

PROBLEM

- Compute the **Maximum Expected Value (MEV)** of a set of two or more independent random variables $X = \{X_1, \dots, X_M\}$ given samples $S = \{S_1, \dots, S_N\}$
- Most RL algorithms need to **approximate MEV**
- A good estimation of MEV is **critical** in many real-world applications

CONTRIBUTIONS

1. We propose the **Weighted Estimator (WE)** to approximate MEV
2. The approximation is done by a **weighted average** of sample means
3. We provide theoretical and empirical **comparisons** of the **performance** of WE and other MEV estimators

MAXIMUM EXPECTED VALUE ESTIMATION

IDEA

- Compute the *Maximum Expected Value*

$$\mu_*(X) = \max_i \mu_i = \max_i \int_{-\infty}^{+\infty} x f_i(x) dx$$

- of independent random variables ($X = \{X_1, X_2, \dots, X_M\}$) whose PDFs f_i are unknown

- **Issue:** μ_* cannot be computed analytically

- Given a set of noisy samples $S = \{S_1, \dots, S_N\}$ retrieved by the unknown distributions of each X_i

GOAL

$$\mu_*(X) \approx \hat{\mu}_*(S)$$

NAÏVE APPROACH

- *Maximum Estimator (ME)*
- i.e., take the maximum of the sample means

$$\hat{\mu}^{ME}(S) = \max_i \hat{\mu}_i(S) \approx \mu_*(X)$$

- **Positive** bias can cause problems in some applications (e.g., Q-Learning)

DOUBLE ESTIMATOR (DE) [Van Hasselt, 2010]

1. Split dataset S in two disjoint sets

$$S^A = \{S_1^A, \dots, S_N^A\} \quad \text{and} \quad S^B = \{S_1^B, \dots, S_N^B\}$$

2. Estimate the maximum index in each set

$$a^* = \arg \max_i \hat{\mu}_i^{ME}(S^A) \quad \text{and} \quad b^* = \arg \max_i \hat{\mu}_i^{ME}(S^B)$$

3. Take the average maximum value

$$\hat{\mu}^{DE}(S) = \frac{\hat{\mu}_{b^*}^{ME}(S^A) + \hat{\mu}_{a^*}^{ME}(S^B)}{2} \approx \mu_*(X)$$

Negative bias may solve ME issues in many applications

WEIGHTED ESTIMATOR (WE)

$$\hat{\mu}^{WE}(S) = \sum_{i=1}^M \hat{\mu}_i(S) w_i^S$$

WE

Weights the sample means by the probability of being the maximum

$$w_i^S = P\left(\hat{\mu}_i(S) = \max_j \hat{\mu}_j(S)\right) = \int_{-\infty}^{+\infty} \underbrace{\hat{f}_i^S(x)}_{\text{PDF unknown}} \prod_{j \neq i} \underbrace{\hat{F}_j^S(x)}_{\text{CDF unknown}} dx$$

- As the number of samples increases, by the **central limit theorem**:

$$\hat{\mu}_i(S) \sim \mathcal{N}\left(\underbrace{\mu_i}_{\text{sample mean}}, \underbrace{\frac{\sigma_i^2}{|S_i|}}_{\text{sample variance}}\right) \quad \text{i.e.,} \quad \hat{f}_i^S = \text{normal distribution}$$

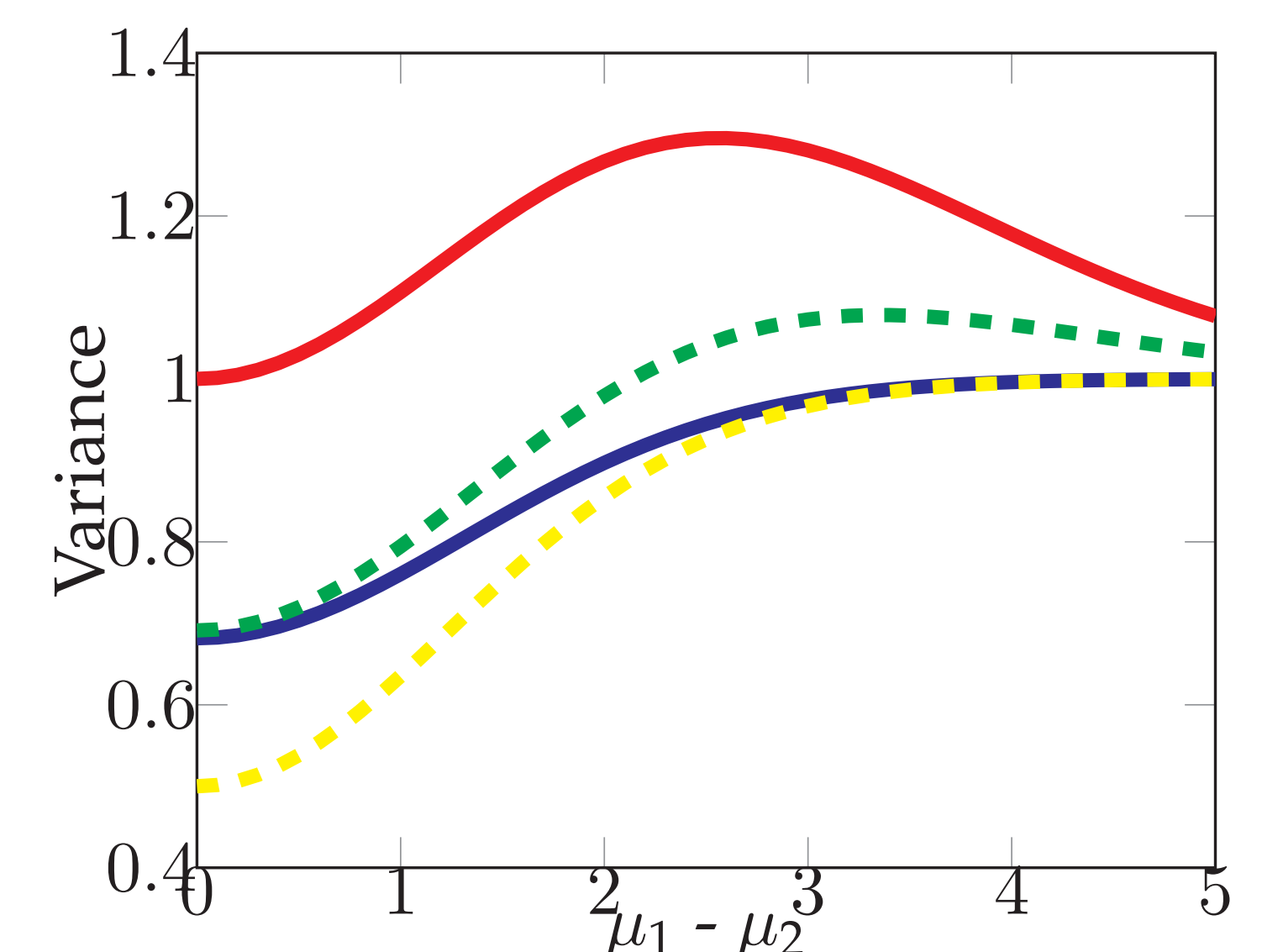
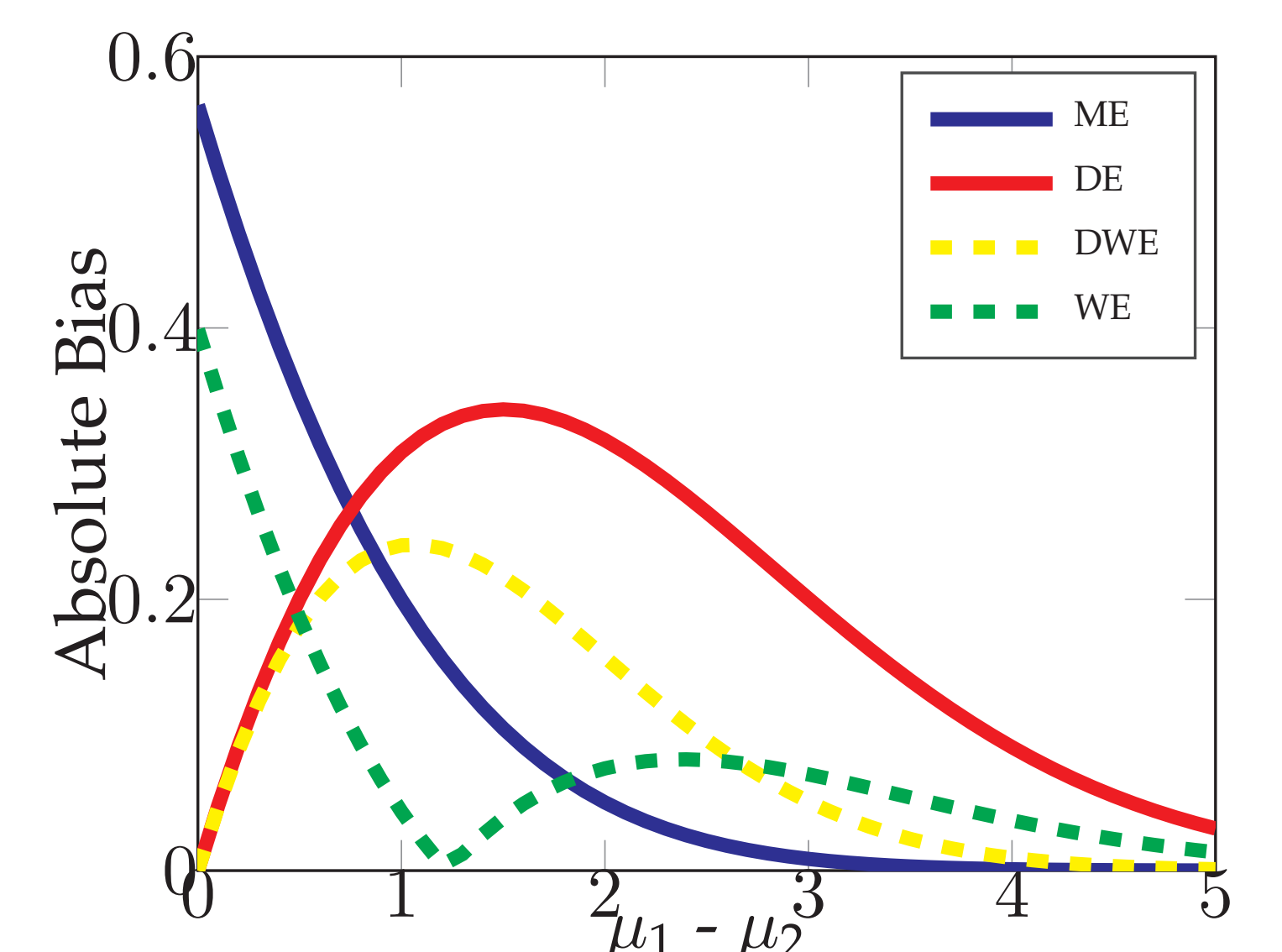
THEORETICAL RESULTS

BIAS AND VARIANCE OF THE ESTIMATORS

$$\text{Bias}(\hat{\mu}^{DE}) \geq -\frac{1}{2} \left(\sqrt{\sum_{i=1}^M \frac{\sigma_i^2}{|S_i^A|}} + \sqrt{\sum_{i=1}^M \frac{\sigma_i^2}{|S_i^B|}} \right)$$

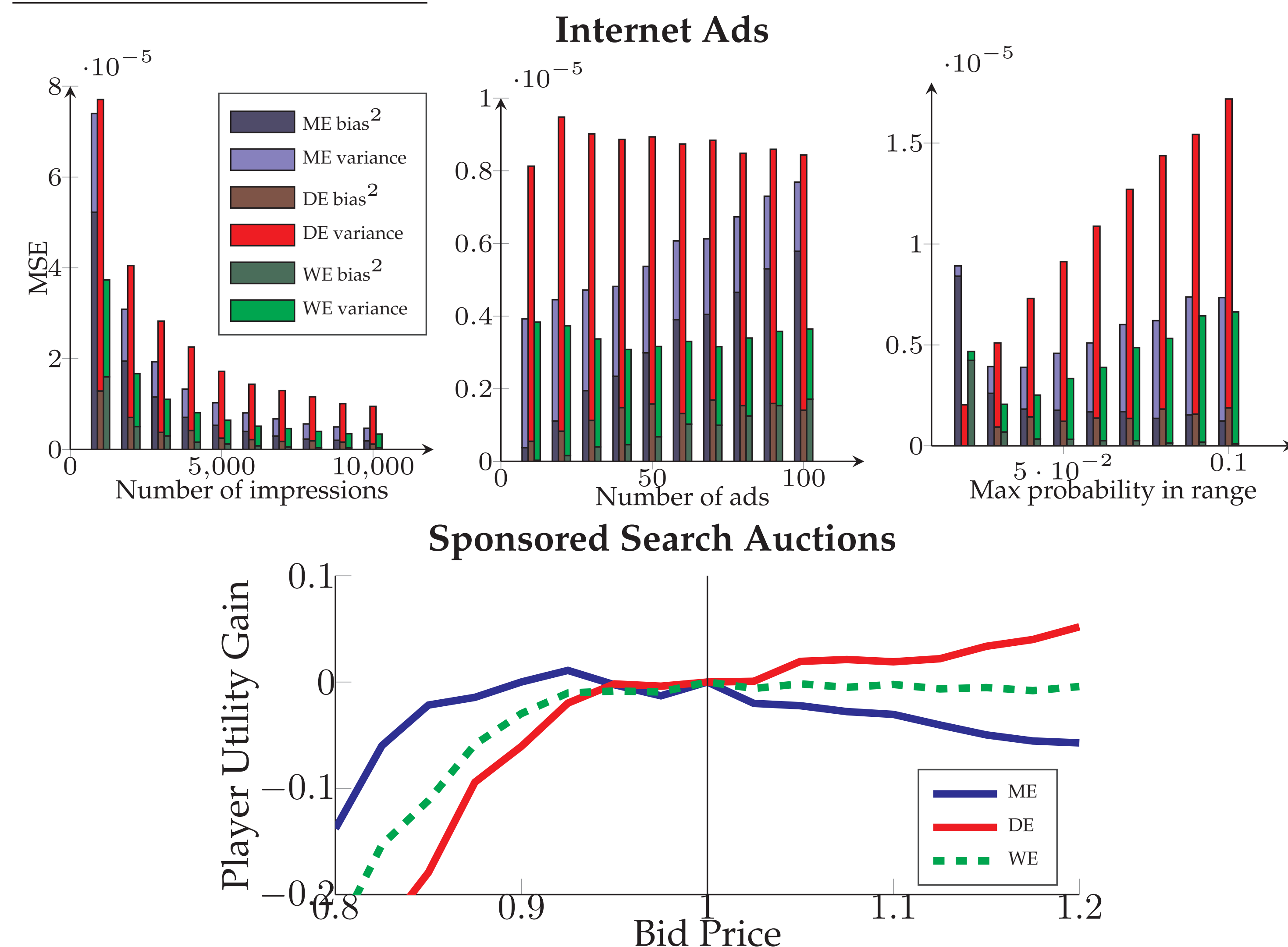
$$\text{Bias}(\hat{\mu}^{WE}) \leq \text{Bias}(\hat{\mu}^{ME}) \leq \sqrt{\frac{M-1}{M} \sum_{i=1}^M \frac{\sigma_i^2}{|S_i|}}$$

$$\text{Var}(\hat{\mu}^{ME, DE, WE}) \leq \sum_{i=1}^M \frac{\sigma_i^2}{|S_i|}$$



EMPIRICAL RESULTS

MULTI-ARMED BANDITS



MARKOV DECISION PROCESSES

