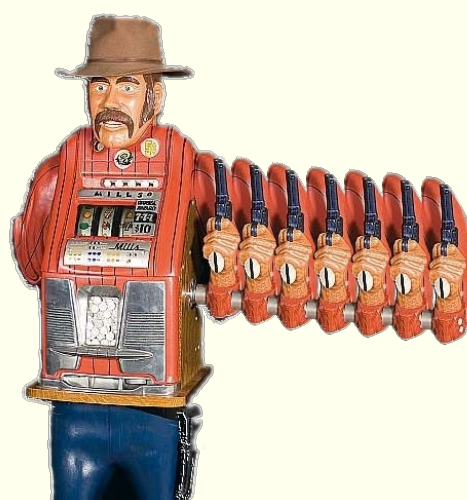


AAAI 2012

David Tolpin, Solomon Eyal Shimony
Ben Gurion University of the Negev,
Beer Sheva, Israel

- A set of K arms.
- Each arm can be pulled multiple times.
- When the i th arm is pulled, a random reward X_i is encountered.
- Simple regret: the reward of the last pull only is collected.
- Cumulative regret: all rewards are accumulated.



- **UCB(c)** pulls arm i that maximizes upper confidence bound b_i on the reward:

$$b_i = \overline{X}_i + \sqrt{\frac{c \log(n)}{n_i}}$$
- UCB is nearly optimal in minimizing the *cumulative regret*.
- **UCT** extends UCB to MCTS by invoking UCB in every node of a rollout.

- A problem-solving agent can perform *base-level* actions from a known set $\{A_i\}$.
- Before committing to an action, the agent may perform a sequence of *meta-level* deliberation actions from a set $\{S_j\}$.
- At any given time there is a base-level action A_α that maximizes the agent's *expected utility*.
- The **value of information** VOI_j is the expected difference between the expected utilities of the new and the old selected base-level action **after meta-level action S_j is taken.**
- The agent selects a meta-level action that **maximizes the VOI**, or A_α if no meta-level action has positive VOI.

- IMG4 Consortium under the MAGNET program of the Israeli Ministry of Trade and Industry
- Israel Science Foundation grant 305/09
- Lynne and William Frankel Center for Computer Sciences
- Paul Ivanier Center for Robotics Research and Production Management

- ## MAIN RESULTS

- Selects an action at **the current root** suitable for minimizing the **simple regret**.
- Then selects actions according to UCB, that approximately minimizes the **cumulative regret**.

Sampling for Simple Regret

1. ε -**greedy** sampling ($\varepsilon = \frac{1}{2}$).
2. Modified version of **UCB** (optimized for *simple regret*).
3. **VOI-aware** sampling:

$$VOI_{\alpha} \approx \frac{\bar{X}_{\beta}}{n_{\alpha} + 1} \exp\left(-2(\bar{X}_{\alpha} - \bar{X}_{\beta})^2 n_{\alpha}\right)$$

$$VOL_i \approx \frac{1 - \bar{X}_\alpha}{n_i + 1} \exp\left(-2(\bar{X}_\alpha - \bar{X}_i)^2 n_i\right), i \neq \alpha$$

- SR+CR outperforms UCT.
- SR+UCT(c) is less dependent on tuning of the exploration factor c .

Diagram illustrating a multi-armed bandit problem structure. A root node branches into several switches, each leading to two arms.

- Switch 1:**
 - Arm 11: $\mu_{11} = 0.6$
 - Arm 12: $\mu_{12} = 1 - \mu_{11} = 0.4$
- Switch 2:**
 - Arm 21: $\mu_{21} = 0.9 = \mu_*$
 - Arm 22: $\mu_{22} = 1 - \mu_{21} = 0.1$
- Switch K:**
 - Arm K1: $\mu_{K1} = \dots$
 - Arm K2: $\mu_{K2} = 1 - \mu_{K1} = \dots$

The diagram shows a sequence of switches and arms, with the second switch and its first arm highlighted in red, indicating a specific path or focus.

Nsamples	Uniform	UCT	1/2-greedy+UCT	UCB[sqrt]+UCT	VOI+UCT
100	0.10	0.08	0.08	0.08	0.08
200	0.10	0.075	0.07	0.07	0.07
400	0.10	0.065	0.055	0.055	0.055
600	0.10	0.055	0.045	0.045	0.04
800	0.10	0.045	0.035	0.035	0.025
1000	0.10	0.035	0.025	0.025	0.015
1200	0.10	0.025	0.018	0.018	0.01
1400	0.10	0.015	0.01	0.01	0.005
1500	0.10	0.01	0.008	0.008	0.004

The diagram on the left shows a sailboat on a grid. A dashed line indicates a path from the sailboat towards the top right corner of the grid, ending at a red square. The diagram on the right is a circular chart showing various sail positions around a central circle. The positions are labeled as follows:

- NO SAIL ZONE (at the top)
- CLOSE HAULED (top-left)
- CLOSE HAULED (top-right)
- CLOSE REACH (middle-left)
- CLOSE REACH (middle-right)
- BEAM REACH (bottom-left)
- BEAM REACH (bottom-right)
- BROAD REACH (bottom-left)
- BROAD REACH (bottom-right)
- RUIN (WIND ON WING) (at the bottom)
- FALLING OFF (bottom-right)
- HEAD-ON (top-left)

Samples per node	Uniform	Uniform+UCT	UCT	1/2-greedy+UCT	UCB[sqrt] + UCT
200	30.2	27.9	27.5	27.3	26.9
300	29.8	27.0	27.4	26.8	26.3
400	29.3	26.3	27.3	26.3	26.0
550	28.9	25.8	27.2	25.8	25.7
800	28.3	25.6	27.1	25.6	25.6
1050	28.0	25.5	27.0	25.5	25.5
1400	27.7	25.5	26.9	25.5	25.5

Exploration Factor	Uniform	Uniform+UCT	UCT	1/2-greedy+UCT	UCB[sqrt] + UCT
1e-05	29.2	26.2	27.6	26.2	26.0
1e-03	29.2	26.2	27.5	26.4	25.9
1e-01	29.2	26.1	27.4	26.1	25.8
1e+01	29.3	26.4	27.0	26.4	26.0
1e+03	28.1	28.1	28.1	28.1	27.8

- Improved MCTS scheme — SR+CR.
- SR+CR performs better than unmodified UCT.
- VOI-aware sampling for minimizing *simple regret*.

- Rational metareasoning in MCTS: theory and VOI estimates.
- Better sampling for non-root nodes.
- Application to Computer Go and other complex domains.

