

Verification of Discrete Upper Confidence Bound algorithm in Isabelle/HOL

Arjan Faber

February 6, 2026

Abstract

This project formally verifies the Upper Confidence Bound (UCB) algorithm in Isabelle/Higher-order Logic (HOL), focusing on its probabilistic guarantees and regret bounds. The work extends Isabelle/HOLs probabilistic framework and explores verification of discrete-time bandit models following [1]. This research advances the formal verification of probabilistic algorithms in reinforcement learning.

theory *MSc-Project-Discrete-Prop15-1*

imports

HOL-Probability.Probability

begin

locale *bandit-problem* =

fixes $A :: 'a \text{ set}$

and $\mu :: 'a \Rightarrow \text{real}$

and $a\text{-star} :: 'a$

assumes *finite-arms*: $\text{finite } A$

and *a-star-in-A*: $a\text{-star} \in A$

and *optimal-arm*: $\forall a \in A. \mu a\text{-star} \geq \mu a$

begin

definition $\Delta :: 'a \Rightarrow \text{real}$ **where**

$\Delta a = \mu a\text{-star} - \mu a$

end

locale *bandit-policy* = *bandit-problem* + *prob-space* +

fixes $\Omega :: 'b \text{ set}$

and $\mathcal{F} :: 'b \text{ set set}$

and $\pi :: \text{nat} \Rightarrow 'b \Rightarrow 'a$

and $N\text{-}n :: \text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{nat}$

assumes *measurable-policy*: $\forall t. \pi t \in \text{measurable } M \text{ (count-space } A)$

and *N-n-def*: $\forall n a \omega. N\text{-}n n a \omega = \text{card } \{t \in \{0..<n\}. \pi (t+1) \omega = a\}$

and *count-assm-pointwise*: $\forall n \omega. (\sum a \in A. \text{real } (N\text{-}n n a \omega)) = \text{real } n$

begin

definition $R\text{-}n :: \text{nat} \Rightarrow 'b \Rightarrow \text{real}$ **where**

$$R\text{-}n\ n\ \omega = n * \mu\ a\text{-}star - (\sum a \in A. \mu\ a * \text{real}\ (N\text{-}n\ n\ a\ \omega))$$

lemma *regret-decomposition-pointwise*:

fixes $n :: \text{nat}$ **and** $\omega :: 'b$

assumes *n-count-assm-pointwise*: $(\sum a \in A. \text{real}\ (N\text{-}n\ n\ a\ \omega)) = \text{real}\ n$

shows $R\text{-}n\ n\ \omega = (\sum a \in A. \Delta\ a * \text{real}\ (N\text{-}n\ n\ a\ \omega))$

<proof>

lemma *integrable-const-fun*:

assumes *finite-measure* M

shows *integrable* M $(\lambda x. c)$

<proof>

lemma *expected-regret*:

assumes *finite* A

and $\forall a \in A. \text{integrable}\ M\ (\lambda \omega. \text{real}\ (N\text{-}n\ n\ a\ \omega))$

shows *expectation* $(\lambda \omega. R\text{-}n\ n\ \omega) = (\sum a \in A. \Delta\ a * \text{expectation}\ (\lambda \omega. \text{real}\ (N\text{-}n\ n\ a\ \omega)))$

<proof>

end

end

theory *Discrete-UCB-Step1*

imports *MSc-Project-Discrete-Prop15-1*

begin

locale *bandit-policy* = *bandit-problem* + *prob-space* +

fixes $\Omega :: 'b\ \text{set}$

and $\mathcal{F} :: 'b\ \text{set}\ \text{set}$

and $\omega :: 'b$

and $\pi :: \text{nat} \Rightarrow 'b \Rightarrow 'a$

and $N\text{-}n :: \text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{nat}$

assumes *measurable-policy*: $\forall t. \pi\ t \in \text{measurable}\ M\ (\text{count-space}\ A)$

and *N-n-def*: $\forall n\ a\ \omega. N\text{-}n\ n\ a\ \omega = \text{card}\ \{t \in \{0..<n\}. \pi\ (t+1)\ \omega = a\}$

and *count-assm-pointwise*: $\forall n\ \omega. (\sum a \in A. \text{real}\ (N\text{-}n\ n\ a\ \omega)) = \text{real}\ n$

begin

lemma *union-eq*:

fixes $a :: 'a$ **and** $n\ k :: \text{nat}$

assumes $k \leq n$

shows $\{t. t < n \wedge \pi\ (t+1)\ \omega = a\} = \{t. t < k \wedge \pi\ (t+1)\ \omega = a\} \cup \{t. k \leq t$

$\wedge t < n \wedge \pi (t+1) \omega = a\}$
 $\langle proof \rangle$

lemma *cardinality-indic-eq*:

fixes $I :: nat \Rightarrow bool$
assumes $finite \{t. k \leq t \wedge t < n\}$
shows $card \{t. k \leq t \wedge t < n \wedge \pi (t+1) \omega = a \wedge I t\} = (\sum t = k..<n. if \pi (t+1) \omega = a \wedge I t then 1 else 0)$
 $\langle proof \rangle$

lemma *ge-rewrite*: $(x::real) \geq y \Longrightarrow y \leq x \langle proof \rangle$

lemma *Nn-expression*:

fixes $a :: 'a$ **and** $s :: nat \Rightarrow real$
and $k :: nat$ **and** $n :: nat$
assumes $a \in A$
and $k \leq n$
and $0 < n$
and $\forall t \in \{0..n\}. 0 < s t$
and $\forall t < n - 1. s t \leq s (t + 1)$
and *init-play-once*: $\forall \omega. a \in A \longrightarrow N-n k a \omega = 1$
and *finite-played-sets*:
 $finite \{t. t < n \wedge \pi (t+1) \omega = a\}$
 $finite \{t. t < k \wedge \pi (t+1) \omega = a\}$
 $finite \{t. k \leq t \wedge t < n \wedge \pi (t+1) \omega = a\}$
shows
 $(N-n n a \omega) = 1 + (\sum t = k..<n. if \pi (t+1) \omega = a \wedge real (N-n t a \omega) < s t then 1 else 0) +$
 $(\sum t = k..<n. if \pi (t+1) \omega = a \wedge real (N-n t a \omega) \geq s t then 1 else 0)$
 $\langle proof \rangle$

lemma *upper-bound-expression-contradiction*:

fixes $a :: 'a$ **and** $s :: nat \Rightarrow real$
and $k :: nat$ **and** $n :: nat$
and $s-n-nat :: nat$
assumes $a \in A$
and $k \leq n$
and $0 < n$
and *non-neg-s*: $\forall t \in \{0..n\}. 0 < s t$
and *base-le*: $s 0 \leq s 1$
and *non-dec*: $\forall t < n - 1. s t \leq s (t + 1)$
and *s-mono*: $\bigwedge t. k \leq t \wedge t \leq n \Longrightarrow s t \leq s n$
and *init-play-once*: $\forall \omega. a \in A \longrightarrow N-n k a \omega = 1$
and *finite-played-sets*:
 $finite \{t. t < n \wedge \pi (t+1) \omega = a\}$
 $finite \{t. t < k \wedge \pi (t+1) \omega = a\}$
 $finite \{t. k \leq t \wedge t < n \wedge \pi (t+1) \omega = a\}$
and *xs-sorted-def*: $xs = sorted-list-of-set \{t \in \{k..<n\}. \pi (t+1) \omega = a \wedge real$

$(N\text{-}n\ t\ a\ \omega) < s\ t\}$
and *s-nat-def*: $s\text{-}n\text{-}nat = nat\ (\lfloor s\ n \rfloor)$
and *len-bound-def*: $s\text{-}n\text{-}nat < length\ xs$
and *distinct-xs*: $distinct\ xs$
and *gt-ineq*: $length\ xs + 1 > \lfloor s\ n \rfloor$
and *N-n-increasing-with-plays*:
 $\forall t\ t'. k \leq t \wedge t < t' \wedge \pi\ (t+1)\ \omega = a \wedge \pi\ (t'+1)\ \omega = a \longrightarrow N\text{-}n\ t'\ a\ \omega \geq$
 $N\text{-}n\ t\ a\ \omega + 1$
and *neg*: $1 + real\ (\sum_{t=k..<n.} if\ \pi\ (t+1)\ \omega = a \wedge real\ (N\text{-}n\ t\ a\ \omega) < s\ t$
*then 1 else 0) > s\ n
and *t-hat* $\in set\ xs$
and *t-hat* $= xs\ !\ s\text{-}n\text{-}nat$
and $real\ (N\text{-}n\ t\text{-}hat\ a\ \omega) \geq real\ s\text{-}n\text{-}nat + 1$*

shows $(real\ (N\text{-}n\ t\text{-}hat\ a\ \omega) \geq \lfloor s\ n \rfloor + 1) \wedge (\pi\ (t\text{-}hat+1)\ \omega = a \wedge (real\ (N\text{-}n\ t\text{-}hat\ a\ \omega) < s\ t\text{-}hat))$

<proof>

lemma *Nn-upper-bound*:

fixes $a :: 'a$ **and** $s :: nat \Rightarrow real$
and $k :: nat$ **and** $n :: nat$
assumes *asm*: $real(1 + (\sum_{t=k..<n.} if\ \pi\ (t+1)\ \omega = a \wedge real\ (N\text{-}n\ t\ a\ \omega) <$
 $s\ t\ then\ 1\ else\ 0)) \leq s\ n$
and *a-in-A*: $a \in A$
and *k-le-n*: $k \leq n$
and *n-pos*: $0 < n$
and *s-pos*: $\forall t \in \{0..n\}. 0 < s\ t$
and *s-nondec*: $\forall t < n - 1. s\ t \leq s\ (t + 1)$
and *init-play-once*: $\forall \omega. a \in A \longrightarrow N\text{-}n\ k\ a\ \omega = 1$
and *finite-played-sets-1*: $finite\ \{t. t < n \wedge \pi\ (t+1)\ \omega = a\}$
and *finite-played-sets-2*: $finite\ \{t. t < k \wedge \pi\ (t+1)\ \omega = a\}$
and *finite-played-sets-3*: $finite\ \{t. k \leq t \wedge t < n \wedge \pi\ (t+1)\ \omega = a\}$
shows $real(N\text{-}n\ n\ a\ \omega) \leq s\ n + real((\sum_{t=k..<n.} if\ \pi\ (t+1)\ \omega = a \wedge s\ t \leq$
 $real\ (N\text{-}n\ t\ a\ \omega)\ then\ 1\ else\ 0))$
<proof>

theorem *ENn-upper-bound*:

assumes
a-in-A: $a \in A$
and *k-le-n*: $k \leq n$
and *n-pos*: $0 < n$

and *s-pos*: $\forall t \in \{0..n\}. 0 < s t$
and *s-nondec*: $\forall t < n. s t \leq s (t + 1)$
and *init-play-once*: $\forall \omega. a \in A \longrightarrow N\text{-}n \text{ } k \text{ } a \text{ } \omega = 1$
and *integrable-Nn*: *integrable* $M (\lambda \omega. \text{real} (N\text{-}n \text{ } n \text{ } a \text{ } \omega))$
and *integrable-rhs-sum*: *integrable* $M (\lambda \omega. s \text{ } n + (\sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0))$
and *integrable-s*: *integrable* $M (\lambda \omega. s \text{ } n)$
and *integrable-indicator-sum*: *integrable* $M (\lambda \omega. \sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0)$
and *linearity*: *integral*^L $M (\lambda \omega. s \text{ } n + (\sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0)) =$
 $\text{integral}^L M (\lambda \omega. s \text{ } n) + \text{integral}^L M (\lambda \omega. \sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0)$
and *pointwise-bound*: $\text{real} (N\text{-}n \text{ } n \text{ } a \text{ } \omega) \leq s \text{ } n + (\sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0)$
and *mono-intgrl*: *integral*^L $M (\lambda \omega. \text{real} (N\text{-}n \text{ } n \text{ } a \text{ } \omega)) \leq \text{integral}^L M (\lambda \omega. s \text{ } n + (\sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0))$
shows
 $\text{expectation} (\lambda \omega. \text{real} (N\text{-}n \text{ } n \text{ } a \text{ } \omega)) \leq$
 $s \text{ } n + \text{expectation} (\lambda \omega. (\sum t = k..<n. \text{if } \pi (t+1) \text{ } \omega = a \wedge s t \leq \text{real} (N\text{-}n \text{ } t \text{ } a \text{ } \omega) \text{ then } 1 \text{ else } 0))$
<proof>

end

end

theory *Discrete-UCB-Step2*

imports *Discrete-UCB-Step1*

begin

locale *bandit-policy* = *bandit-problem* + *prob-space* +

fixes $\Omega :: 'b \text{ set}$

and $\mathcal{F} :: 'b \text{ set set}$

and $\omega :: 'b$

fixes $\pi :: \text{nat} \Rightarrow 'b \Rightarrow 'a$

and $N\text{-}n :: \text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{nat}$

and $Z :: \text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$

and $\delta :: \text{real}$

and $q :: \text{real}$

assumes *finite-A*: *finite* A

and *a-in-A*: $a \in A$

and *measurable-policy*: $\forall t. \pi \text{ } t \in \text{measurable } M (\text{count-space } A)$

and *N-n-def*: $\forall n \text{ } a \text{ } b. N\text{-}n \text{ } n \text{ } a \text{ } b = \text{card} \{t \in \{0..<n\}. \pi (t+1) \text{ } b = a\}$

and *delta-pos*: $0 < \delta$

and *delta-less1*: $\delta < 1$

and *q-pos*: $q > 0$

begin

definition *sample-mean-Z* :: $\text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$ **where**

$$\text{sample-mean-Z } n \ a \ b \equiv (1 / \text{real } n) * (\sum_{i < n}. Z \ i \ a \ b)$$

definition *M-val* :: $\text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$ **where**

$$\begin{aligned} M\text{-val } t \ a \ b &\equiv (\text{if } N\text{-n } (t+1) \ a \ b = 0 \ \text{then } 0 \\ &\quad \text{else } (\sum_{s < t}. \text{if } \pi \ s \ b = a \ \text{then } Z \ s \ a \ b \ \text{else } 0) / \text{real } (N\text{-n } t \ a \ b)) \end{aligned}$$

definition *U* :: $\text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$ **where**

$$U \ t \ a \ b \equiv M\text{-val } t \ a \ b + q * \text{sqrt } (\ln (1 / \delta) / (2 * \text{real } (\max 1 (N\text{-n } t \ a \ b))))$$

definition *A-t-plus-1* :: $\text{nat} \Rightarrow 'b \Rightarrow 'a$ **where**

$$A\text{-t-plus-1 } t \ b \equiv (\text{SOME } a. a \in A \wedge (\forall a'. a' \in A \longrightarrow U \ t \ a \ b \geq U \ t \ a' \ b))$$

lemma (in *finite-measure*) *finite-measure-mono*:

assumes $A \subseteq B \ B \in \text{sets } M$ **shows** $\text{measure } M \ A \leq \text{measure } M \ B$

<proof>

theorem *union-bound*:

fixes $E \ F \ G :: 'b \ \text{set}$

assumes $E \subseteq F \cup G$

and $E \in \text{events } F \in \text{events } G \in \text{events}$

shows $\text{prob } E \leq \text{prob } F + \text{prob } G$

<proof>

theorem *hoeffding-iid-bound-ge-general*:

fixes $a :: 'a$ **and** $n :: \text{nat}$ **and** $\varepsilon :: \text{real}$ **and** $\mu\text{-hat} :: \text{real}$ **and** $l \ u :: \text{real}$

assumes *a-in*: $a \in A$

and *eps-pos*: $\varepsilon \geq 0$

and *bounds*: $\forall i < n. \forall \omega \in \Omega. l \leq Z \ i \ a \ \omega \wedge Z \ i \ a \ \omega \leq u$

and *mu-def*: $\mu\text{-hat} = (\sum_{i < n}. \text{expectation } (\lambda \omega. Z \ i \ a \ \omega))$

and $u - l \neq 0$

and *n-pos*: $n > 0$

and *space-M*: $\text{space } M = \Omega$

and *sets-M*: $\text{sets } M = \mathcal{F}$

and *indep*: $\text{indep-vars } (\lambda \cdot. \text{borel}) (\lambda i. (\lambda \omega. Z \ i \ a \ \omega)) \{i. i < n\}$

and *rv*: $\forall i < n. \text{random-variable borel } (\lambda \omega. Z \ i \ a \ \omega)$

shows $\text{prob } \{\omega \in \Omega. (\sum_{i < n}. Z \ i \ a \ \omega) \geq \mu\text{-hat} + \varepsilon\}$

$$\leq \exp (- 2 * \varepsilon^2 / (\text{real } n * (u - l)^2))$$

<proof>

theorem *hoeffding-iid-bound-le-general*:

fixes $a :: 'a$ **and** $n :: \text{nat}$ **and** $\varepsilon :: \text{real}$ **and** $\mu\text{-hat} :: \text{real}$ **and** $l \ u :: \text{real}$

assumes *a-in*: $a \in A$

and *eps-pos*: $\varepsilon \geq 0$

and bounds: $\forall i < n. \forall \omega \in \Omega. l \leq Z i a \omega \wedge Z i a \omega \leq u$
and mu-def: $\mu\text{-hat} = (\sum i < n. \text{expectation } (\lambda \omega. Z i a \omega))$
and u - l \neq 0
and n-pos: $n > 0$
and space-M: $\text{space } M = \Omega$
and sets-M: $\text{sets } M = \mathcal{F}$
and indep: $\text{indep-vars } (\lambda \cdot. \text{borel}) (\lambda i. (\lambda \omega. Z i a \omega)) \{i. i < n\}$
and rv: $\forall i < n. \text{random-variable borel } (\lambda \omega. Z i a \omega)$
shows $\text{prob } \{\omega \in \Omega. (\sum i < n. Z i a \omega) \leq \mu\text{-hat} - \varepsilon\}$
 $\leq \text{exp } (- 2 * \varepsilon^2 / (\text{real } n * (u - l)^2))$
 <proof>

theorem hoeffding-iid-ge-delta-bound:

fixes $a :: 'a$ **and** $n :: \text{nat}$ **and** $\delta\text{-hat} :: \text{real}$ **and** $\mu\text{-hat} :: \text{real}$ **and** $l u :: \text{real}$
assumes $a\text{-in}: a \in A$
and delta-bound: $0 < \delta\text{-hat} \wedge \delta\text{-hat} \leq 1$
and bounds: $\forall i < n. \forall \omega \in \Omega. l \leq Z i a \omega \wedge Z i a \omega \leq u$
and mu-def: $\mu\text{-hat} = (\sum i < n. \text{expectation } (\lambda \omega. Z i a \omega))$
and n-pos: $n > 0$
and eps-pos: $\varepsilon \geq 0$
and u-minus-l-nonzero: $u - l \neq 0$
and space-M: $\text{space } M = \Omega$
and sets-M: $\text{sets } M = \mathcal{F}$
and indep: $\text{indep-vars } (\lambda \cdot. \text{borel}) (\lambda i. (\lambda \omega. Z i a \omega)) \{i. i < n\}$
and rv: $\forall i < n. \text{random-variable borel } (\lambda \omega. Z i a \omega)$
and eps-expression: $\varepsilon = \text{sqrt } ((\text{real } n * (u - l)^2 * \ln (1 / \delta\text{-hat})) / 2)$
shows $\text{prob } \{\omega \in \Omega. (\sum i < n. Z i a \omega) \geq \mu\text{-hat} + \varepsilon\} \leq \delta\text{-hat} \wedge$
 $\text{prob } \{\omega \in \Omega. (\sum i < n. Z i a \omega) \leq \mu\text{-hat} - \varepsilon\} \leq \delta\text{-hat}$
 <proof>

lemma add-le-iff:

fixes $x y z :: \text{real}$
shows $x \leq y - z \iff x - y \leq -z$
 <proof>

lemma max-Suc-0-eq-1: $\text{max } (\text{Suc } 0) x = \text{max } 1 x$
 <proof>

theorem ucb-suboptimal-bound-set:

fixes $t :: \text{nat}$
and $a :: 'a$
and $\Delta :: 'a \Rightarrow \text{real}$
assumes $\text{finite-A}: \text{finite } A$
and $a\text{-in-A}: a \in A$
and $a\text{-star-in-A}: a\text{-star} \in A$
and $\text{argmax-exists}: A \neq \{\}$
and $\text{subopt-gap}: \Delta a > 0$
and $a\text{-not-opt}: \exists a'. a' \in A \wedge \Delta a' > 0$
and $\text{delta-a}: \Delta a = \mu a\text{-star} - \mu a$

and ω -in- Ω : $\omega \in \Omega$
and asm : $\omega \in \{\omega \in \Omega. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a\}$
and $setopt$: $\forall \omega \in \Omega. \exists a\text{-}max \in A. \forall b \in A. U\ t\ b\ \omega \leq U\ t\ a\text{-}max\ \omega$
and $A\text{-}t\text{-}plus\text{-}1\text{-}maximizes$:
 $\bigwedge t\ \omega\ a. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a \implies a \in A \wedge (\forall b \in A. U\ t\ a\ \omega \geq U\ t\ b\ \omega)$
shows $\{\omega \in \Omega. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a\} \subseteq$
 $\{\omega \in \Omega. U\ t\ a\text{-}star\ \omega \leq \mu\ a\text{-}star\} \cup \{\omega \in \Omega. \mu\ a\text{-}star \leq U\ t\ a\ \omega\}$
 $\langle proof \rangle$

theorem $ucb\text{-}suboptimal\text{-}bound\text{-}prob\text{-}statement$:

fixes $t :: nat$ **and** $a :: 'a$ **and** $\Delta :: 'a \Rightarrow real$
assumes $finite\text{-}A$: $finite\ A$
and $a\text{-}star\text{-}in\text{-}A$: $a\text{-}star \in A$
and $argmax\text{-}exists$: $A \neq \{\}$
and $subopt\text{-}gap$: $\Delta\ a > 0$
and $a\text{-}not\text{-}opt$: $\exists a'. a' \in A \wedge \Delta\ a > 0$
and ω -in- Ω : $\omega \in \Omega$
and asm : $\omega \in \{\omega \in \Omega. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a\}$
and $setopt$: $\forall \omega \in \Omega. \exists a\text{-}max \in A. \forall b \in A. U\ t\ b\ \omega \leq U\ t\ a\text{-}max\ \omega$
and $A\text{-}t\text{-}plus\text{-}1\text{-}maximizes$:
 $\bigwedge t\ \omega\ a. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a \implies a \in A \wedge (\forall b \in A. U\ t\ a\ \omega \geq U\ t\ b\ \omega)$
and $a\text{-}in\text{-}A$: $a \in A$
and $omega\text{-}in$: $\omega \in \Omega$
and $subopt\text{-}gap$: $\Delta\ a > 0$
and $delta\text{-}a$: $\Delta\ a = \mu\ a\text{-}star - \mu\ a$
and $H\text{-}def$: $H = (2 * \ln (1 / \delta)) / (\Delta\ a)^{\wedge}2$
and $E\text{-}def$: $E = \{\omega \in \Omega. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a\} \cap \{\omega \in \Omega. H \leq real\ (N\text{-}n\ t\ a\ \omega)\}$
and $F\text{-}def$: $F = \{\omega \in \Omega. U\ t\ a\text{-}star\ \omega \leq \mu\ a\text{-}star\} \cap \{\omega \in \Omega. H \leq real\ (N\text{-}n\ t\ a\ \omega)\}$
and $G\text{-}def$: $G = \{\omega \in \Omega. \mu\ a\text{-}star \leq U\ t\ a\ \omega\} \cap \{\omega \in \Omega. H \leq real\ (N\text{-}n\ t\ a\ \omega)\}$
and $meas\text{-}sets$: $E \in sets\ M\ F \in sets\ M\ G \in sets\ M$
and $prob\text{-}inter$: $prob\ (F \cap G) \equiv enn2real\ (emeasure\ M\ (F \cap G))$

shows $prob\ (\{\omega \in \Omega. A\text{-}t\text{-}plus\text{-}1\ t\ \omega = a\} \cap \{\omega \in \Omega. H \leq real\ (N\text{-}n\ t\ a\ \omega)\}) \leq$
 $prob\ (\{\omega \in \Omega. U\ t\ a\text{-}star\ \omega \leq \mu\ a\text{-}star\} \cap \{\omega \in \Omega. H \leq real\ (N\text{-}n\ t\ a\ \omega)\})$
 $+$
 $prob\ (\{\omega \in \Omega. U\ t\ a\ \omega \geq \mu\ a\text{-}star\} \cap \{\omega \in \Omega. H \leq real\ (N\text{-}n\ t\ a\ \omega)\})$
 $\langle proof \rangle$

lemma $U\text{-}le\text{-}\mu\text{-}pointwise$:

$U\ t\ a\text{-}star\ \omega \leq \mu\ a\text{-}star \iff$
 $M\text{-}val\ t\ a\text{-}star\ \omega - \mu\ a\text{-}star \leq$
 $- q * sqrt\ (\ln\ (1 / \delta) / (2 * real\ (max\ 1\ (N\text{-}n\ t\ a\text{-}star\ \omega))))$
 $\langle proof \rangle$

lemma *U-ge-μ-pointwise:*

assumes *delta-a:* $\Delta a = \mu a\text{-star} - \mu a$

shows

$U t a \omega \geq \mu a\text{-star} \longleftrightarrow$

$M\text{-val } t a \omega - \mu a \geq \Delta a - q * \text{sqrt} (\ln (1 / \delta) / (2 * \text{real} (\max 1 (N\text{-n } t a \omega))))$

<proof>

theorem *hoeffding-iid-bound-le:*

fixes *a* :: 'a **and** *n* :: nat **and** ε :: real **and** $\mu\text{-hat}$:: real **and** *l-hat* *u-hat* :: real

and *I* :: nat set

and *X-new* :: nat \Rightarrow 'b \Rightarrow real

and *a-bound* *b-bound* :: nat \Rightarrow real

assumes *a-in:* $a \in A$

and $b \in \Omega$

and *eps-pos:* $\varepsilon \geq 0$

and *eps:* $\varepsilon = \text{abs} (u\text{-hat} - l\text{-hat}) * \text{sqrt} (((\text{real } n) / 2) * \ln (1 / \delta))$

and $\delta \geq 0 \wedge \delta \leq 1$

and *t-eq-n:* $t = n$

and $c > 0$

and *bounds:* $\forall j < t. \forall \omega \in \Omega. \forall a \in A. l\text{-hat} \leq Z\text{-hat } j a \omega \wedge Z\text{-hat } j a \omega \leq u\text{-hat}$

and *mu-def:* $\mu\text{-hat} = (\sum j < t. \text{expectation} (\lambda \omega. Z\text{-hat } j a \omega))$

and $u\text{-hat} - l\text{-hat} \neq 0$

and *t-pos:* $t > 0$

and $\forall t < n. N\text{-n } t a\text{-star } b > 0$

and *n-pos:* $n > 0$

and *M-val* $t a\text{-star } b \equiv (\sum s < t. \text{if } \pi s b = a\text{-star} \text{ then } Z s a\text{-star } b \text{ else } 0) / \text{real} (N\text{-n } t a\text{-star } b)$

and *widths:* $(\sum i \in I. (b\text{-bound } i - a\text{-bound } i)^2) = (\text{real } n) * (u\text{-hat} - l\text{-hat})^2$

and *space-M:* $\text{space } M = \Omega$

and *sets-M:* $\text{sets } M = \mathcal{F}$

and *indep:* *indep-vars* $(\lambda \cdot. \text{borel}) (\lambda j. (\lambda \omega. Z j a \omega)) \{j. j < t\}$

and *rv:* $\forall j < t. \text{random-variable borel} (\lambda \omega. Z j a \omega)$

and $\forall j < t. Z\text{-hat } j a\text{-star } \omega = c * (\text{if } \pi j b = a\text{-star} \text{ then } Z j a\text{-star } b \text{ else } 0)$

and *indep-interval-bounded-random-variables* $M I X\text{-new } a\text{-bound } b\text{-bound}$

and *indep-loc:* *indep-interval-bounded-random-variables* $M \{j. j < t\}$

$(\lambda j. (\lambda \omega. Z\text{-hat } j a\text{-star } \omega))$

$(\lambda j. l\text{-hat}) (\lambda j. u\text{-hat})$

and *H:* *Hoeffding-ineq* $M \{j. j < t\}$

$(\lambda j. (\lambda \omega. Z\text{-hat } j a\text{-star } \omega))$

$(\lambda j. l\text{-hat}) (\lambda j. u\text{-hat})$

and *sum-integrals-eq:* $(\sum j \in \{j. j < t\}. \text{integral}^L M (\lambda \omega. Z\text{-hat } j a\text{-star } \omega)) = \mu\text{-hat}$

and *rewriting:* $\text{prob} \{\omega \in \Omega. (\sum j < t. Z\text{-hat } j a\text{-star } \omega) - (\sum j < t. \text{expectation} (Z\text{-hat } j a\text{-star})) \leq -\varepsilon\} =$

$\text{prob} \{x \in \text{space } M. (\sum j < t. Z\text{-hat } j a\text{-star } x) \leq (\sum j < t. \text{expectation} (Z\text{-hat } j a\text{-star})) - \varepsilon\}$

shows $\text{prob } \{\omega \in \Omega. (\sum j < n. Z\text{-hat } j \text{ a-star } \omega) \leq \mu\text{-hat} - \varepsilon\}$
 $\leq \delta$

<proof>

theorem *hoeffding-iid-bound-ge*:

fixes $a :: 'a$ **and** $n :: \text{nat}$ **and** $\varepsilon :: \text{real}$ **and** $\mu\text{-hat} :: \text{real}$ **and** $l\text{-hat } u\text{-hat} :: \text{real}$
and $I :: \text{nat set}$
and $X\text{-new} :: \text{nat} \Rightarrow 'b \Rightarrow \text{real}$
and $a\text{-bound } b\text{-bound} :: \text{nat} \Rightarrow \text{real}$
assumes $a\text{-in}: a \in A$
and $b \in \Omega$
and $\text{eps-pos}: \varepsilon \geq 0$
and $\text{eps}: \varepsilon = \text{abs } (u\text{-hat} - l\text{-hat}) * \text{sqrt } (((\text{real } n) / 2) * \ln (1 / \delta))$
and $\delta \geq 0 \wedge \delta \leq 1$
and $t\text{-eq-n}: t = n$
and $c > 0$
and $\text{bounds}: \forall j < t. \forall \omega \in \Omega. \forall a \in A. l\text{-hat} \leq Z\text{-hat } j \text{ a } \omega \wedge Z\text{-hat } j \text{ a } \omega \leq$
 $u\text{-hat}$
and $\mu\text{-def}: \mu\text{-hat} = (\sum j < t. \text{expectation } (\lambda\omega. Z\text{-hat } j \text{ a } \omega))$
and $u\text{-hat} - l\text{-hat} \neq 0$
and $t\text{-pos}: t > 0$
and $\forall t < n. N\text{-n } t \text{ a-star } b > 0$
and $n\text{-pos}: n > 0$
and $M\text{-val } t \text{ a } b \equiv (\sum s < t. \text{if } \pi \text{ s } b = a \text{ then } Z \text{ s } a \text{ b else } 0) /$
 $\text{real } (N\text{-n } t \text{ a } b)$
and $\text{widths}: (\sum i \in I. (b\text{-bound } i - a\text{-bound } i)^2) = (\text{real } n) * (u\text{-hat} -$
 $l\text{-hat})^2$
and $\text{space-M}: \text{space } M = \Omega$
and $\text{sets-M}: \text{sets } M = \mathcal{F}$
and $\text{indep}: \text{indep-vars } (\lambda\cdot. \text{borel}) (\lambda j. (\lambda\omega. Z \text{ j } a \omega)) \{j. j < t\}$
and $\text{rv}: \forall j < t. \text{random-variable borel } (\lambda\omega. Z \text{ j } a \omega)$
and $\forall j < t. Z\text{-hat } j \text{ a } \omega = c * (\text{if } \pi \text{ j } b = a \text{ then } Z \text{ j } a\text{-star } b \text{ else } 0)$
and $\text{indep-interval-bounded-random-variables } M \text{ I } X\text{-new } a\text{-bound } b\text{-bound}$
and $\text{indep-loc}: \text{indep-interval-bounded-random-variables } M \{j. j < t\}$
 $(\lambda j. (\lambda\omega. Z\text{-hat } j \text{ a } \omega))$
 $(\lambda j. l\text{-hat}) (\lambda j. u\text{-hat})$
and $H: \text{Hoeffding-ineq } M \{j. j < t\}$
 $(\lambda j. (\lambda\omega. Z\text{-hat } j \text{ a } \omega))$
 $(\lambda j. l\text{-hat}) (\lambda j. u\text{-hat})$
and $\text{sum-integrals-eq}: (\sum j \in \{j. j < t\}. \text{integral}^L M (\lambda\omega. Z\text{-hat } j \text{ a } \omega)) =$
 $\mu\text{-hat}$
and $\text{rewriting}: \text{prob } \{\omega \in \Omega. (\sum j < t. Z\text{-hat } j \text{ a } \omega) - (\sum j < t. \text{expectation}$
 $(Z\text{-hat } j \text{ a})) \geq \varepsilon\} =$
 $\text{prob } \{x \in \text{space } M. (\sum j < t. Z\text{-hat } j \text{ a } x) \geq (\sum j < t. \text{expectation } (Z\text{-hat } j$
 $a)) + \varepsilon\}$
shows $\text{prob } \{\omega \in \Omega. (\sum j < n. Z\text{-hat } j \text{ a } \omega) \geq \mu\text{-hat} + \varepsilon\} \leq \delta$
<proof>

```

end
end
theory Discrete-UCB-Step3
  imports Discrete-UCB-Step2

begin

locale bandit-policy = bandit-problem + prob-space +
  fixes
     $\Omega :: 'b \text{ set}$ 
    and  $\mathcal{F} :: 'b \text{ set set}$ 
    and  $\pi :: \text{nat} \Rightarrow 'b \Rightarrow 'a$ 
    and  $N\text{-}n :: \text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{nat}$ 
    and  $Z :: \text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$ 
    and  $\delta :: \text{real}$ 
    and  $q :: \text{real}$ 
  assumes fin-A: finite A
    and  $\omega \in \Omega$ 
    and a-in-A:  $a \in A$ 
    and measurable-policy:  $\forall t. \pi t \in \text{measurable } M \text{ (count-space } A)$ 
    and N-n-def:  $\forall n a \omega. N\text{-}n n a \omega = \text{card } \{t \in \{0..<n\}. \pi (t+1) \omega = a\}$ 
    and count-assm-pointwise:  $\forall n \omega. (\sum a \in A. \text{real } (N\text{-}n n a \omega)) = \text{real } n$ 
    and delta-pos:  $0 < \delta$ 
    and delta-less1:  $\delta < 1$ 
    and q-pos:  $q > 0$ 

begin

definition sample-mean-Z ::  $\text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$  where
  sample-mean-Z n a  $\omega = (1 / \text{real } n) * (\sum i < n. Z i a \omega)$ 

definition M-fun ::  $\text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$  where
  M-fun t a  $\omega = (\text{if } N\text{-}n (t+1) a \omega = 0 \text{ then } 0$ 
    else  $(\sum s < t. (\text{if } \pi s \omega = a \text{ then } Z s a \omega \text{ else } 0)) / \text{real } (N\text{-}n t a \omega))$ )

definition U ::  $\text{nat} \Rightarrow 'a \Rightarrow 'b \Rightarrow \text{real}$  where
  U t a  $\omega = M\text{-}fun t a \omega + q * \text{sqrt } (\ln (1 / \delta) / (2 * \text{real } (\max 1 (N\text{-}n t a \omega))))$ 

definition A-t-plus-1 ::  $\text{nat} \Rightarrow 'b \Rightarrow 'a$  where
  A-t-plus-1 t  $\omega = (\text{SOME } a. a \in A \wedge (\forall a'. a' \in A \longrightarrow U t a \omega \geq U t a' \omega))$ 

definition prob-eq-Ex ::  $'b \text{ set} \Rightarrow \text{bool}$  where
  prob-eq-Ex E  $\equiv \text{prob } E = \text{expectation } (\lambda \omega. \text{indicator } E \omega)$ 

theorem proposition-15-7:
  assumes

```

a-in-A: $a \in A$
and $\omega \in \Omega$
and subopt-gap: $\Delta a > 0$
and a-not-opt: $\exists a' \in A. \Delta a' > 0$
and delta-a: $\forall a \in A. \Delta a = \mu a\text{-star} - \mu a$
and $k \leq n$
and from-UCB-step1: $\forall a \in A. \text{expectation} (\lambda\omega. \text{real} (N-n \ n \ a \ \omega)) \leq$
 $s \ n + \text{expectation} (\lambda\omega. (\sum_{t=k..<n} \text{if } \pi (t+1) \ \omega = a \wedge s \ t \leq \text{real} (N-n \ t$
 $a \ \omega) \text{ then } 1 \text{ else } 0))$
and from-UCB-step2: $\forall a \in A. \forall t \in \{k..<n\}. \text{prob} (\{\omega \in \Omega. A\text{-}t\text{-plus-1 } t \ \omega$
 $= a\} \cap$
 $\{\omega \in \Omega. N\text{-}n \ t \ a \ \omega \geq (2 * \ln (1 / \delta)) / (\Delta a)^{\wedge 2}\}) \leq 2 * \delta$
and eps-pos: $\varepsilon > 0$
and t-gt0: $\forall t \in \{k..<n\}. t > 0$
and choice-delta: $\forall t \in \{k..<n\}. \delta = 1 / (\text{real } t \ \text{powr } \varepsilon)$
and s-form: $\forall a \in A. \forall u. s \ u = (2 * \varepsilon * \ln (\text{real } u)) / ((\Delta a)^{\wedge 2})$
and subset-meas: $\forall a \in A. \forall t \in \{k..<n\}. \forall \omega \in \Omega. \{\omega. \pi (t+1) \ \omega = a \wedge 2 * \varepsilon$
 $* \ln (\text{real } t) / (\Delta a)^{\wedge 2} \leq N\text{-}n \ t \ a \ \omega\} \subseteq \Omega$
and prob-eq-E-assm: $\forall a \in A. \forall t \in \{k..<n\}. \text{prob} \{\omega. \pi (Suc \ t) \ \omega = a \wedge 2 * \varepsilon$
 $* \ln (\text{real } t) / (\Delta a)^{\wedge 2} \leq \text{real} (N\text{-}n \ t \ a \ \omega)\} =$
 $\text{prob} (\{\omega. \pi (Suc \ t) \ \omega = a \wedge 2 * \varepsilon * \ln (\text{real } t) / (\Delta a)^{\wedge 2} \leq \text{real}$
 $(N\text{-}n \ t \ a \ \omega)\} \cap \text{space } M)$
and finiteness: $\forall t \in \{k..<n\}. \forall a \in A. \text{emeasure } M \{\omega. \pi (t+1) \ \omega = a \wedge 2 * \varepsilon * \ln (\text{real}$
 $t) / (\Delta a)^{\wedge 2} \leq \text{real} (N\text{-}n \ t \ a \ \omega)\} < \infty$
and measurable-set: $\forall t \in \{k..<n\}. \forall a \in A. \{\omega. \pi (t+1) \ \omega = a \wedge 2 * \varepsilon * \ln (\text{real}$
 $t) / (\Delta a)^{\wedge 2} \leq \text{real} (N\text{-}n \ t \ a \ \omega)\} \in \text{sets } M$

and eq-sets-optimum:
 $\forall a \in A. \forall t \in \{k..<n\}. \{\omega. \pi (t+1) \ \omega = a \wedge 2 * \varepsilon * \ln (\text{real } t) / (\Delta$
 $a)^{\wedge 2} \leq \text{real} (N\text{-}n \ t \ a \ \omega)\} =$
 $\{\omega \in \Omega. A\text{-}t\text{-plus-1 } t \ \omega = a\} \cap \{\omega \in \Omega. N\text{-}n \ t \ a \ \omega \geq (2 * \varepsilon * \ln$
 $(\text{real } t)) / (\Delta a)^{\wedge 2}\}$

shows

$\forall a \in A. \text{expectation} (\lambda\omega. \text{real} (N-n \ n \ a \ \omega)) \leq s \ n + (\sum_{t=k..<n} 2 / (\text{real}$
 $t \ \text{powr } \varepsilon))$

<proof>

theorem theorem-15-4:

assumes

a-in-A: $a \in A$
and finite A and $\forall a \in A. \text{integrable } M (\lambda\omega. \text{real} (N-n \ n \ a \ \omega))$
and $\omega\text{-in-}\Omega: \omega \in \Omega$
and subopt-gap: $\forall a \in A. \Delta a > 0$
and a-not-opt: $\exists a' \in A. \Delta a' > 0$
and delta-a: $\forall a \in A. \Delta a = \mu a\text{-star} - \mu a$
and $k \leq n$
and n-count-assm-pointwise: $(\sum_{a \in A} \text{real} (N-n \ n \ a \ \omega)) = \text{real } n$

and expected-regret-prop-15-1: $expectation (\lambda\omega. R-n n \omega) = (\sum a \in A. \Delta a * expectation (\lambda\omega. real (N-n n a \omega)))$
and from-UCB-step1: $\forall a \in A. expectation (\lambda\omega. real (N-n n a \omega)) \leq s n + expectation (\lambda\omega. (\sum t = k..<n. if \pi (t+1) \omega = a \wedge s t \leq real (N-n t a \omega) then 1 else 0))$
and from-UCB-step2: $\forall a \in A. \forall t \in \{k..<n\}. prob (\{\omega \in \Omega. A-t-plus-1 t \omega = a\} \cap \{\omega \in \Omega. N-n t a \omega \geq (2 * \ln (1 / \delta)) / (\Delta a)^2\}) \leq 2 * \delta$
and eps-pos: $\varepsilon > 0$
and t-gt0: $\forall t \in \{k..<n\}. t > 0$
and choice-delta: $\forall t \in \{k..<n\}. \delta = 1 / (real t powr \varepsilon)$
and s-form: $\forall a \in A. \forall u. s u = (2 * \varepsilon * \ln (real u)) / ((\Delta a)^2)$
and subset-meas: $\forall a \in A. \forall t \in \{k..<n\}. \forall \omega \in \Omega. \{\omega. \pi (t+1) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq N-n t a \omega\} \subseteq \Omega$
and prob-eq-E-assm: $\forall a \in A. \forall t \in \{k..<n\}. prob \{\omega. \pi (Suc t) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega)\} = prob (\{\omega. \pi (Suc t) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega)\} \cap space M)$
and finiteness: $\forall t \in \{k..<n\}. \forall a \in A. emeasure M \{\omega. \pi (t+1) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega)\} < \infty$
and measurable-set: $\forall t \in \{k..<n\}. \forall a \in A. \{\omega. \pi (t+1) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega)\} \in sets M$

and eq-sets-optimum:
 $\forall a \in A. \forall t \in \{k..<n\}. \{\omega. \pi (t+1) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega)\} = \{\omega \in \Omega. A-t-plus-1 t \omega = a\} \cap \{\omega \in \Omega. N-n t a \omega \geq (2 * \varepsilon * \ln (real t)) / (\Delta a)^2\}$
and assms-lin-expect: $\forall a \in A. expectation (\lambda\omega. \sum t = k..<n. (if \pi (t+1) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega) then 1 else 0)) = (\sum t = k..<n. expectation (\lambda\omega. indicat-real \{\omega. \pi (t+1) \omega = a \wedge 2 * \varepsilon * \ln (real t) / (\Delta a)^2 \leq real (N-n t a \omega)\} \omega))$
and mono-sum-sets:
 $(\forall a \in A. \Delta a * expectation (\lambda\omega. real (N-n n a \omega)) \leq \Delta a * (s n + (\sum t = k..<n. 2 / (real t powr \varepsilon)))) \implies (\sum a \in A. \Delta a * expectation (\lambda\omega. real (N-n n a \omega))) \leq (\sum a \in A. \Delta a * (s n + (\sum t = k..<n. 2 / (real t powr \varepsilon))))$

shows $expectation (\lambda\omega. R-n n \omega) \leq (\sum a \in A. \Delta a * ((2 * \varepsilon * \ln (real n)) / ((\Delta a)^2) + (\sum t = k..<n. 2 / (real t powr \varepsilon))))$
{proof}

end
end

References

- [1] R. Rebeschini. Lecture 15: Stochastic multi-armed bandit problem and algorithms. 2022. Available at <https://www.stats.ox.ac.uk/~rebeschini/teaching/AFoL/22/material/lecture15.pdf>.