



Mathematical
Institute

Multi-armed bandits under uncertainty aversion

TANUT TREETATANTHIPLOET

*Mathematical Institute
University of Oxford*

Based on joint work with Prof. Samuel N. Cohen

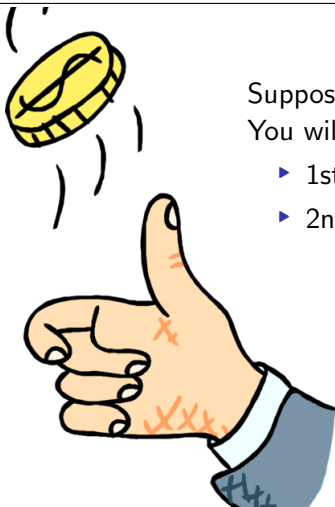
12th European Summer School in Financial Mathematics
Padova, September 2019

Oxford
Mathematics



Toy Example

Gambling is a way of buying hope on credit. – Alan Wykes



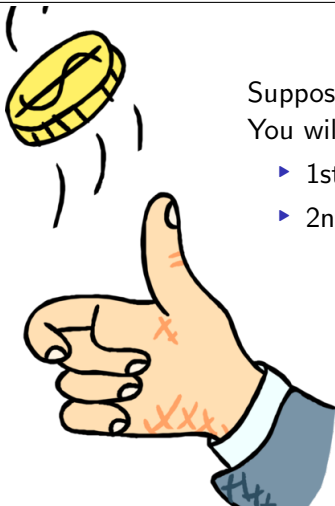
Suppose there are two biased coins.

You will gain £1 for a head

- ▶ 1st coin: 3 tosses with 2 heads.
- ▶ 2nd coin: 3000 tosses with 2000 heads.

Toy Example

Gambling is a way of buying hope on credit. – Alan Wykes



Suppose there are two biased coins.

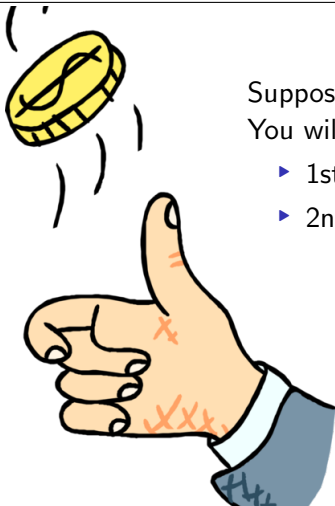
You will gain £1 for a head

- ▶ 1st coin: 3 tosses with 2 heads.
- ▶ 2nd coin: 3000 tosses with 2000 heads.

- ▶ We are biased toward a less uncertain choice.

Toy Example

Gambling is a way of buying hope on credit. – Alan Wykes



Suppose there are two biased coins.

You will gain £1 for a head

- ▶ 1st coin: 3 tosses with 2 heads.
- ▶ 2nd coin: 3000 tosses with 2000 heads.

- ▶ We are biased toward a less uncertain choice.
- ▶ What shall we do if we need to repeat for a million tosses?

Multi-armed bandits problem

Gambling is a way of buying hope on credit. – Alan Wykes



- Suppose there are M slot machines.
- One machine can be played at a time.
- Each machine may have its own state.
- The machines are independent.
- Playing a machine generates a cost and its state may evolve.

Multi-armed bandits problem

Gambling is a way of buying hope on credit. – Alan Wykes



- Suppose there are M slot machines.
- One machine can be played at a time.
- Each machine may have its own state.
- The machines are independent.
- Playing a machine generates a cost and its state may evolve.

	Distribution	Evolving state
Risky bandit	Known	Yes
Stationary bandit	Unknown	No
Non-stationary bandit	Unknown	Yes

Gittins' index theorem as a risk model

An optimist is a guy that has never had much experience. –Don Marquis

Consider a **risky bandit problem** with known distribution.

- ▶ Costs $(h^{(m)}(t))$ are not IID.
- ▶ Objective: Minimise $\mathbb{E}\left(\sum_{n=0}^{\infty} \beta^n h^{(\rho_n)}(t_n^\rho)\right)$.

There exist **indices** associated to each machine which can be **evaluated independently** such that the optimal policy is to play at each epoch a machine of **the lowest index** $\gamma^{(m)}$. (Gittins, 1979)

Gittins' index theorem as a risk model

An optimist is a guy that has never had much experience. –Don Marquis

Consider a **risky bandit problem** with known distribution.

- ▶ Costs $(h^{(m)}(t))$ are not IID.
- ▶ Objective: Minimise $\mathbb{E}\left(\sum_{n=0}^{\infty} \beta^n h^{(\rho_n)}(t_n^\rho)\right)$.

There exist **indices** associated to each machine which can be **evaluated independently** such that the optimal policy is to play at each epoch a machine of **the lowest index** $\gamma^{(m)}$. (Gittins, 1979)

By modeling costs $h^{(m)}(t)$ under a Bayesian perspective, we can show that

$$\gamma^{(m)} \approx \bar{x}^{(m)} - \sigma^{(m)} \psi\left(\frac{1}{(n^{(m)} + 1)(1 - \beta)}\right) \quad (\text{Brezzi and Lai, 2002})$$

where $\bar{x}^{(m)}$ and $\sigma^{(m)}$ are posterior mean and s.d. and ψ is positive and nondecreasing.

Optimistic Analogy for Reinforcement learning

An optimist is a guy that has never had much experience. –Don Marquis

For **stationary bandit problem**,

- ▶ The m^{th} machine generates IID costs $\left(h^{(m)}(t)\right)_{t \in \mathbb{N}}$ with mean $\mu^{(m)}$.

$$\rho^* = \arg \min_m \left(\text{Est. of } \mu^{(m)} - \underbrace{\text{Learning Premium}}_{\text{decreases in } n^{(m)}} \right).$$

i.e. We have less learning reward when we are more certain about our estimator.

Uncertainty Aversion

If you expect the worst, you'll never be disappointed – Sarah Dessen

Suppose we will only play once.

- ▶ This is equivalent to setting $\beta = 0$.
- ▶ Gittins' objective: Minimise $\mathbb{E}\left(h^{(\rho_0)}(t_0^\rho)\right) = \bar{x}^{(\rho_0)}$.

Uncertainty Aversion

If you expect the worst, you'll never be disappointed – Sarah Dessen

Suppose we will only play once.

- ▶ This is equivalent to setting $\beta = 0$.
- ▶ Gittins' objective: Minimise $\mathbb{E}\left(h^{(\rho_0)}(t_0^\rho)\right) = \bar{x}^{(\rho_0)}$.

No accounting for uncertainty!

Uncertainty Aversion

If you expect the worst, you'll never be disappointed – Sarah Dessen

Suppose we will only play once.

- ▶ This is equivalent to setting $\beta = 0$.
- ▶ Gittins' objective: Minimise $\mathbb{E}\left(h^{(\rho_0)}(t_0^\rho)\right) = \bar{x}^{(\rho_0)}$.

No accounting for uncertainty!

Including uncertainty aversion, we expect to have

$$\gamma^{(m)} = \text{Est. of } \mu^{(m)} + \underbrace{\left(-\text{Learning premium} + \text{Uncertainty aversion}\right)}_{\text{Uncertainty valuation}}.$$

Time-consistent Nonlinear Expectation

If you expect the worst, you'll never be disappointed. – Sarah Dessen

Classical control problem under uncertainty aversion:

$$\inf_{\rho} \sup_{\mathbb{Q}} \mathbb{E}^{\mathbb{Q}} \left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^{\rho}) \right) =: \inf_{\rho} \mathcal{E} \left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^{\rho}) \right).$$

Time-consistent Nonlinear Expectation

If you expect the worst, you'll never be disappointed. – Sarah Dessen

Classical control problem under uncertainty aversion:

$$\inf_{\rho} \sup_{\mathbb{Q}} \mathbb{E}^{\mathbb{Q}} \left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^{\rho}) \right) =: \inf_{\rho} \mathcal{E} \left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^{\rho}) \right).$$

We say a system of operator $\mathcal{E}(\cdot | \mathcal{F}_t) : L^{\infty}(\mathbb{P}, \mathcal{F}_T) \rightarrow L^{\infty}(\mathbb{P}, \mathcal{F}_t) : t = 0, 1, \dots, T$ is an (\mathcal{F}_t) -consistent coherent nonlinear expectation if it satisfies strict monotonicity, positive homogeneity, subadditivity and Lebesgue property (lower semi-continuity) and

- (\mathcal{F}_t) -consistency: for $t \leq t' \leq T$,

$$\mathcal{E} \left(\mathcal{E}(X | \mathcal{F}_{t'}) | \mathcal{F}_t \right) = \mathcal{E}(X | \mathcal{F}_t).$$

Time-consistent Nonlinear Expectation

If you expect the worst, you'll never be disappointed. – Sarah Dessen

Classical control problem under uncertainty aversion:

$$\inf_{\rho} \sup_{\mathbb{Q}} \mathbb{E}^{\mathbb{Q}} \left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^{\rho}) \right) =: \inf_{\rho} \mathcal{E} \left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^{\rho}) \right).$$

We say a system of operator $\mathcal{E}(\cdot | \mathcal{F}_t) : L^{\infty}(\mathbb{P}, \mathcal{F}_T) \rightarrow L^{\infty}(\mathbb{P}, \mathcal{F}_t) : t = 0, 1, \dots, T$ is an (\mathcal{F}_t) -consistent coherent nonlinear expectation if it satisfies strict monotonicity, positive homogeneity, subadditivity and Lebesgue property (lower semi-continuity) and

- (\mathcal{F}_t) -consistency: for $t \leq t' \leq T$,

$$\mathcal{E}(\mathcal{E}(X | \mathcal{F}_{t'}) | \mathcal{F}_t) = \mathcal{E}(X | \mathcal{F}_t).$$

The filtration (\mathcal{F}_t) must be identified in advance.

Our information structures

If you expect the worst, you'll never be disappointed. – Sarah Dessen

We have M bandits, each with a filtered space $(\Omega^{(m)}, (\mathcal{F}_t^{(m)}), \mathbb{P}^{(m)})$ and a consistent coherent nonlinear expectation

$$\mathcal{E}^{(m)}(X|\mathcal{F}_t^{(m)}) = \operatorname{ess\,sup}_{\mathbb{Q} \in \mathcal{Q}^{(m)}} \mathbb{E}^{\mathbb{Q}}(X|\mathcal{F}_t^{(m)})$$

(See Follmer & Schied, 2016).

- ▶ Define the *orthant space* by $\bar{\Omega} = \bigotimes_m \Omega^{(m)}$, similarly $\bar{\mathbb{P}}$, and

$$\bar{\mathcal{F}}(\underline{s}) = \bigotimes_m \mathcal{F}^{(m)}(s^{(m)}) : \underline{s} = (s^{(1)}, s^{(2)}, \dots, s^{(m)}).$$

- ▶ Define the orthant nonlinear expectation by

$$\mathfrak{E}_s(Y) = \operatorname{ess\,sup}_{\mathbb{Q} \in \bar{\mathcal{Q}}} \mathbb{E}^{\mathbb{Q}}(Y|\bar{\mathcal{F}}(\underline{s})) : \bar{\mathcal{Q}} = \{\mathbb{Q} = \bigotimes_m \mathbb{Q}^{(m)} \text{ for } \mathbb{Q}^{(m)} \in \mathcal{Q}^{(m)}\}.$$

Time consistency

When the Facts Change, I Change My Mind. What Do You Do, Sir?—Keynes

One may want to minimise

$$\mathfrak{E}_0\left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)\right) =: \mathfrak{E}_0(H^\rho)$$

but \mathfrak{E}_s does not satisfy time-consistency.

Time consistency

When the Facts Change, I Change My Mind. What Do You Do, Sir?—Keynes

One may want to minimise

$$\mathfrak{E}_0\left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)\right) =: \mathfrak{E}_0(H^\rho)$$

but \mathfrak{E}_s does not satisfy time-consistency.

- ▶ No DPP \Rightarrow Curse of dimensionality.

Time consistency

When the Facts Change, I Change My Mind. What Do You Do, Sir?—Keynes

One may want to minimise

$$\mathfrak{E}_0\left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)\right) =: \mathfrak{E}_0(H^\rho)$$

but \mathfrak{E}_s does not satisfy time-consistency.

- ▶ No DPP \Rightarrow Curse of dimensionality.
- ▶ Inconsistency in decision making. i.e. ‘Optimal’ strategies may not be followed in the future.

Time consistency

When the Facts Change, I Change My Mind. What Do You Do, Sir?–Keynes

One may want to minimise

$$\mathfrak{E}_0\left(\sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)\right) =: \mathfrak{E}_0(H^\rho)$$

but \mathfrak{E}_s does not satisfy time-consistency.

- ▶ No DPP \Rightarrow Curse of dimensionality.
- ▶ Inconsistency in decision making. i.e. ‘Optimal’ strategies may not be followed in the future.
- ▶ Using an indifference valuation perspective can be helpful.

'Gittins' optimality'

When the Facts Change, I Change My Mind. What Do You Do, Sir?—Keynes

- ▶ $\mathfrak{E}_0(H^\rho - \mathfrak{E}_0(H^\rho)) = 0$.
- ▶ $\min_\rho \mathfrak{E}_0(H^\rho) := C^{\rho^*} \leq C^\rho$ where $C^\rho \in \mathbb{R}$ and $\mathfrak{E}_0(H^\rho - C^\rho) = 0$.

'Gittins' optimality'

When the Facts Change, I Change My Mind. What Do You Do, Sir?—Keynes

- ▶ $\mathfrak{E}_0(H^\rho - \mathfrak{E}_0(H^\rho)) = 0.$
- ▶ $\min_\rho \mathfrak{E}_0(H^\rho) := C^{\rho^*} \leq C^\rho$ where $C^\rho \in \mathbb{R}$ and $\mathfrak{E}_0(H^\rho - C^\rho) = 0.$

Definition

We say Y^ρ is a compensator of a cost H^ρ (under strategy ρ) if

$$\mathfrak{E}_0(H^\rho - Y^\rho) = 0.$$

We say ρ^* is a Gittins' optimum if there exists a compensator family $\{Y^\rho\}_\rho$ such that $Y^{\rho^*} \leq Y^\rho.$

- ▶ In a dynamic setting, we require Y^ρ to be predictable and supercompensate at all times.

The main theorem

There is no more miserable human being than one in whom nothing is habitual but indecision. —W. James

For a machine m , we have $(\mathcal{F}_t^{(m)})_{t \geq 0}$ -adapted random costs $h^{(m)}(t)$.

Theorem

A *Gittins' optimal* strategy for the cost $H^\rho := \sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)$ can be given by always playing a machine with the *minimal index* $\gamma^{(m)}(s)$ where

$$\gamma^{(m)}(s) := \text{ess inf} \left\{ \gamma : \text{ess inf}_{\tau \in \mathcal{T}^{(m)}(s)} \mathcal{E}^{(m)} \left(\sum_{t=1}^{\tau} \beta^t (h^{(m)}(s+t) - \gamma) \mid \mathcal{F}_s^{(m)} \right) \leq 0 \right\}.$$

The main theorem

There is no more miserable human being than one in whom nothing is habitual but indecision. —W. James

For a machine m , we have $(\mathcal{F}_t^{(m)})_{t \geq 0}$ -adapted random costs $h^{(m)}(t)$.

Theorem

A *Gittins' optimal* strategy for the cost $H^\rho := \sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)$ can be given by always playing a machine with the *minimal index* $\gamma^{(m)}(s)$ where

$$\gamma^{(m)}(s) := \text{ess inf} \left\{ \gamma : \text{ess inf}_{\tau \in \mathcal{T}^{(m)}(s)} \mathcal{E}^{(m)} \left(\sum_{t=1}^{\tau} \beta^t (h^{(m)}(s+t) - \gamma) \mid \mathcal{F}_s^{(m)} \right) \leq 0 \right\}.$$

- ▶ $\gamma^{(m)}$ can be found by solving a nonlinear optimal stopping problem independently for each machine.

The main theorem

There is no more miserable human being than one in whom nothing is habitual but indecision. —W. James

For a machine m , we have $(\mathcal{F}_t^{(m)})_{t \geq 0}$ -adapted random costs $h^{(m)}(t)$.

Theorem

A *Gittins' optimal* strategy for the cost $H^\rho := \sum_{n=1}^L \beta^n h^{(\rho_n)}(t_n^\rho)$ can be given by always playing a machine with the *minimal index* $\gamma^{(m)}(s)$ where

$$\gamma^{(m)}(s) := \text{ess inf} \left\{ \gamma : \text{ess inf}_{\tau \in \mathcal{T}^{(m)}(s)} \mathcal{E}^{(m)} \left(\sum_{t=1}^{\tau} \beta^t (h^{(m)}(s+t) - \gamma) \mid \mathcal{F}_s^{(m)} \right) \leq 0 \right\}.$$

- ▶ $\gamma^{(m)}$ can be found by solving a nonlinear optimal stopping problem independently for each machine.
- ▶ Playing based on indices yields a consistent decision. (i.e. an optimal decision follows through.)

Bernoulli Bandit

But to us, probability is the very guide of life. —Joseph Butler

- ▶ Bernoulli bandit: $h(t) \sim B(1, \theta) : \theta$ is unknown.
- ▶ We model the uncertainty by

$$\mathcal{E}_{(t)}^k(\cdot) := \operatorname{ess\,sup}_{\theta \in \Theta_t^k} \mathbb{E}^\theta(\cdot)$$

where Θ_t^k is a posterior credible set of size k (under an improper prior).

- ▶ We extend it by $\mathcal{E}^k(\cdot | \mathcal{F}_t) := \mathcal{E}_{(t)}^k(\dots \mathcal{E}_{(T-1)}^k(\cdot) \dots)$.

Bernoulli Bandit

But to us, probability is the very guide of life. —Joseph Butler

- ▶ Bernoulli bandit: $h(t) \sim B(1, \theta) : \theta$ is unknown.
- ▶ We model the uncertainty by

$$\mathcal{E}_{(t)}^k(\cdot) := \operatorname{ess\,sup}_{\theta \in \Theta_t^k} \mathbb{E}^\theta(\cdot)$$

where Θ_t^k is a posterior credible set of size k (under an improper prior).

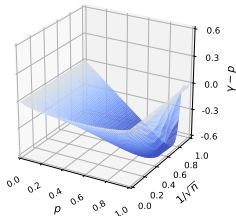
- ▶ We extend it by $\mathcal{E}^k(\cdot | \mathcal{F}_t) := \mathcal{E}_{(t)}^k(\dots \mathcal{E}_{(T-1)}^k(\cdot) \dots)$.
- ▶ It can be shown that

$$\gamma(t) = \gamma_{k,\beta}\left(p_t, \frac{1}{\sqrt{n_t}}\right) =: p_t + \text{Uncertainty valuation.}$$

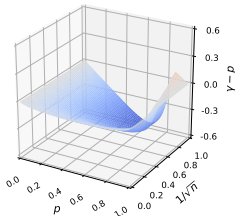
Difference between γ and p (Uncertainty valuation)

But to us, probability is the very guide of life. —Joseph Butler

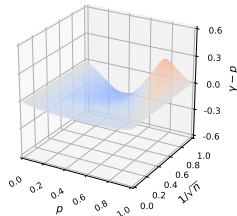
$\beta = 0.9999, k = 0.01$



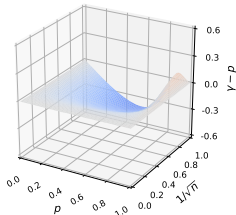
$\beta = 0.9999, k = 0.5$



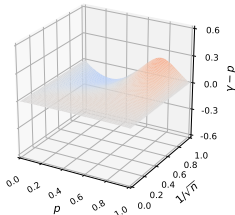
$\beta = 0.9999, k = 0.8$



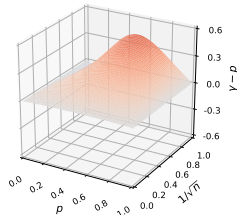
$\beta = 0.95, k = 0.01$



$\beta = 0.95, k = 0.5$



$\beta = 0.95, k = 0.8$



Monte-Carlo Simulation

Facts are stubborn things, but statistics are pliable.—Mark Twain

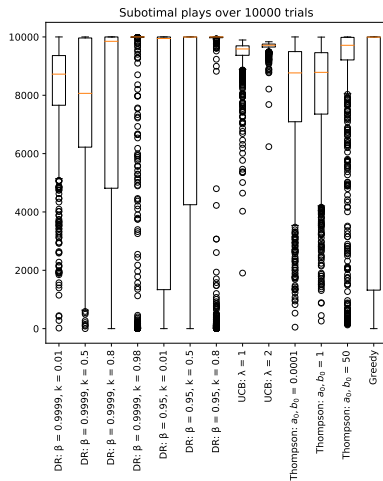
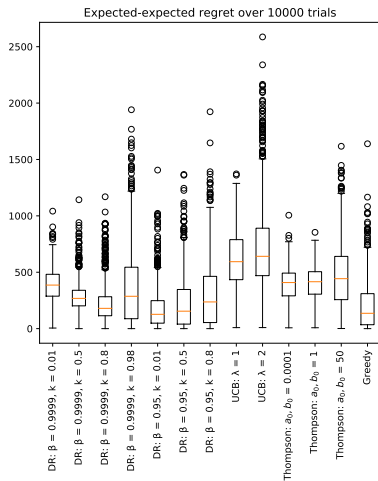
- ▶ Randomly choose a and b independently from $\Gamma(1, 1/100)$.
- ▶ Take 50 samples from $\text{Beta}(a, b)$ and use them as true probabilities of the 50 Bernoulli bandits.
- ▶ Evaluate algorithms over 10000 trials/simulation.
- ▶ Start with initial information of size 10 from each bandit.

We will consider the performance of each algorithm by considering

- ▶ Expected-regret: $R(L) := \sum_{n=1}^L (\theta^{(\rho_n)} - \theta^*)$.
- ▶ Suboptimal plays: $N_V(L) := \sum_{n=1}^L \mathbb{I}(\theta^{(\rho_n)} \neq \theta^*)$.

Performance

Facts are stubborn things, but statistics are pliable.—Mark Twain



Conclusion

The finest studies are like the finest of anything else: They cost big bucks. —Charles Wheelan

Contribution

- ▶ We propose an alternative optimality criterion to address consistent decision making under uncertainty over multiple filtrations.
- ▶ We derive an index which only involves a one dimensional (time-consistent) robust problem which is computationally tractable.
- ▶ Our model takes into account the desire to learn and uncertainty aversion.

Reference

- ▶ S.N. Cohen and T. Treetanthiploet, Gittins' theorem under uncertainty, arXiv:1907.05689.

Q & A

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ Observe that $\gamma^{(m)}(t)$ is the minimum compensated reward which could encourage us to pay the cost from time t until the optimal stopping time τ^* .
- ▶ By the minimality of $\gamma^{(m)}(t)$, we must have zero total return under the optimal stopping, i.e.

$$\mathcal{E}^{(m)}\left(\sum_{s=t+1}^{\tau^*} \beta^s (h^{(m)}(s) - \gamma^{(m)}(t)) \middle| \mathcal{F}_t^{(m)}\right) = 0$$

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ Observe that $\gamma^{(m)}(t)$ is the minimum compensated reward which could encourage us to pay the cost from time t until the optimal stopping time τ^* .
- ▶ By the minimality of $\gamma^{(m)}(t)$, we must have zero total return under the optimal stopping, i.e.

$$\mathcal{E}^{(m)}\left(\sum_{s=t+1}^{\tau^*} \beta^s (h^{(m)}(s) - \gamma^{(m)}(t)) \middle| \mathcal{F}_t^{(m)}\right) = 0$$

- ▶ In particular, at time τ^* , we need to increase the compensated reward to encourage further play.

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ We increase the reward minimality by considering the reward process

$$\Gamma^{(m)}(t) = \max_{0 \leq \theta \leq t-1} \gamma^{(m)}(\theta).$$

- ▶ This minimal reward encourage us to pay the random cost $h^{(m)}(t)$ until the horizon $T^{(m)}$ with the return 0.

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ We increase the reward minimality by considering the reward process

$$\Gamma^{(m)}(t) = \max_{0 \leq \theta \leq t-1} \gamma^{(m)}(\theta).$$

- ▶ This minimal reward encourage us to pay the random cost $h^{(m)}(t)$ until the horizon $T^{(m)}$ with the return 0.
- ▶ Since this reward encourage us to continue at any point in time, for every $t \geq 0$,

$$\mathcal{E}^{(m)}\left(\sum_{s=t+1}^{T^{(m)}} \beta^s (h^{(m)}(s) - \Gamma^{(m)}(s)) \middle| \mathcal{F}_t^{(m)}\right) \leq 0$$

with equality at $t = 0$.

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ Now imagine that we are taking a break from the m -th arm to play another arm. Once we return to play the m -th arm we face a further rescaling of the discount factor. In particular, we have the discount factor $\alpha(s)\beta^s$ instead of β^s where α is decreasing.

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ Now imagine that we are taking a break from the m -th arm to play another arm. Once we return to play the m -th arm we face a further rescaling of the discount factor. In particular, we have the discount factor $\alpha(s)\beta^s$ instead of β^s where α is decreasing.
- ▶ By the robust representation theorem, for any $\epsilon > 0$, show that there exists $\mathbb{Q} \in \mathcal{Q}^{(m)}$ such that

$$\mathbb{E}^{\mathbb{Q}} \left[\sum_{s=1}^{T^{(m)}} \alpha(s)\beta^s (h^{(m)}(s) - \Gamma^{(m)}(s)) \right] \geq -\epsilon$$

for every predictable decreasing process α in $[0, 1]$.

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ Now imagine that we are taking a break from the m -th arm to play another arm. Once we return to play the m -th arm we face a further rescaling of the discount factor. In particular, we have the discount factor $\alpha(s)\beta^s$ instead of β^s where α is decreasing.
- ▶ By the robust representation theorem, for any $\epsilon > 0$, show that there exists $\mathbb{Q} \in \mathcal{Q}^{(m)}$ such that

$$\mathbb{E}^{\mathbb{Q}} \left[\sum_{s=1}^{T^{(m)}} \alpha(s) \beta^s (h^{(m)}(s) - \Gamma^{(m)}(s)) \right] \geq -\epsilon$$

for every predictable decreasing process α in $[0, 1]$.

- ▶ In particular,

$$\mathbb{E} \left(\sum_n \beta^n (h^{(\rho_n)}(t_n^\rho) - \Gamma^{(\rho_n)}(t_n^\rho)) \right) \geq 0.$$

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ If we take a break (i.e. leave the arm) only when

$$\mathcal{E}^{(m)}\left(\sum_{s=t+1}^{T^{(m)}} \beta^s (h^{(m)}(s) - \Gamma^{(m)}(s)) \middle| \mathcal{F}_t^{(m)}\right) = 0$$

i.e. when $\gamma^{(m)}$ reaches a new maximum, the total cost must be zero.

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ If we take a break (i.e. leave the arm) only when

$$\mathcal{E}^{(m)} \left(\sum_{s=t+1}^{T^{(m)}} \beta^s (h^{(m)}(s) - \Gamma^{(m)}(s)) \middle| \mathcal{F}_t^{(m)} \right) = 0$$

i.e. when $\gamma^{(m)}$ reaches a new maximum, the total cost must be zero.

- ▶ In particular,

$$\mathbb{E} \left(\sum_n \beta^n (h^{(\rho_n^*)}(t_n^*) - \Gamma^{(\rho_n^*)}(t_n^*)) \right) = 0.$$

Sketch idea of the proof

Everything should be made as simple as possible, but not simpler—Albert Einstein

- ▶ If we take a break (i.e. leave the arm) only when

$$\mathcal{E}^{(m)}\left(\sum_{s=t+1}^{T^{(m)}} \beta^s (h^{(m)}(s) - \Gamma^{(m)}(s)) \middle| \mathcal{F}_t^{(m)}\right) = 0$$

i.e. when $\gamma^{(m)}$ reaches a new maximum, the total cost must be zero.

- ▶ In particular,

$$\mathfrak{E}\left(\sum_n \beta^n (h^{(\rho_n^*)}(t_n^*) - \Gamma^{(\rho_n^*)}(t_n^*))\right) = 0.$$

- ▶ Finally, as $t \mapsto \Gamma^{(m)}(t)$ is increasing, ρ^* minimise $\sum_{n=1}^N \beta^n \Gamma^{(\rho_n^*)}(t_n^*)$ for all N .