

Problem set up

- MDP $\mu = (S, A, P, \rho, T) \in \mathcal{M}$
- Utility $\mathcal{U} = \sum_{t=1}^T r_t$

For a fixed MDP $\mu \in \mathcal{M}$, we define the

- Utility $\mathcal{U}(\pi, \mu) = E_{\mu}^{\pi}[\mathcal{U}]$
- Optimal Utility $\mathcal{U}^*(\mu) = \max_{\pi} \mathcal{U}(\pi, \mu)$

For a distribution β over MDPs, we define the

- Utility $\mathcal{U}(\pi, \beta) = E_{\beta}^{\pi}[\mathcal{U}] = \int_{\mathcal{M}} \mathcal{U}(\pi, \mu) d\beta(\mu)$
- Bayes-optimal utility $\mathcal{U}^*(\beta) = \sup_{\pi} \mathcal{U}(\pi, \beta)$

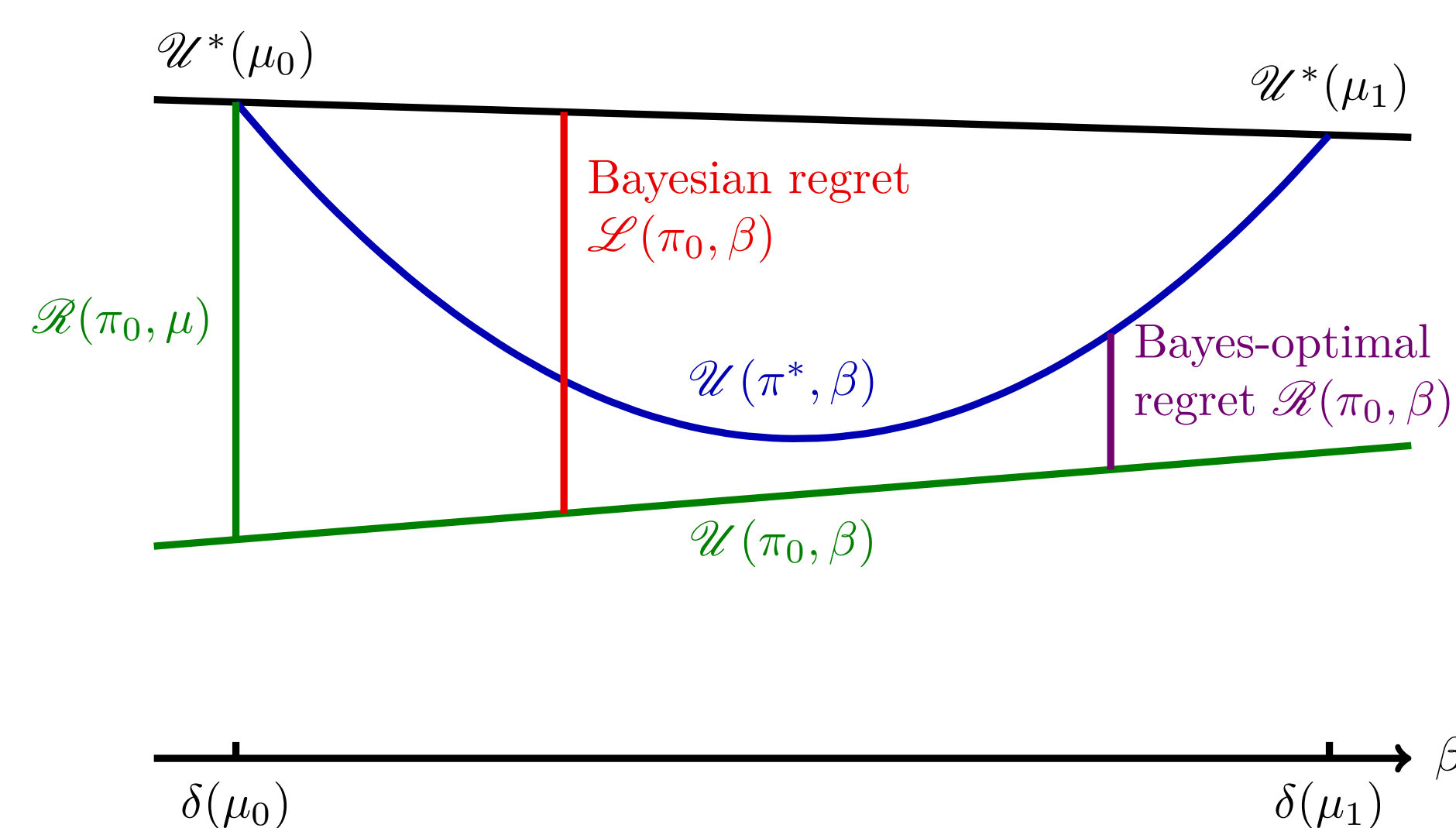
Interpretation of β :

- The agent's subjective belief about which MDP is the most likely a priori.
- The MDP is actually drawn randomly from distribution.

Suppose Nature chooses β arbitrarily or adversarially, then we are interested in finding

$$\max_{\pi} \min_{\beta} \mathcal{U}(\pi, \beta)$$

However, for an unrestricted set of priors, Nature could pick a prior such that all rewards are zero, thus trivially achieving minimal utility. Instead we consider regret.



Notions of regret

For a fixed MDP $\mu \in \mathcal{M}$, we define

$$\mathcal{R}(\pi, \mu) = \mathcal{U}^*(\mu) - \mathcal{U}(\pi, \mu)$$

For a prior β , we define the

- Bayes-optimal regret

$$\mathcal{R}(\pi, \beta) = \mathcal{U}^*(\beta) - \mathcal{U}(\pi, \beta)$$

- Bayesian regret (comparing against an oracle)

$$\mathcal{L}(\pi, \beta) = E_{\mu \sim \beta}[\mathcal{R}(\pi, \mu)] = \int_{\mathcal{M}} \mathcal{U}^*(\mu) - \mathcal{U}(\pi, \mu) d\beta(\mu)$$

Of course, we always have $\mathcal{R}(\pi, \beta) \leq \mathcal{L}(\pi, \beta)$.

Minimax game against Nature

We define minimax games with respect to the Bayes-optimal regret $\min_{\pi} \max_{\beta} \mathcal{R}(\pi, \beta)$ and Bayesian regret $\min_{\pi} \max_{\beta} \mathcal{L}(\pi, \beta)$.

Corollary (value of the game)

The minimax game with respect to the utility and the Bayesian regret have a value, i.e. it holds that

$$\max_{\pi} \min_{\beta} \mathcal{U}(\pi, \beta) = \min_{\beta} \max_{\pi} \mathcal{U}(\pi, \beta), \quad \min_{\pi} \max_{\beta} \mathcal{L}(\pi, \beta) = \max_{\beta} \min_{\pi} \mathcal{L}(\pi, \beta).$$

Lemma

The minimax game with respect to the Bayes-optimal regret may not have a value, i.e.,

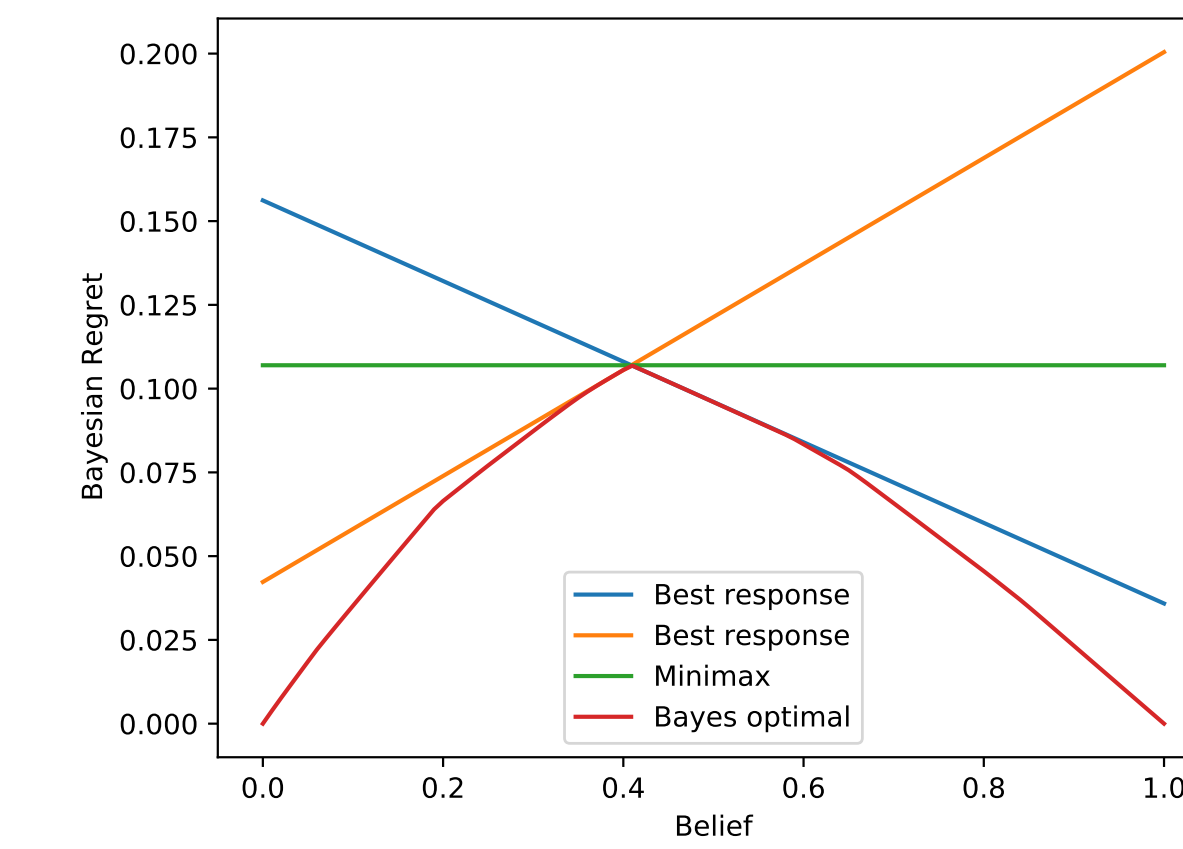
$$\min_{\pi} \max_{\beta} \mathcal{R}(\pi, \beta) < \max_{\beta} \min_{\pi} \mathcal{R}(\pi, \beta).$$

Lemma (Bayesian regret of the Bayes-optimal policy)

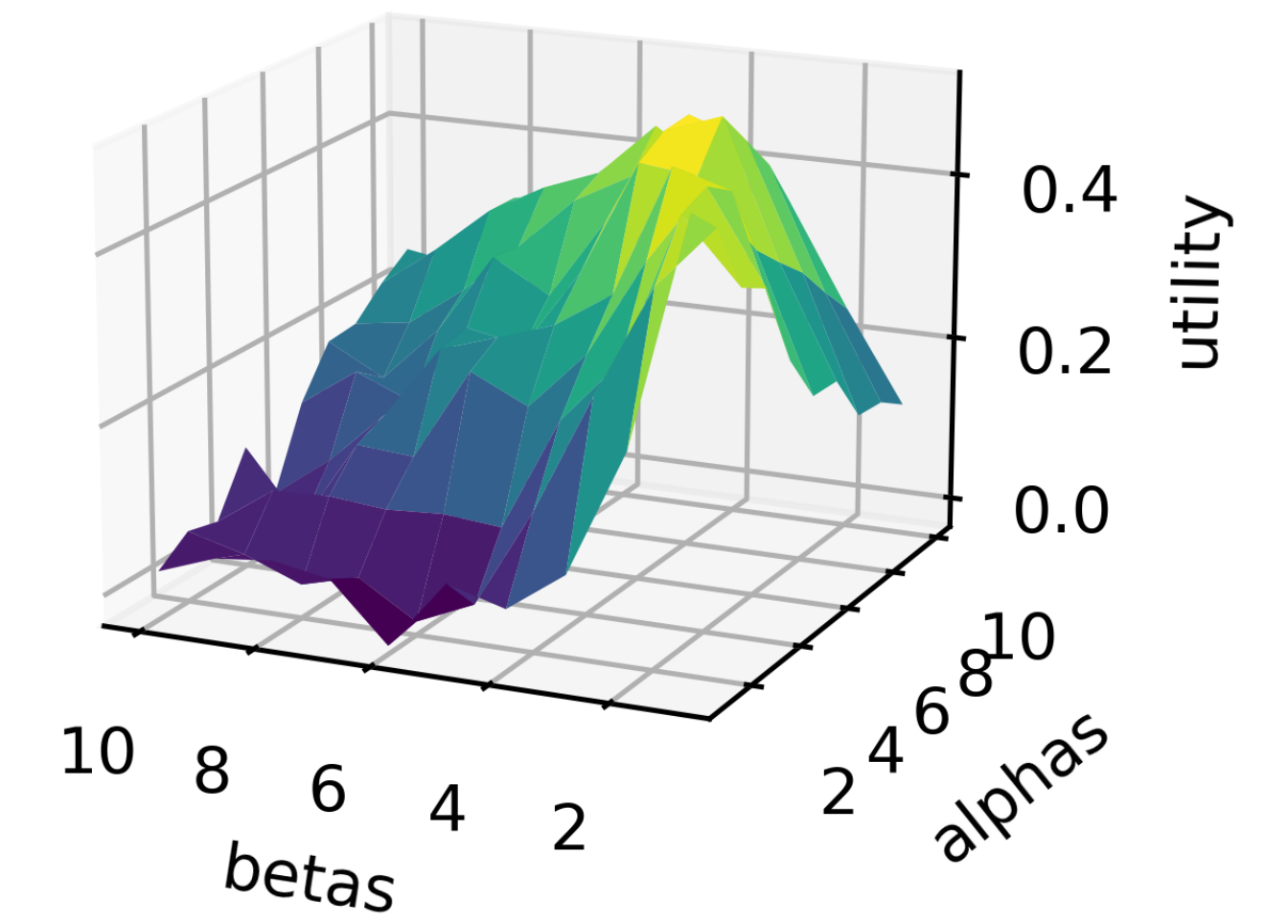
The worst-case Bayesian regret of the Bayes-optimal policy equals the minimax Bayesian regret, i.e.

$$\max_{\beta} \mathcal{L}(\pi^*(\beta), \beta) = \min_{\pi} \max_{\beta} \mathcal{L}(\pi, \beta)$$

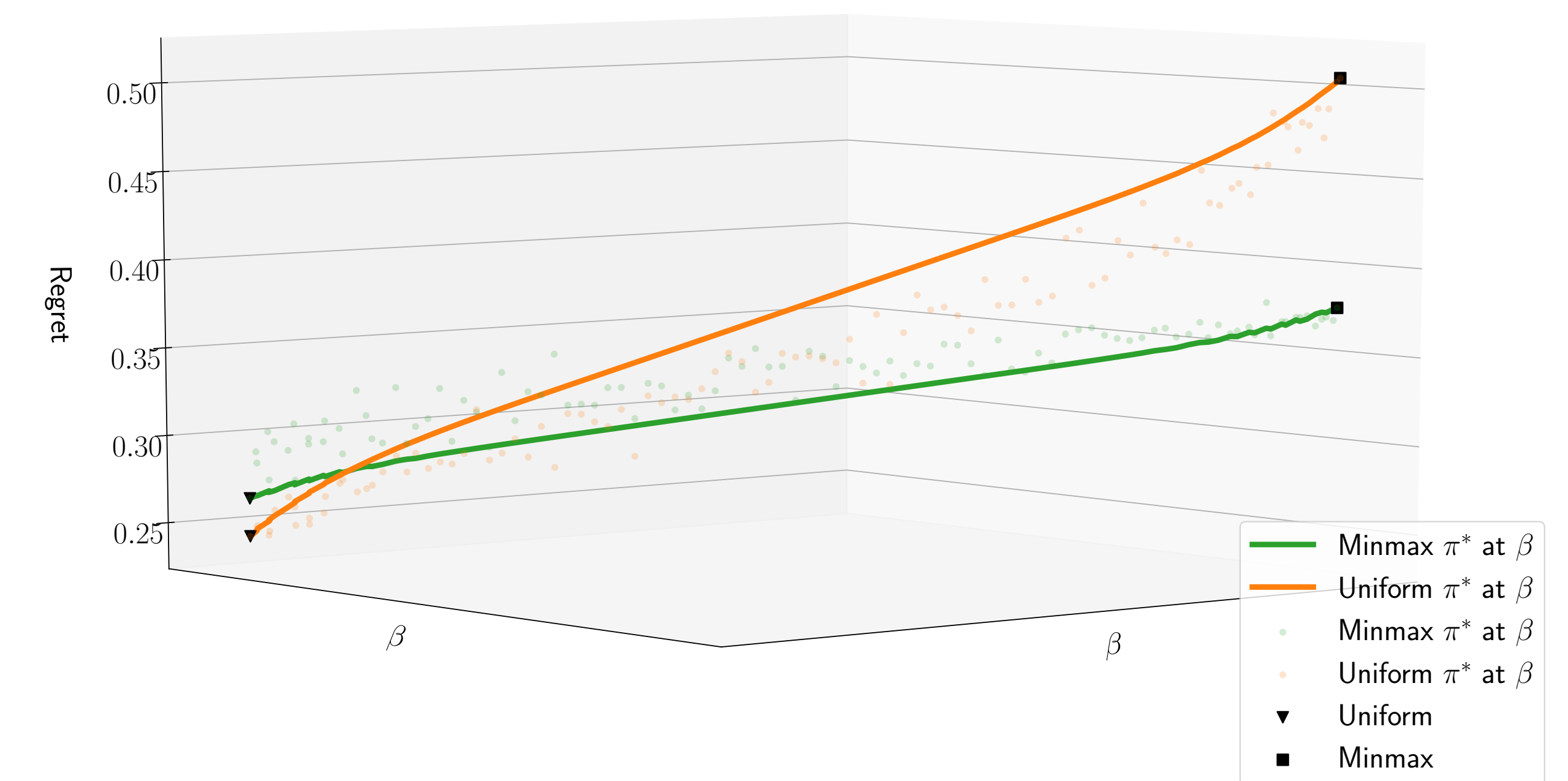
Experiments



(a) Regret for a two MDP task for a variety of policies.



(b) Bayesian regret of the Bayes-optimal policy in two armed Bernoulli bandit task. The first arms prior is fixed to $\text{Beta}(4, 2)$ while the other is given by the values on the x- and y-axis.



(c) t-SNE embeddings for (approximately) minimax and uniform beliefs with their corresponding Bayesian regret \mathcal{L} .