$V_A^{\pi,\Delta r_{\text{sw}}}(\hat{s},h=0) \leq V_A^{\pi,\Delta r^*}(\hat{s},h=0) \leq V_A^{\pi_{\text{sw}},\Delta r^*}(\hat{s},h=0) = V_A^{\pi_{\text{sw}},\Delta r_{\text{sw}}}(\hat{s},h=0).$$

The first inequality holds because the construction, $\Delta r_{\rm sw}(s,a,h) \leq \Delta r^*(s,a,h)$ for all (s,a,h), which implies, for any action policy π , $V_A^{\pi,\Delta r_{\rm sw}}(\hat{s},h=0) \leq V_A^{\pi,\Delta r^*}(\hat{s},h=0)$. The second inequality follows from (A.4). The final equality holds because $\Delta r^*(s,a,h) = \Delta r_{\rm sw}(s,a,h)$ whenever $\pi_{\rm sw}(a|s,h)>0$.

By combining (a), (b), and (c), we can conclude $\overline{V}_A^{\Delta r^*}(\hat{s},h=0)=\overline{V}_A^{\Delta r_{\rm sw}}(\hat{s},h=0)$ follows.

Next, we show that the principal's worst-case reward under $\Delta r_{\rm sw}$ is no worse than under Δr^* . Suppose an action policy π is globally ϵ -IC under $\Delta r_{\rm sw}$. By definition, this implies that the agent's value under $\Delta r_{\rm sw}$ achieves $\overline{V}_A^{\Delta r_{\rm sw}}(\hat{s},h=0)-\epsilon$. From the previous discussion, we know that $\overline{V}_A^{\Delta r^*}(\hat{s},h=0)=\overline{V}_A^{\Delta r_{\rm sw}}(\hat{s},h=0)$ and $\Delta r_{\rm sw}\leq \Delta r^*$. Consequently, the action policy π is also globally ϵ -IC under Δr^* . The converse, however, does not necessarily hold.

Thus, relative to Δr^* , the scheme $\Delta r_{\rm sw}$ reduces the set of agent policies that are globally ϵ -IC. Moreover, for any given action policy, the principal's payoff under $\Delta r_{\rm sw}$ is at least as large as under Δr^* . It follows that the principal attains at least the same worst-case reward under $\Delta r_{\rm sw}$ as under Δr^* . Since Δr^* is optimal, the principal's worst-case reward is identical under both schemes, thereby establishing the first part of the proposition.

Action Policy under Δr^* To prove the latter part of Proposition 4.2, the key idea is to analyze the agent's response under a given subsidy scheme Δr . Recall from the preliminaries that we define the adversarial response $\pi_{\Delta r}$ to a subsidy scheme Δr as

$$\pi_{\Delta r} \in \operatorname*{arg\,min}_{\pi \in \Pi(\Delta r)} V_P^{\pi,\Delta r}(\hat{s},h=0).$$

In the following discussion, we refer to $\pi_{\Delta r}$ as the agent's adversarial action policy. Further, we begin by analyzing the agent's behavior in a single-period instance and then extend the results to the

multi-period case. Formally, when H=1, there is only one state \hat{s} , so we can omit (s,h) in the expressions. The optimization problem for agent then becomes

$$\min_{\pi} \sum_{a} \pi(a) r_P^{\Delta r}(a) \quad \text{s.t. } \sum_{a} \pi(a) r_A^{\Delta r}(a) \geq \max_{a} r_A^{\Delta r}(a) - \epsilon, \pi(a) \geq 0, \ \sum_{a} \pi(a) = 1.$$

As this is a linear program, we apply the KKT conditions to analyze the optimal solution. The Lagrangian function is

$$\begin{split} \mathcal{L}(\pi;\alpha,\beta,V) &= \sum_{a} \pi(a) r_P^{\Delta r}(a) + \alpha \Big(\max_{a} r_A^{\Delta r}(a) - \epsilon - \sum_{a} \pi(a) r_A^{\Delta r}(a) \Big) \\ &+ \sum_{a} \beta(a) (-\pi(a)) + V \Big(1 - \sum_{a} \pi(a) \Big) \\ &= \sum_{a} \pi(a) \Big(r_P^{\Delta r}(a) - \alpha r_A^{\Delta r}(a) - V + \beta(a) \Big) + V + \alpha \Big(\max_{a} r_A^{\Delta r}(a) - \epsilon \Big). \end{split}$$

The resulting dual program is as followed:

$$\max_{\alpha, V} V + \alpha \Big(\max_{a} r_A^{\Delta r}(a) - \epsilon \Big) \quad \text{s.t.} \quad \alpha \ge 0, \ V \le r_P(a) - \alpha r_A(a).$$

Let $\alpha^{\Delta r}$ and $V^{\Delta r}$ denote the optimal dual values under the subsidy scheme Δr , and let $\text{OTP}^{\Delta r}$ denote the final principal value under the same subsidy scheme. By complementary slackness, for any action a such that

$$r_P^{\Delta r}(a) - \alpha^{\Delta r} r_A^{\Delta r}(a) = V^{\Delta r} = \min_{a} \left(r_P^{\Delta r}(a) - \alpha^{\Delta r} r_A^{\Delta r}(a) \right),$$

we have $\pi_{\Delta r}(a) \geq 0$. We refer to such actions as the *candidate actions* $a \in \overline{A}$, since they can potentially be chosen by the agent given subsidy scheme Δr .

However in certain problem instances, there exist candidate actions that do not appear in any agent's adversarial action policy. For example, consider $\epsilon=1$ and two actions: the first action has principal reward 0 and agent reward 1, while the second action has principal reward 1 and agent reward 2. The unique agent's adversarial action policy deterministically selects the first action, yet setting $\alpha=2$ would include both actions as candidate actions. In general, based on the value of optimal dual variable α^* under optimal subsidy scheme Δr^* , we claim there are three possible scenarios:

- Case 1: $\alpha^* = 0$. Every candidate action at this point attains the minimum principal value, so any action satisfying the globally ϵ -IC constraint can be deterministically chosen.
- Case 2: $\alpha^* > 0$, and all candidate actions satisfy the globally ϵ -IC constraint. Complementary slackness implies the agent's value is exactly $\max_a \overline{r}_A(a) \epsilon$, so only actions attaining this value can be chosen. This case coincides with the above example.
- Case 3: $\alpha^* > 0$, and some candidate actions have agent reward below $\max_a \overline{r}_A(a) \epsilon$. Then, the agent can mix actions above and below this threshold to form an agent's adversarial action policy, leading to the fact that any candidate action may become part of agent response.

To further explain the agent's behavior pattern in Case 2 and Case 3, we use the following equations to show that when $\alpha^{\Delta r} > 0$, as long as the action policy distribution π is supported only on candidate actions and achieves an agent value exactly equal to $\max_a \overline{r}_A(a) - \epsilon$, the policy π constitutes one possible adversarial action policy of the agent. In other words, we only need to consider how to organize the policy distribution supported on candidate actions so as to achieve an agent value exactly equal to $\max_a \overline{r}_A(a) - \epsilon$.

$$\begin{split} \sum_{a} \pi(a) r_P^{\Delta r}(a) &= \sum_{a} \pi(a) (r_P^{\Delta r}(a) - \alpha^{\Delta r} r_A^{\Delta r}(a)) + \alpha^{\Delta r} \sum_{a} \pi(a) r_A^{\Delta r}(a) \\ &= \sum_{a} \pi(a) V^{\Delta r} + \alpha \Delta r (\max_{a} r_A^{\Delta r}(a) - \epsilon) \\ &= \text{OPT}^{\Delta r} \end{split}$$

We now extend to the multi-period case. Observe that under subsidy scheme Δr , any multi-period policy π_i can be viewed as a single-period action a_i with principal reward $V_P^{\pi_i,\Delta r}(\hat{s},h=0)$ and agent reward $V_A^{\pi_i,\Delta r}(\hat{s},h=0)$, yielding a single-period instance with infinitely many actions. Although this transformation is generally intractable, it provides a useful framework for analyzing the properties of action policies. Suppose the single-period agent' adversarial action policy is π_s^* , the multi-period agent's adversarial action policy can be recovered as $\pi^*(a|s,h) = \sum_i \pi_s^*(a_i) \pi_i(a|s,h)$ for any state s, action a and timestep b, where $\pi_s^*(a_i)$ denotes the probability assigned to policy π_i . Accordingly, we define the *candidate policy* $\overline{\pi} \in \overline{\Pi}$ under subsidy scheme Δr as

$$V_P^{\overline{\pi},\Delta r}(\hat{s},h=0) - \alpha^{\Delta r} V_A^{\overline{\pi},\Delta r}(\hat{s},h=0) = \min_{\pi} \big\{ V_P^{\pi,\Delta r}(\hat{s},h=0) - \alpha^{\Delta r} V_A^{\pi,\Delta r}(\hat{s},h=0) \big\},$$

analogous to the single-period candidate actions, so that the adversarial agent's action policy can be expressed as a convex combination of these candidate policies.

We next establish that, under the optimal subsidy scheme Δr^* from Theorem 4.1, **every** action policy π , including π_{sw} , qualifies as a candidate. Recall from (4.3) that

$$\begin{split} \Delta r^*(s,a,h) &= V_{x^*}^*(s,h) - Q_{x^*}^*(s,a,h) \\ &= V_{x^*}^*(s,h) - x^* r_{\text{sw}}(s,a,h) + r_P(s,a,h) - \sum_{s' \in \mathcal{S}} P(s'|s,a,h) V_{x^*}^*(s',h+1) \\ &= (1-x^*) r_P(s,a,h) - x^* r_A(s,a,h) + V_{x^*}^*(s,h) - \sum_{s' \in \mathcal{S}} P(s'|s,a,h) V_{x^*}^*(s',h+1), \end{split}$$

where $x^* = \frac{\alpha^*}{1+\alpha^*}$. For any action policy, under optimal subsidy scheme Δr^* substituting Δr^* into the expression $V_P^{\pi,\Delta r}(\hat{s},h=0) - \alpha^* V_A^{\pi,\Delta r}(\hat{s},h=0)$ yields

$$V_{P}^{\pi,\Delta r^{*}}(\hat{s}, h = 0) - \alpha^{*}V_{A}^{\pi,\Delta r^{*}}(\hat{s}, h = 0)$$

$$= \mathbb{E}_{\pi} \left[\sum_{t=0}^{H-1} r_{P}(s_{t}, a_{t}, t) - \alpha^{*}r_{A}(s_{t}, a_{t}, t) - (1 + \alpha^{*})\Delta r^{*}(s_{t}, a_{t}, t) \right]$$

$$= \mathbb{E}_{\pi} \left[\sum_{t=0}^{H-1} r_{P}(s_{t}, a_{t}, t) - \frac{x^{*}}{1 - x^{*}}r_{A}(s_{t}, a_{t}, t) - \frac{1}{1 - x^{*}}\Delta r^{*}(s_{t}, a_{t}, t) \right]$$

$$= \mathbb{E}_{\pi} \left[\sum_{t=0}^{H-1} r_{P}(s_{t}, a_{t}, t) - \frac{x^{*}}{1 - x^{*}}r_{A}(s_{t}, a_{t}, t) - \frac{1}{1 - x^{*}}((1 - x^{*})r_{P}(s_{t}, a_{t}, t) - x^{*}r_{A}(s_{t}, a_{t}, t)) \right]$$

$$- \frac{1}{1 - x^{*}} \left(V_{x^{*}}^{*}(s_{t}, t) - \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_{t}, a_{t}, t + 1)V_{x^{*}}^{*}(s_{t+1}, t + 1) \right)$$

$$= -\frac{1}{1 - x^{*}}V_{x^{*}}^{*}(\hat{s}, h = 0).$$
(A.6)

Thus, every action policy π yields the same value

$$V_P^{\pi,\Delta r^*}(\hat{s}, h = 0) - \alpha^* V_A^{\pi,\Delta r^*}(\hat{s}, h = 0).$$

Consequently, all action policies qualify as candidate policies, and we have

$$\min_{\pi} \left\{ V_P^{\pi, \Delta r^*}(\hat{s}, h = 0) - \alpha^* V_A^{\pi, \Delta r^*}(\hat{s}, h = 0) \right\} = -\frac{1}{1 - x^*} V_{x^*}^*(\hat{s}, h = 0).$$

Having established this, we note that not all candidate policies necessarily receive positive probability in the support of an agent's adversarial action policy, as illustrated in the single-period analysis. Similarly, we distinguish cases based on the value of the optimal dual variable α^* : when $\alpha^* = 0$, the corresponding optimal solution is $x^* = 0$, making the proposition trivial; When $\alpha^* > 0$, we

assert that Case 2 will never happen under optimal subsidy scheme Δr^* as there exists a candidate policy $\hat{\pi}$ whose agent value is strictly smaller than $\overline{V}_A^{\Delta r^*}(\hat{s},h=0)-\epsilon$. Moreover, this policy $\hat{\pi}$ can be combined with $\pi_{\rm sw}$ to construct the agent's adversarial action policy.

In detail, we construct the candidate policy $\hat{\pi}$ by analyzing the derivative of the objective function. Applying the envelope theorem (Milgrom & Segal, 2002), the derivative of the objective with respect to x in Theorem 4.1 is given by

$$F'(x) = V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_x}(\hat{s}, h = 0) - \frac{\epsilon}{(1 - x)^2},$$
(A.7)

where $\pi_x = \arg\max_{\pi} \{xV_{\rm sw}^{\pi}(\hat{s},h=0) - V_P^{\pi,\Delta r=0}(\hat{s},h=0)\}$. In general, the objective function may have a finite number of non-differentiable points, arising from the potential non-uniqueness of $V_{\rm sw}^{\pi_x}(\hat{s},h=0)$. Nevertheless, since the set of sub-differentials can be fully characterized by the derivative expression, for simplicity and clarity we do not distinguish between derivatives and sub-derivatives, and we treat stationary points by directly setting F'(x)=0.

Since the objective function is concave, the optimal solution x^* and the corresponding policy π_{x^*} can be characterized by the vanishing derivative condition. Furthermore, the requirement $\alpha^* > 0$ ensures that $x^* \in (0,1)$, implying that x^* , as an interior optimum, necessarily exists as a stationary point. Consequently, imposing F'(x) = 0 yields

$$V_{\rm sw}^*(\hat{s},h=0) - V_{\rm sw}^{\pi_x*}(\hat{s},h=0) \; = \; \frac{\epsilon}{(1-x^*)^2}.$$

In general, π_{x^*} can be represented as any convex combination of some action policies $\widetilde{\pi}_x$ maximizing $xV_{\rm sw}^{\widetilde{\pi}_x}(\hat{s},h=0)-V_P^{\widetilde{\pi}_x,\Delta r=0}(\hat{s},h=0)$. Since every action policy qualifies as a candidate policy under Δr^* , π_{x^*} can equivalently be viewed as a convex combination of candidate policies. Hence, there exists a candidate policy $\hat{\pi}$ such that

$$\hat{\pi} = \arg\max_{\pi} \{ x V_{\text{sw}}^{\pi}(\hat{s}, h = 0) - V_{P}^{\pi, \Delta r = 0}(\hat{s}, h = 0) \},$$

with

$$V_{\rm sw}^*(\hat{s},h=0) - V_{\rm sw}^{\hat{\pi}}(\hat{s},h=0) \ge \frac{\epsilon}{(1-x^*)^2}.$$

To upper bound the agent value of $\hat{\pi}$, using equations (A.3) and (A.5), we deduce that for action policy $\hat{\pi}$,

$$\overline{V}_{A}^{\Delta r^{*}}(\hat{s}, h = 0) - V_{A}^{\hat{\pi}, \Delta r^{*}}(\hat{s}, h = 0) = (1 - x^{*}) \Big(V_{\text{sw}}^{*}(\hat{s}, h = 0) - V_{\text{sw}}^{\hat{\pi}}(\hat{s}, h = 0) \Big), \tag{A.8}$$

which, combined with the preceding inequality, implies

$$\overline{V}_A^{\Delta r^*}(\hat{s}, h = 0) - V_A^{\hat{\pi}, \Delta r^*}(\hat{s}, h = 0) \ge \frac{\epsilon}{1 - r^*} \ge \epsilon.$$

Thus, by mixing $\pi_{\rm sw}$ and $\hat{\pi}$ with weight p, we construct an agent's adversarial action policy whose value is exactly $\overline{V}_A^{\Delta r^*}(\hat{s},h=0) - \epsilon$. Moreover, the mixing weight on $\pi_{\rm sw}$ must satisfy $p \geq x^*$. The derivation of this lower bound on p is as follows:

First, by the definition of the mixed policy and dual variable $\alpha^* > 0$,

$$pV_{A}^{\pi_{\text{sw}},\Delta r^{*}}(\hat{s},h=0) + (1-p)V_{A}^{\hat{\pi},\Delta r^{*}}(\hat{s},h=0) = \overline{V}_{A}^{\Delta r^{*}}(\hat{s},h=0) - \epsilon.$$

Since $V_A^{\pi_{\rm sw},\Delta r^*}(\hat{s},h=0)=\overline{V}_A^{\Delta r^*}(\hat{s},h=0)$, this equality can be rewritten as

$$p\epsilon + (1-p)\Big(V_A^{\hat{\pi},\Delta r^*}(\hat{s},h=0) - \left(\overline{V}_A^{\Delta r^*}(\hat{s},h=0) - \epsilon\right)\Big) = 0.$$

Next, using the inequality $\overline{V}_A^{\Delta r^*}(\hat{s},h=0) - V_A^{\hat{\pi},\Delta r^*}(\hat{s},h=0) \geq (1-x^*)^{-1}\epsilon$, we obtain

$$p\epsilon + (1-p)\left(\frac{-\epsilon}{1-x^*} + \epsilon\right) \ge 0.$$

Finally, dividing both sides by $\epsilon > 0$ gives $p \geq x^*$.

Action Policy under $\Delta r_{\rm sw}$ When the optimal subsidy scheme is shifted from Δr^* to $\Delta r_{\rm sw}$, our primary objective—towards establishing the latter part of the proposition—is to verify that both $\hat{\pi}$ and $\pi_{\rm sw}$ continue to satisfy the requirements of candidate policies. We first claim that the optimal solution x^* and the optimal dual variable α^* remain unchanged under this modification. By the first part of Proposition 4.2, the principal's optimal value is preserved in this transition. Furthermore, recalling from Theorem 4.1 that the objective function F(x) is concave, it follows that both x^* and α^* remain optimal.

Then, to determine whether an action policy $\overline{\pi}$, such as $\hat{\pi}$ and $\pi_{\rm sw}$, qualifies as a candidate policy under $\Delta r_{\rm sw}$, it suffices to verify whether

$$V_{P}^{\overline{\pi}, \Delta r_{sw}}(\hat{s}, h = 0) - \alpha^{*} V_{A}^{\overline{\pi}, \Delta r_{sw}}(\hat{s}, h = 0)$$

$$= \min_{\pi} \left\{ V_{P}^{\overline{\pi}, \Delta r_{sw}}(\hat{s}, h = 0) - \alpha^{*} V_{A}^{\pi, \Delta r_{sw}}(\hat{s}, h = 0) \right\}.$$
(A.9)

For clarity, we define

$$g(\pi; \Delta r) \triangleq V_P^{\pi, \Delta r}(\hat{s}, h = 0) - \alpha^* V_A^{\pi, \Delta r}(\hat{s}, h = 0).$$

The proof proceeds in two steps. First, we establish a lower bound for the right-hand side of (A.9) as

$$\min_{\pi} g(\pi; \Delta r_{\text{sw}}) = \min_{\pi} \left\{ V_P^{\pi, \Delta r_{\text{sw}}}(\hat{s}, h = 0) - \alpha^* V_A^{\pi, \Delta r_{\text{sw}}}(\hat{s}, h = 0) \right\} \geq -\frac{1}{1 - x^*} V_{x^*}^*(\hat{s}, h = 0).$$

Second, we show that for the action policies $\hat{\pi}$ and π_{sw} ,

$$g(\hat{\pi}; \Delta r_{\text{sw}}) = g(\pi_{\text{sw}}; \Delta r_{\text{sw}}) = -\frac{1}{1 - x^*} V_{x^*}^*(\hat{s}, h = 0).$$

In the first step, note that since $\Delta r^* \geq \Delta r_{\rm sw}$, we have for any action policy π ,

$$V_P^{\pi,\Delta r^*}(\hat{s},h=0) \leq V_P^{\pi,\Delta r_{\text{sw}}}(\hat{s},h=0), \quad V_A^{\pi,\Delta r^*}(\hat{s},h=0) \geq V_A^{\pi,\Delta r_{\text{sw}}}(\hat{s},h=0).$$

As $\alpha^* \geq 0$, it follows that for any π .

$$g(\pi, \Delta r_{\text{sw}}) = \min_{\pi} \left\{ V_P^{\pi, \Delta r_{\text{sw}}}(\hat{s}, h = 0) - \alpha^* V_A^{\pi, \Delta r_{\text{sw}}}(\hat{s}, h = 0) \right\}$$

$$\geq \min_{\pi} \left\{ V_P^{\pi, \Delta r^*}(\hat{s}, h = 0) - \alpha^* V_A^{\pi, \Delta r^*}(\hat{s}, h = 0) \right\}$$

$$= -\frac{1}{1 - x^*} V_{x^*}^*(\hat{s}, h = 0).$$
(A.10)

In the second step, to prove $\hat{\pi}$ is a candidate policy and attains the minimum $-\frac{1}{1-x^*}V_{x^*}^*(\hat{s},h=0)$, we connect $g(\hat{\pi};\Delta r_{\rm sw})$ and $-\frac{1}{1-x^*}V_{x^*}^*(\hat{s},h=0)$ via $g(\hat{\pi};\Delta r=0)$. Substituting $x^*=\frac{\alpha^*}{1+\alpha^*}$, we obtain

$$\begin{split} g(\hat{\pi};\Delta r = 0) &= V_P^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) - \alpha^* V_A^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) \\ &= V_P^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) - \frac{x^*}{1-x^*} V_A^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) \\ &= -\frac{1}{1-x^*} \left((x^*-1) V_P^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) + x^* V_A^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) \right) \\ &= -\frac{1}{1-x^*} \left(x^* V_{\rm sw}^{\hat{\pi}}(\hat{s},h = 0) - V_P^{\hat{\pi},\Delta r = 0}(\hat{s},h = 0) \right) \\ &= -\frac{1}{1-x^*} \max_{\pi} \left\{ x^* V_{\rm sw}^{\pi}(\hat{s},h = 0) - V_P^{\pi,\Delta r = 0}(\hat{s},h = 0) \right\} \\ &= \min_{\pi} \left\{ V_P^{\pi,\Delta r = 0}(\hat{s},h = 0) - \frac{x^*}{1-x^*} V_A^{\pi,\Delta r = 0}(\hat{s},h = 0) \right\} \\ &= -\frac{1}{1-x^*} V_{x^*}^*(\hat{s},h = 0). \end{split} \tag{A.12}$$

Equality (A.11) follows since

$$\hat{\pi} = \arg\max_{\pi} \{ x^* V_{\text{sw}}^{\pi}(\hat{s}, h = 0) - V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0) \},$$

while (A.12) corresponds to the definition of $V_x^*(\hat{s}, h = 0)$ in Theorem 4.1:

$$V_x^*(\hat{s},h=0) \triangleq -(1-x) \cdot \min_{\pi} \left\{ V_P^{\pi,\Delta r=0}(\hat{s},h=0) - \frac{x}{1-x} V_A^{\pi,\Delta r=0}(\hat{s},h=0) \right\}.$$

We next turn to discuss the equality between $g(\hat{\pi}, \Delta r_{\rm sw})$ and $g(\hat{\pi}, \Delta r = 0)$. Combining (A.10) with $g(\hat{\pi}, \Delta r = 0) = -\frac{1}{1-x^*}V_{x^*}^*(\hat{s}, h = 0)$, we obtain

$$g(\hat{\pi}, \Delta r_{\rm sw}) \ge g(\hat{\pi}, \Delta r = 0)$$

To establish equality, it suffices to show $g(\hat{\pi}, \Delta r_{\rm sw}) \leq g(\hat{\pi}, \Delta r = 0)$. Since $\Delta r_{\rm sw}(s, a, h) \geq \Delta r(s, a, h)$ for all (s, a, h), we have

$$V_P^{\hat{\pi}, \Delta r_{\rm sw}}(\hat{s}, h = 0) \leq V_P^{\hat{\pi}, \Delta r = 0}(\hat{s}, h = 0), \quad V_A^{\hat{\pi}, \Delta r_{\rm sw}}(\hat{s}, h = 0) \geq V_A^{\hat{\pi}, \Delta r = 0}(\hat{s}, h = 0).$$

With $\alpha^* > 0$, it follows that

$$g(\hat{\pi}, \Delta r_{\text{sw}}) = V_P^{\hat{\pi}, \Delta r_{\text{sw}}}(\hat{s}, h = 0) - \alpha^* V_A^{\hat{\pi}, \Delta r_{\text{sw}}}(\hat{s}, h = 0)$$

$$\leq V_P^{\hat{\pi}, \Delta r = 0}(\hat{s}, h = 0) - \alpha^* V_A^{\hat{\pi}, \Delta r = 0}(\hat{s}, h = 0)$$

$$= g(\hat{\pi}, \Delta r = 0).$$

Thus $g(\hat{\pi}, \Delta r_{\text{sw}}) = g(\hat{\pi}, \Delta r = 0)$.

For $\pi_{\rm sw}$, recall that after changing from Δr^* to $\Delta r_{\rm sw}$, its agent value remains unchanged. Since the social welfare is unaffected by the subsidy scheme, the principal's value also remains the same. Hence,

$$g(\pi_{\text{sw}}, \Delta r_{\text{sw}}) = V_P^{\pi_{\text{sw}}, \Delta r_{\text{sw}}}(\hat{s}, h = 0) - \alpha^* V_A^{\pi_{\text{sw}}, \Delta r_{\text{sw}}}(\hat{s}, h = 0)$$
(A.13)

$$=V_{P}^{\pi_{\rm sw},\Delta r^{*}}(\hat{s},h=0)-\alpha^{*}V_{A}^{\pi_{\rm sw},\Delta r^{*}}(\hat{s},h=0) \tag{A.14}$$

$$= -\frac{1}{1 - x^*} V_{x^*}^*(\hat{s}, h = 0), \tag{A.15}$$

where the last equality follows from (A.6).

Finally, to bound the probability weight, recall that both the maximum agent value and the agent value of π_{sw} remain unchanged, i.e.

$$V_A^{\pi_{\rm sw}, \Delta r_{\rm sw}}(\hat{s}, h=0) = V_A^{\pi_{\rm sw}, \Delta r^*}(\hat{s}, h=0) = \overline{V}_A^{\Delta r_{\rm sw}}(\hat{s}, h=0) = \overline{V}_A^{\Delta r^*}(\hat{s}, h=0).$$

Meanhile, for $\hat{\pi}$, its deviation from the maximum agent value is still bounded as

$$\begin{split} V_A^{\hat{\pi},\Delta r_{\text{sw}}}(\hat{s},h=0) &\leq V_A^{\hat{\pi},\Delta r^*}(\hat{s},h=0) \\ &\leq \overline{V}_A^{\Delta r^*}(\hat{s},h=0) - \frac{\epsilon}{1-x^*} \\ &= \overline{V}_A^{\Delta r_{\text{sw}}}(\hat{s},h=0) - \frac{\epsilon}{1-x^*}. \end{split}$$

Therefore, by the same reasoning as under Δr^* , we conclude that $\pi_{\rm sw}$ and $\hat{\pi}$ constitute an optimal adversarial pair for the agent, with $\pi_{\rm sw}$ assigned a probability weight of at least x^* .

A.6 PROOF OF PROPOSITION 4.3

According to the proof Theorem 4.1, recall that the optimal subsidy scheme Δr^* given optimal dual variable (α^*, V^*) is

$$\Delta r^*(s, a, h) = \frac{1}{1 + \alpha^*} \Big(-V^*(s, h) + \sum_{s' \in \mathcal{S}} P(s'|s, a, h) V^*(s', h+1) + r_P(s, a, h) - \alpha r_A(s, a, h) \Big).$$

Therefore, the agent's action policy $\pi_{\Delta r^*}$ under Δr^* satisfies

$$V_P^{\pi_{\Delta r^*}, \Delta r^*}(\hat{s}, h = 0) = \frac{1}{1 + \alpha^*} V^*(\hat{s}, h = 0) + \frac{\alpha^*}{1 + \alpha^*} V_{\text{sw}}^*(\hat{s}, h = 0) - \alpha^* \epsilon,$$

$$V_A^{\pi_{\Delta r^*}, \Delta r^*}(\hat{s}, h = 0) \ge \overline{V}_A^{\Delta r^*}(\hat{s}, h = 0) - \epsilon = \frac{1}{1 + \alpha^*} (V_{\text{sw}}^*(\hat{s}, h = 0) - V^*(\hat{s}, h = 0)) - \epsilon.$$

Based on the KKT condition, when $x^* \in (0,1)$ hence $\alpha^* > 0$, the above inequality becomes the equality. Therefore, combining the above two equations, we can conclude that

$$\delta_{\text{sw}} = V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{\Delta_r^*}}(\hat{s}, h = 0) = (1 + \alpha^*)\epsilon = \frac{\epsilon}{1 - r^*}.$$

The key to analyze such gap is to examine the relationship between ϵ and x^* . By applying the envelope theorem (Milgrom & Segal, 2002), the derivative of the objective function respect to x in Theorem 4.1 is

$$F'(x) = V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_x}(\hat{s}, h = 0) - \frac{\epsilon}{(1 - x)^2},$$
(A.16)

where $\pi_x = \arg\max_{\pi} \{xV_{\rm sw}^{\pi}(\hat{s},h=0) - V_P^{\pi,\Delta r=0}(\hat{s},h=0)\}$. In general, the objective function contains a finite number of non-differentiable points, which arise from the non-uniqueness of $V_{\rm sw}^{\pi_x}(\hat{s},h=0)$. Nevertheless, since the set of sub-differential can be fully characterized by the expression of the derivative, for the sake of simplicity and clarity we do not distinguish between derivatives and sub-derivatives, and in the discussion of stationary points we directly set the derivative F'(x) to 0.

Meanwhile, since the objective function is concave, the optimal solution x^* and the corresponding policy π_{x^*} can be characterized by the condition that the derivative vanishes. Moreover, the requirement $\alpha^* > 0$ ensures that $x^* \in (0,1)$, which implies that x^* , as an interior optimum, necessarily exists as a stationary point. Consequently, imposing the condition F'(x) = 0 yields

$$1 - x^* = \sqrt{\frac{\epsilon}{V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{x^*}}(\hat{s}, h = 0)}}.$$
 (A.17)

Since π_{x^*} is an action policy, we can immediately have a trivial bound as

$$V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{x^*}}(\hat{s}, h = 0) \ge V_{\text{sw}}^*(\hat{s}, h = 0) - \min_{\pi} V_{\text{sw}}^{\pi}(\hat{s}, h = 0),$$

However, in fact, we can obtain a tighter constant bound on $V_{\rm sw}^*(\hat{s},h=0) - V_{\rm sw}^{\pi_{x^*}}(\hat{s},h=0)$. Denote π_A as the action policy that attains the minimum principal value in the absence of a subsidy, i.e.,

$$\pi_A = \arg\min_{\pi} V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0),$$

and if multiple action policies achieve this minimum, we select π_A to be the one among them that maximizes social welfare. We claim that

$$V_{\text{cw}}^*(\hat{s}, h = 0) - V_{\text{cw}}^{\pi_{x^*}}(\hat{s}, h = 0) > V_{\text{cw}}^*(\hat{s}, h = 0) - V_{\text{cw}}^{\pi_{A}}(\hat{s}, h = 0).$$

To prove this, we show that for any $x \in [0, 1)$,

$$V_{\text{sw}}^{\pi_x}(\hat{s}, h = 0) \ge V_{\text{sw}}^{\pi_A}(\hat{s}, h = 0).$$

Recall that

$$\pi_x = \arg\max_{\pi} \{xV_{\text{sw}}^{\pi}(\hat{s}, h = 0) - V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0)\}.$$

By simple rearrangement and using x>0, we obtain

$$V_{\text{sw}}^{\pi_x}(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_A}(\hat{s}, h = 0) \ge \frac{1}{r} \left(V_P^{\pi_x, \Delta r = 0}(\hat{s}, h = 0) - V_P^{\pi_A, \Delta r = 0}(\hat{s}, h = 0) \right)$$

As π_A minimizes the principal value with $\Delta r = 0$, the right-hand side is nonnegative, which immediately implies

$$V_{\text{sw}}^{\pi_x}(\hat{s}, h = 0) > V_{\text{sw}}^{\pi_A}(\hat{s}, h = 0),$$

as desired. Moreover, this bound is tight: based on the definition of π_A , there always exists some $x_0 \in (0,1)$ such that for any $x \in [0,x_0]$, the policy π_A maximizes $xV_{\rm sw}^\pi(s,h) - V_P^{\pi,\Delta r=0}(s,h)$.

Having established a bound on $V_{\rm sw}^*(\hat{s},h=0)-V_{\rm sw}^{\pi_{x^*}}(\hat{s},h=0)$, we can immediately derive a corresponding bound on $1-x^*$:

$$1 - x^* \le \sqrt{\frac{\epsilon}{V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_A}(\hat{s}, h = 0)}}.$$

Finally, applying the condition F'(x) = 0 once more, we obtain

$$\delta_{\text{sw}} = \frac{\epsilon}{1 - x^*} = (V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_A}(\hat{s}, h = 0)) (1 - x^*)$$
$$\leq \sqrt{(V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_A}(\hat{s}, h = 0)) \cdot \epsilon},$$

demonstrating that the social welfare gap $\delta_{\rm sw}$ is upper bounded by $O(\sqrt{\epsilon})$.

A.6.1 THE DEPENDENCE OF δ_{SW} on ϵ

When analyzing the relationship between $\delta_{\rm sw}$ and ϵ , a straightforward observation is that the effectiveness of the limited subsidy diminishes as ϵ approaches infinity, allowing the agent to bypass the globally ϵ -IC constraint and achieve a trivial minimization of the leader's value. In this scenario, even as ϵ continues to increase, the social welfare gap remains unchanged. Therefore, our analysis focuses on relatively small values of ϵ , examining how the social welfare gap grows $\delta_{\rm sw}$ as ϵ increases and $x^* \in (0,1)$.

The key issue in equation (A.17) is that changes in ϵ may simultaneously affect both x^* and the term $V^*_{\rm sw}(\hat s,h=0)-V^{\pi_{x^*}}_{\rm sw}(\hat s,h=0)$. However, a closer analysis of the derivative reveals that ϵ cannot influence these two quantities simultaneously. Before presenting the detailed argument, we first establish the following lemma, which shows that x^* is monotone in ϵ .

Lemma A.1. The optimal solution $x^* \in (0,1)$ is monotonically non-decreasing in ϵ .

Proof. Recall that

$$V_x^*(\hat{s}, h = 0) \triangleq \max_{\pi} \{xV_{\text{sw}}^{\pi}(\hat{s}, h = 0) - V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0)\}.$$

Since $V_x^*(\hat{s},h=0)$ is the maximum of finitely many linear functions in x, its derivative $V_{\rm sw}^{\pi_x}(\hat{s},h=0)$ is non-decreasing in x. Consequently, $V_{\rm sw}^*(\hat{s},h=0)-V_{\rm sw}^{\pi_{x^*}}(\hat{s},h=0)$ is non-increasing in x. Further, for $x^*\in(0,1)$, imposing the stationarity condition F'(x)=0 yields

$$\epsilon = (1 - x^*)^2 (V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{x^*}}(\hat{s}, h = 0)).$$

As x^* increases, both $(1-x^*)$ and $V^*_{\rm sw}(\hat s,h=0)-V^{\pi_{x^*}}_{\rm sw}(\hat s,h=0)$ decrease, making the right-hand side of the above equation monotonically non-increasing in x^* . Thus, for the equality to hold, an increase in ϵ must be matched by an increase in x^* , which proves the claim.

We then analyze the two possible cases of x^* given ϵ , corresponding to $\delta_{\rm sw}$ scaling as $O(\sqrt{\epsilon})$ or $O(\epsilon)$, respectively:

- When the piecewise linear function $V_x^*(s,h)$ is differentiable at x^* , every $\pi_{x^*} \in \Pi_{x^*}$ yields the same social welfare $V_{\mathrm{sw}}^{\pi_{x^*}}(\hat{s},h=0)$. Owing to this uniqueness, there exists a small interval δ such that for all $x \in (x^* \delta, x^* + \delta)$, the policy π_{x^*} remains unchanged. Consequently, x^* is linearly related to $\sqrt{\epsilon}$, and the social welfare gap scales as $O(\sqrt{\epsilon})$.
- When $V_x^*(s,h)$ is non-differentiable at x^* , different π_{x^*} induce different levels of social welfare. In this case as F(x) is concave, within Π_{x^*} , there exist two policies $\pi_{x^*}^l$ and $\pi_{x^*}^r$ achieving the minimum and maximum social welfare, respectively, which define the left-and right-hand derivatives around x^* . We claim as long as

$$V_{\text{sw}}^*(\hat{s},h=0) - V_{\text{sw}}^{\pi_x^l}(\hat{s},h=0) \geq \frac{\epsilon'}{(1-x^*)^2} \geq V_{\text{sw}}^*(\hat{s},h=0) - V_{\text{sw}}^{\pi_x^r}(\hat{s},h=0),$$

the optimal solution for ϵ' remains x^* , and the social welfare gap scales as $O(\epsilon)$. To see this, consider the counterexample. Recall from the proof of Lemma A.1 that $V_{\rm sw}^{\pi_{x^*}}(\hat{s},h=0)$ is

non-decreasing in x. Suppose the solution for ϵ' is x'. If $x' < x^*$, then

$$F'(x') = V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{x'}}(\hat{s}, h = 0) - \frac{\epsilon}{(1 - x')^2}$$

$$\geq V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{x^*}^l}(\hat{s}, h = 0) - \frac{\epsilon}{(1 - x')^2}$$

$$> V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_{x^*}^l}(\hat{s}, h = 0) - \frac{\epsilon}{(1 - x^*)^2}$$

$$> 0.$$

which shows x' is not a stationary point, and hence not optimal for a concave function. The proof for $x' > x^*$ follows analogously by replacing $\pi^l_{x^*}$ with $\pi^r_{x^*}$ and showing F'(x') < 0.

A.7 Proof of Theorem 5.1

A.7.1 TECHNICAL LEMMAS

As our proof involves iteratively optimizing the subsidy scheme for the single-period problem instance, we begin by introducing the following definitions for clarity.

Definition A.1. Given tolerance ϵ , the problem instance $I = ((r_P(a_i), r_A(a_i))_i)$ is a single-period problem where H = 1, the agent is a globally ϵ -IC agent, and r_P, r_A are the principal and agent reward functions for actions $a_1, \dots, a_{|A|}$. Given ϵ and problem instance I:

- $V_{P}^{\Delta r}(I)$ denotes the **principal value** under subsidy scheme Δr .
- $V_A^{\Delta r}(I)$ denotes the maximal agent value under subsidy scheme Δr .
- $V_P^*(I)$ denotes the **optimal principal value** under the optimal subsidy scheme Δr^* .
- $V_A^*(I)$ denotes the **maximal agent value** under the optimal subsidy scheme Δr^* .

Based on the above definitions, we establish several useful properties of the subsidy scheme in the following lemmas, which will be employed in the proof of NP-hardness. Intuitively, the first lemma describes how the principal's value is determined when the agent adversarially reallocates probabilities in response to a given reward transfer. The second lemma characterizes the optimal reward transfer and the corresponding principal and agent values in a simple two-action instance. The third lemma analyzes how the optimal principal value in a three-action instance relates to the optimal values of its two-action sub-instances, providing useful bounds for iterative constructions.

Lemma A.2. Let I = ((a,b),(c,d)) with a > c and $b \ge d$. Under the subsidy scheme $\Delta r = 0$, the principal's final value is

$$V_P^{\Delta r=0}(I) = \frac{ab - a\epsilon - ad + c\epsilon}{b - d}.$$

Proof. To adversarially minimize the principal value under the constraint of global ϵ -IC, it's obvious that the agent will choose to mix the first and the second action so that $V_A^{\Delta r=0}(I)=b-\epsilon$. By denoting the probability weight on the first action as p, we have

$$pb + (1-p)d = b - \epsilon \implies p = \frac{b - d - \epsilon}{b - d}.$$

Substituting into the principal's value formula,

$$V_P^{\Delta r=0}(I) = pa + (1-p)c = \frac{ab - a\epsilon - ad + c\epsilon}{b - d}.$$

Lemma A.3. Let I = ((a,0),(0,0)) with $a > \epsilon$. Then the principal's optimal value is

$$V_{P}^{*}(I) = (\sqrt{a} - \sqrt{\epsilon})^{2}$$
.

achieved by setting the reward transfer for the first action as $\sqrt{a\epsilon}$. Correspondingly, the agent's value under this optimal transfer is

$$V_A^*(I) = \sqrt{a\epsilon}$$
.

Proof. Clearly, the principal will allocate a positive subsidy to the first action only if doing so can yield a principal value exceeding zero. Consider a reward transfer $x \ge \epsilon$ assigned to the first action under a subsidy scheme Δr . For a given x, the principal's value is

$$V_P^{\Delta r}(I) = a - x - \frac{a\epsilon}{r}.$$

Maximizing over $x \geq \epsilon$, we obtain the optimal principal value

$$V_P^*(I) = \max_{x > \epsilon} V_P^{\Delta r}(I) = a + \epsilon - 2\sqrt{a\epsilon} = (\sqrt{a} - \sqrt{\epsilon})^2.$$

Consequently, the agent's value under the optimal agent value is

$$V_{\Lambda}^{*}(I) = \sqrt{a\epsilon}$$
.

Lemma A.4. Let $I = (A_1, A_2, A_3)$ with $A_i = (r_P(a_i), r_A(a_i))$. Suppose for $i \in \{2, 3\}$, $r_P(a_1) + r_A(a_1) > r_P(a_i) + r_A(a_i) + \epsilon$. Define $I' = (A_1, A_2)$ and $I'' = (A_1, A_3)$. Then $V_P^*(I) < \min\{V_P^*(I'), V_P^*(I'')\}$.

Proof. According to Proposition 4.2, for any instance, the optimal reward transfer assigns a nonzero reward only to the first action. Applying concave maximization to instance I, we recall that the objective function is

$$F(x) = xV_{\text{sw}}^*(\hat{s}, h = 0) - V_x^*(\hat{s}, h = 0) - \frac{x}{1 - x}\epsilon,$$

where $\pi_x = \arg\max_{\pi} \{xV_{\text{sw}}^{\pi}(\hat{s}, h = 0) - V_P^{\pi, \Delta r = 0}(\hat{s}, h = 0)\}$, and its derivative is

$$F'(x) = V_{\text{sw}}^*(\hat{s}, h = 0) - V_{\text{sw}}^{\pi_x}(\hat{s}, h = 0) - \frac{\epsilon}{(1 - x)^2}.$$

By observation, we have

$$F'(0) = V_{sw}^*(\hat{s}, h = 0) - V_{sw}^{\pi_0}(\hat{s}, h = 0) - \epsilon,$$

where π_0 minimizes the principal's value, and $F'(x) \to -\infty$ as $x \to 1$. Therefore, there are two cases:

- There exists π_0 that chooses a_1 in instance I. Substituting into the derivative, we obtain F'(0) < 0, and the optimal solution is $x^* = 0$, which implies that the optimal subsidy scheme is $\Delta r^* = 0$. Similarly, the optimal schemes for I' and I'' are also $\Delta r^* = 0$. Since action a_1 yields the highest social welfare, it also provides the largest agent value. Consequently, in all three instances, the agent will deterministically select action a_1 , resulting in $V_P^*(I) = V_P^*(I') = V_P^*(I'')$.
- π_0 chooses an action other than a_1 . In this case, based on the condition that for $i \in \{2,3\}$, $r_P(a_1) + r_A(a_1) > r_P(a_i) + r_A(a_i) + \epsilon$, we have F'(0) > 0 and $x^* \in (0,1)$ for instances I, I' and I''. According to Proposition 4.2, the optimal subsidy assigns a positive transfer only to action a_1 , and there exists an agent's adversarial action policy that mixes action a_1 with other actions in all three instances. Suppose the optimal subsidy schemes for instances I, I', I'' are $\Delta r^*(I)$, $\Delta r^*(I')$, and $\Delta r^*(I'')$, respectively. Then we have

$$\begin{split} V_P^*(I) &= \min\{V_P^{\Delta r^*(I)}(A_1,A_2), V_P^{\Delta r^*(I)}(A_1,A_3)\}, \\ V_P^*(I') &= V_P^{\Delta r^*(I')}(A_1,A_2), \\ V_P^*(I'') &= V_P^{\Delta r^*(I'')}(A_1,A_3). \end{split}$$

By the definition of the optimal subsidy scheme, we then obtain

$$V_P^{\Delta r^*(I)}(A_1, A_2) \le V_P^{\Delta r^*(I')}(A_1, A_2), \quad V_P^{\Delta r^*(I)}(A_1, A_3) \le V_P^{\Delta r^*(I'')}(A_1, A_3),$$

which completes the proof.

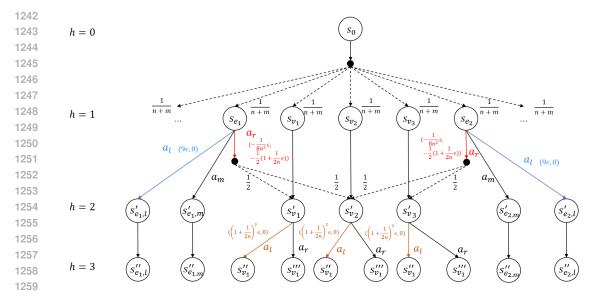


Figure 4: Illustrative Construction of the NP-Hardness Reduction Instance

A.7.2 PROOF OF THEOREM 5.1

We prove the hardness by a reduction from the Maximum Independent Set problem.

Construction Given a graph $G = \langle V, E \rangle$, we construct a corresponding problem instance as illustrated in Figure 4. Let n = |V| and m = |E|, and assume $n \ge 2$.

Throughout the proof, we refer to an action with $r_P=0$ and $r_A=0$ as a blank action. The time horizon is set to H=4, with a global initial state \hat{s} that has only one blank action, which transitions uniformly to the vertex states $s_{v_1}, s_{v_2}, \ldots, s_{v_n}$ and the edge states $s_{e_1}, s_{e_2}, \ldots, s_{e_m}$.

For each vertex state s_v , there is a blank action leading to an intermediate state s_v' for padding. The subsidy scheme at this intermediate state encodes whether the corresponding vertex v is included in the maximum independent set. Two actions, a_l and a_r , are available from s_v' with the following specifications:

- a_l : $r_P(s_v', a_l, h=2) = \left(1 + \frac{1}{2n}\right)^2 \epsilon$, $r_A(s_v', a_l, h=2) = 0$, deterministically transitioning to state s_v'' .
- a_r : a blank action that deterministically transitions to state s_v''' .

For each edge $e \in E$ connecting vertices v_i and v_j , the corresponding edge state s_e has three actions a_l , a_m , and a_r designed to enforce the independent set constraints:

- a_l : $r_P(s_e, a_l, h = 1) = 9\epsilon$, $r_A(s_e, a_l, h = 1) = 0$, deterministically transitioning to state $s'_{e,l}$, which is then followed by a padding state $s''_{e,l}$.
- a_m : a blank action that deterministically transitions to state $s'_{e,m}$, which is then followed by a padding state $s''_{e,m}$.
- a_r : $r_P(s_e,a_r,h=1)=-\frac{1}{8n^2}\epsilon$, $r_A(s_e,a_r,h=1)=-\frac{1}{2}\left(1+\frac{1}{2n}\right)\epsilon$, which transitions with equal probability to states s'_{v_i} and s'_{v_j} .

Under this MDP construction, we claim that there exists an independent set of size k in G if and only if there exists a subsidy scheme Δr in the MDP that allows the principal to achieve a reward of

$$\frac{\frac{k}{4n^2} + 4m}{n+m} \epsilon.$$

If direction Given a size-k independent set $V^* \subset V$ in graph G, we construct a subsidy scheme Δr that achieves a principal value of $\frac{\frac{k}{4n^2} + 4m}{n+m} \epsilon$. The scheme is defined as follows:

- For each $v \in V^*$, set $\Delta r(s_v', a_l, h = 2) = 1 + \frac{1}{2n}\epsilon$ and $\Delta r(s_v', a_r, h = 2) = 0$.
- For each $v \notin V^*$, set zero subsidy for both actions a_l and a_r at $(s'_n, h = 2)$.
- For each edge $e \in E$, set $\Delta r(s_e, a_l, h = 1) = 3\epsilon$ and leave actions a_m and a_r with zero subsidy.
- No subsidy is applied to any other actions.

With this subsidy scheme Δr , the agent greedily minimizes the principal value from bottom to top. First, consider a vertex state s'_v and let π denote the agent's action policy under Δr . By utilizing Lemma A.3, we obtain:

- For $v\in V^*$, $V_P^{\pi_{\Delta r},\Delta r}(s'_v,h=2)=\frac{1}{4n^2}\epsilon$ and $\overline{V}_A^{\Delta r}(s'_v,h=2)=1+\frac{1}{2n}\epsilon$.
- For $v \notin V^*$, $V_P^{\pi_{\Delta r}, \Delta r}(s'_v, h = 2) = 0$ and $\overline{V}_A^{\Delta r}(s'_v, h = 2) = 0$.

Next, consider an edge state s_e , where e connects vertices v_1 and v_2 . There are two scenarios depending on whether one of the endpoints is in the independent set:

- If one endpoint is in V^* (i.e., $v_1 \in V^*$ or $v_2 \in V^*$), the agent faces a single-period problem instance $((6\epsilon, 3\epsilon), (0, 0), (0, 0))$. Since two actions have identical principal and agent rewards, Lemma A.2 implies that the resulting principal value in $(s_\epsilon, h = 1)$ is 4ϵ .
- If neither endpoint is in V^* (i.e., $V_A \notin V^*$ and $v_2 \notin V^*$), the agent faces the instance $((6\epsilon, 3\epsilon), (0, 0), (-\frac{1}{8n^2}\epsilon, -\frac{1}{2}(1+\frac{1}{2n})\epsilon))$. Following the analysis in Lemma A.4, the agent chooses a mixture between (a_1, a_2) or (a_1, a_3) . Applying Lemma A.2, we find with n>1

$$\begin{split} V_P^{\Delta r = 0}((6\epsilon, 3\epsilon), (0, 0)) - V_P^{\Delta r = 0}\left((6\epsilon, 3\epsilon), \left(-\frac{1}{8n^2}\epsilon, -\frac{1}{2}\left(1 + \frac{1}{2n}\right)\epsilon\right)\right) \\ &= 4\epsilon - \frac{15 + \frac{3}{2n} - \frac{1}{8n^2}}{\frac{7}{2} + \frac{1}{4n}}\epsilon \\ &= \frac{(14 + \frac{1}{n}) - (15 + \frac{3}{2n} - \frac{1}{8n^2})}{\frac{7}{2} + \frac{1}{4n}}\epsilon \\ &= \frac{-1 - \frac{1}{2n} + \frac{1}{8n^2}}{\frac{7}{2} + \frac{1}{4n}}\epsilon \\ &< 0. \end{split}$$

which confirms that the final principal value remains 4ϵ .

Therefore, the total principal value under subsidy scheme Δr is

$$V_P^{\pi_{\Delta r}, \Delta r}(\hat{s}, h = 0) = \frac{k \cdot \frac{1}{4n^2} \epsilon + (n - k) \cdot 0 + 4\epsilon \cdot m}{n + m} = \frac{\frac{k}{4n^2} + 4m}{n + m} \epsilon.$$

Only if direction Suppose a subsidy scheme Δr achieves $\frac{\frac{k}{4n^2}+4m}{n+m}\epsilon$. for the principal. We show that this implies the existence of a size-k independent set $V^* \subset V$ in G.

We first upper bound the maximum principal value achievable under any subsidy scheme Δr . There are two primary sources of principal value:

- Vertex states s_v' : by Lemma A.3, each vertex contributes at most $\frac{1}{4n^2}\epsilon$.
- Edge states s_e : by Lemma A.3 and Lemma A.4, each edge contributes at most $V_P^*((9\epsilon,0),(0,0))=4\epsilon$.

Consequently, to attain the claimed principal value, at least k vertex states must yield positive contributions. We claim that these vertices form an independent set. To see this, suppose otherwise: let v_1 and v_2 be connected by an edge \overline{e} . Since both s'_{v_1} and s'_{v_2} have nonzero principal values, the principal must provide a reward transfer of at least 1 on action a_l at both s'_{v_1} and s'_{v_2} . Then, for action a_r at $s_{\overline{e}}$, the agent's expected value is at least $(\frac{1}{2} - \frac{1}{4n})\epsilon$, while the principal's value is at most $\frac{1}{8n^2}\epsilon$. We can upper bound the principal value from $(s_{\overline{e}}, h = 1)$ under any subsidy scheme Δr as

$$\begin{split} V_P^{\pi_{\Delta r},\Delta r}(s_{\overline{e}},h=1) &\leq V_P^* \left((9\epsilon,0), (\frac{1}{8n^2}\epsilon, (\frac{1}{2} - \frac{1}{4n})\epsilon \right) \\ &= \frac{1}{8n^2}\epsilon + V_0^* \left((9\epsilon - \frac{1}{8n^2}\epsilon,0), (0, (\frac{1}{2} - \frac{1}{4n}))\epsilon \right) \\ &\leq \frac{1}{8n^2}\epsilon + V_0^* \left((9\epsilon,0), (0, (\frac{1}{2} - \frac{1}{4n}))\epsilon \right) \\ &= \frac{1}{8n^2}\epsilon + V_0^* ((\frac{17}{2} + \frac{1}{4n})\epsilon,0), (0,0)) \\ &= \frac{1}{8n^2}\epsilon + \left(\sqrt{(\frac{17}{2} + \frac{1}{4n})\epsilon} - \sqrt{\epsilon} \right)^2 \\ &= \left(\frac{1}{8n^2} + (\frac{19}{2} + \frac{1}{4n}) - 2\sqrt{\frac{17}{2} + \frac{1}{4n}} \right)\epsilon \end{split}$$
(A.18)

Inequality (A.18) follows directly from Lemma A.4 together with the observation that, in any single-period problem instance, simultaneously decreasing the principal's reward and increasing the agent's reward for an action can only reduce the optimal principal value. In equality (A.19), It is evident that subtracting the same value from the principal reward of each action and then summing afterwards does not affect the optimal solution. Inequality (A.20) arises from the fact that there is a pure principal reward increase in the first action. In equation (A.21), to obtain a strictly positive principal value, at least $(\frac{1}{2} - \frac{1}{4n})\epsilon$ must be subsidized on the first action. After such subsidy, as both actions now have the same agent reward, we set zero reward for both actions to sustain the relative value. Meanwhile, we can apply lemma A.3 to find the to optimal principal value.

Next, we upper bound the total principal reward across all sources:

- Principal reward from agent visiting s_v : at most $\frac{1}{4n^2}\epsilon$ per vertex, for at most n vertices.
- Principal reward from agent visiting s_e for edges e ∈ E \ {ē}: at most 4ε per edge, for at most m − 1 edges.
- Principal reward from agent visiting $s_{\overline{e}}$: at most $\left(\frac{1}{8n^2} + (\frac{19}{2} + \frac{1}{4n}) 2\sqrt{\frac{17}{2} + \frac{1}{4n}}\right)\epsilon$.

Summing over all contributions, the total principal value is

$$V_P^{\pi_{\Delta r}, \Delta r}(\hat{s}, h = 0) \le \frac{1}{n+m} \left[\left(\frac{1}{8n^2} + \left(\frac{19}{2} + \frac{1}{4n} \right) - 2\sqrt{\frac{17}{2} + \frac{1}{4n}} \right) \epsilon + (m-1) \cdot 4\epsilon + n \cdot \frac{\epsilon}{4n^2} \right].$$

Comparing with the claimed value, for $n \geq 2$ we obtain

$$\frac{\frac{k}{4n^2} + 4m}{n+m} \epsilon - V_P^{\pi_{\Delta r}, \Delta r}(\hat{s}, h = 0)$$

$$= \frac{\epsilon}{n+m} \left(\left(\frac{k}{4n^2} + 4 \right) - \left(\frac{1}{8n^2} + \frac{19}{2} + \frac{1}{4n} - 2\sqrt{\frac{17}{2} + \frac{1}{4n}} \right) - \frac{1}{4n} \right) \tag{A.22}$$

$$\geq \frac{\epsilon}{n+m} \left(2\sqrt{\frac{17}{2}} - \frac{1}{8n^2} - \frac{11}{2} - \frac{1}{2n} \right) \tag{A.23}$$

$$\geq \frac{\epsilon}{n+m} \left(2\sqrt{\frac{17}{2}} - \frac{1}{32} - \frac{11}{2} - \frac{1}{4} \right) \tag{A.24}$$

$$>0$$
 (A.25)

Here, inequality (A.22) follows from neglecting the $\frac{k}{4n^2}$ term, and inequality (A.23) is obtained by substituting n=2. This contradiction demonstrates that any set of k vertices yielding positive principal value must form an independent set, thereby completing the proof.