# Gittins Index Tutorial

## Overview

What is the Gittins index?

Solution to the Bayesian/Markovian multi-armed bandit problem

How should one understand the Gittins index?      What other problems does the Gittins index solve?

∝ super slick proof      →      combinatorial optimization

ΔV dynamic programming      →      "superprocesses"

6o⌐ queueing theory glasses      →      arrivals, long-run avg cost

First: golf!

## Intro

Markov-process multi-golf



Multi-golf:
- Every turn, hit one ball
- When hit, ball moves to random neighbor
- Goal: minimize $\underbrace{\text{time (\# hits)}}_{\Phi^\pi}$ to get either ball • to flag

  ↳ minimize $E[\Phi^\pi]$
  policy $\pi$

**Question:** Which ball to hit first?

Use expected # hits to score one ball?

Let $S_x$ = time from • at $x$ to ⚑ (if not interrupted)
    ("swings")

$$E[S_{2,1}] = 1 + \tfrac{1}{2} E[S_{2,2}] = 3$$
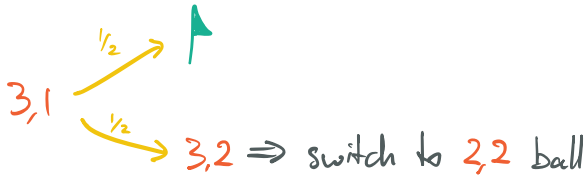$$E[S_{2,2}] = 1 + E[S_{2,1}] = 4$$

$$E[S_{3,1}] = 5$$
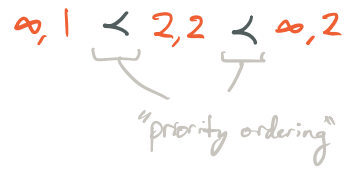$$E[S_{3,2}] = 8$$
$$E[S_{3,3}] = 9$$

$$E[S_{2,2}] < E[S_{3,1}]$$

... but should hit 3,1 first!

Why 3,1 before 2,2?

3,1 →$\frac{1}{2}$→ 🚩
3,1 →$\frac{1}{2}$→ 3,2 ⟹ switch to 2,2 ball
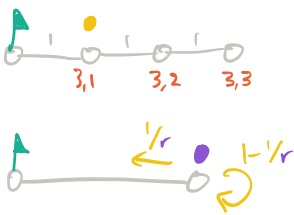
Extreme case: ∞,1 vs 2,2

∞,1 < 2,2 < ∞,2

"priority ordering"

# Takeaway:

- $E[S_x]$ is relevant if we would hit 🟡 from $x$ to 🚩 without interruption
- But we might interrupt if another 🟡 is better
⟹ $E[S_x]$ is not enough!

Generalization: any Markov chain/process, any cost/time per edge

# Key idea: Calibrate each state $x$ to a simple benchmark

🚩 · 🟡 · ∘ · ∘
  3,1  3,2  3,3

🚩 ←$\frac{1}{r}$— 🟣 —$1-\frac{1}{r}$→ ↺

Pay-to-win golf
- Every turn, hit ball 🟡 (just one)
  OR give up by paying $r$
- Goal: minimize total cost to get a 🚩:

$\underset{\text{policy } \pi}{\text{minimize}} \left( E[\text{cost from hitting}] + E[\text{cost from giving up}] \right)$

$= \underset{\text{state sets } Y}{\text{minimize}} \left( \underbrace{E[\text{time to exit } Y]}_{= S_x(Y)} + r \underbrace{P[\text{when out of } Y, \text{ not at } 🚩]}_{= h_x(Y)} \right)$

(starting from $x$)     (starting from $x$)

$E[S_x(Y)] = s_x(Y)$     "had enough"

# Def: Pay-to-win golf cost fns:

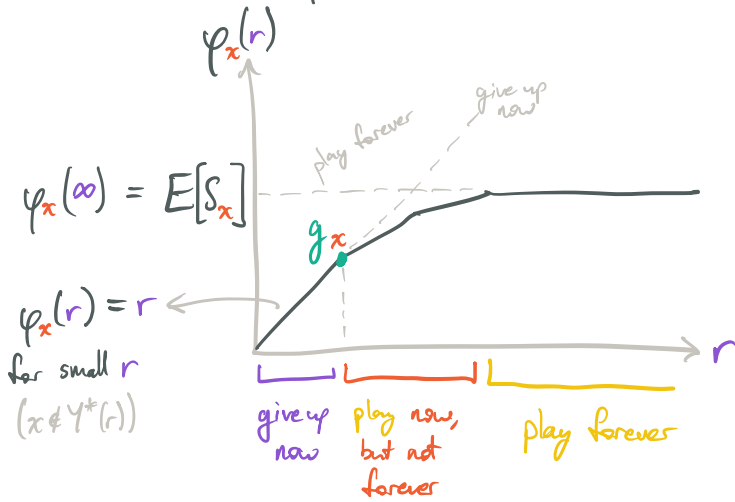$\varphi_x(r, Y) = s_x(Y) + r h_x(Y)$

$Y^*(r) = \underset{Y}{\text{argmax }} \varphi_x(r, Y)$

$\varphi_x(r) = \varphi_x(r, Y^*(r))$

$= \underbrace{s_x(Y^*(r))}_{} + r \underbrace{h_x(Y^*(r))}_{}$

$= s_x(r) + r h_x(r)$

For short → $\varphi = s + rh$

How does pay-to-win golf go?
- Play while in $Y^*(r)$
- Give up if we exit $Y^*(r)$

What does $\varphi_x(r)$ look like? How can we use it to "rate" state $x$?

$$\varphi_x(r)$$

play forever

give up now

$$\varphi_x(\infty) = E[S_x]$$

$g_x$

$$\varphi_x(r) = r$$
for small $r$
$(x \notin Y^*(r))$

give up now | play now, but not forever | play forever

$r$

**Def:** The **Gittins rank** of state $x$ is:

$$g_x = \inf \{ r \geq 0 \mid \varphi_x(r) < r \}$$

$$= \inf \{ r \geq 0 \mid x \in Y^*(r) \}$$

$$= \inf \{ r \geq 0 \mid \text{play at least once} \}$$

(The **Gittins index** is $\frac{1}{g_x}$)

**Exercise:** Find $g_{2,1}, \; g_{2,2}, \; g_{3,1}, \; g_{3,2}, \; g_{3,3}, \; g_{\infty,1}, \; g_{\infty,2}, \; g_{\infty,k}, \; g_{n,k}$

Hint: In each case, what is $\underline{Y^*(g_x)? \; S_x(g_x)? \; h_x(g_x)?}$

$\hookrightarrow$ Include $x \in Y^*(g_x)$

**Lem:** $g_x = \inf\limits_{Y} \dfrac{S_x(Y)}{1 - h_x(Y)} = \dfrac{S_x(g_x)}{1 - h_x(g_x)}$   $\rightarrow$ short for $Y^*(g_x)$

**Pf:** $g_x < r \iff \exists Y, \; S_x(Y) + r \, h_x(Y) < r \iff \exists Y, \; \dfrac{S_x(Y)}{1 - h_x(Y)} < r$

$\hookrightarrow$ Can pick $Y^*(g_x)$ if $g_x < r$; optimal as $r \to g_x$

(Going to be less color-coded hereafter...)

# Gittins optimality proof 1

α   Super slick proof
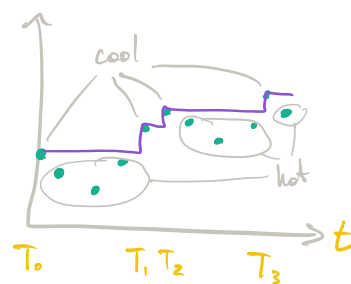
Gittins policy always plays ball at _minimal_ $g_x$

Optimal for one stochastic ball, one deterministic ball (pay-to-win golf)

Two stochastic balls?

**Key idea:** Reduce two-ball case to $\underset{1}{\overset{random}{\underline{sequence}}}$ of pay-to-win problems

**Def:** Consider trajectory $X_x(t)$ with $X_x(0) = x$

- The **Gittins envelope** is $G_x(t) = \underset{0 \le u \le t}{max} \, g_{X(u)}$

- The **Gittins maximum** is $G_x = G(\infty)$

- Ball is **cool** if $G(t) = g_{X(t)}$, else **hot**

(Often omit $x$ subscript)



Consider following sequence of times

- $T_0 = 0$, play through pay-to-win with $r = g_x = G(T_0)$

  We **give up** iff $G_x > r$, and if so...

- $T_1 = S_x(G(T_0))$ is when we give up

  $\vdots$

- $T_{i+1} = S_x(G(T_i)) = T_i + S_{X(T_i)}\left(g_{X(T_i)}\right)$ while $G_x > g_{X(T_i)}$

$T_i$ is $i^{th}$ time ball is cool

Consider $\varphi_x(g_x) = S_x(g_x) + g_x h_x(g_x) = g_x$

← by def. of $g_x$ (or lemma)

$$\Rightarrow g_x(1 - h_x(g_x)) = S_x(g_x)$$

$$\Rightarrow E\left[G_x \mathbb{1}(\text{get to } \Lambda)\right] = E\left[\text{time spent playing}\right]$$

Intuition: imagine instead of pay-to-win, we get a reward when done

More generally: $E\left[G_x \mathbb{1}(\text{get to } A \text{ in round } i) \mid T_i, X(T_i)\right] = E\left[G(T_i) \underbrace{\mathbb{1}(A_i)}_{=G_i} \mid T_i, X(T_i)\right]$

$\parallel$

$E\left[\text{time played in heat } i \mid \cdots\right] = E\left[\underbrace{T_{i+1} - T_i}_{=S_i} \mid T_i, \underbrace{X(T_i)}_{=X_i}\right]$

**Lem:** $E\left[G_i \mathbb{1}(A_i) \mid X_i, T_i\right] = E\left[S_i \mid X_i, T_i\right]$

**Def:** The **pseudotime** of heat $i$ is $G_i \mathbb{1}(A_i) = G_x \mathbb{1}(A_i)$

So $E\left[\text{pseudotime of heat}\right] = E\left[\text{time of heat}\right]$

But pseudotime is all at the end of heat, if at all, so

$E\left[\text{pseudotime of incomplete heat}\right] = 0 \leq E\left[\text{time of incomplete heat}\right]$

**Lem:** For any stopping time $\tau$, $E\left[\text{pseudotime by } \tau\right] \leq E\left[\text{time by } \tau\right] = E[\tau]$
   Becomes = if always cool at $\tau$

All of this is still for one ball... but still true for two balls interleaved!

**Lem:** For any policy $\pi$, $E[\underbrace{\text{pseudotime of } \pi}_{=\Psi^\pi}] \leq E[\underbrace{\text{time of } \pi}_{=\Phi^\pi}]$
   Becomes = if $\pi$ only interrupts a ball when it's cool

We have:
$$E[\Psi^\pi] \leq E[\Phi^\pi]$$
$$E[\Psi^{\text{Gittins}}] = E[\Phi^{\text{Gittins}}]$$
*Same reasoning for any # of jobs*

Suffices:
$$E[\Psi^{\text{Gittins}}] \leq E[\Psi^\pi]$$

**Question:** Is this true? Why?
   ◦ Yes!
   ◦ Because $\Psi^\pi \geq \min\{G_{x_1}, G_{x_2}\} = \Psi^{\text{Gittins}}$

**Thm:** Gittins policy is optimal, w/
$$E\left[\Phi_{\vec{x}}^{\text{Gittins}}\right] = E\left[\min_j G_{x_j}\right]$$

**Extension:** First $k$ balls, set of balls is matroid base, ... [Singh '18, GJSS '19]

# Gittins optimality proof 2

## ΔV Dynamic programming

Let $\Phi(\vec{x}) = E\left[\Phi_{\vec{\pi}}^{Gittins}\right] = E\left[\min_j G_{x_j}\right]$

$\Phi$ is cost-to-go fn of multi-golf

Can we show $\Phi$ satisfies Bellman eq?

Two steps:
- Write Bellman eq
- Think about $G_x$ to understand $\Phi$

## Exercise: What is $P\left[G_x > r\right]$?

$G_x > r \iff$ give up in pay-to-win from $x$ w/ cost $r$

$\Rightarrow P\left[G_x > r\right] = P\left[\text{give up}\right] = h_x(r)$

*or time to give up*

## Exercise: What is $E\left[\min\{G_x, r\}\right]$?

$\min\{G_x, r\}$ is pseudotime of pay-to-win golf, playing optimally

$\underbrace{\qquad}_{G_{give-up} = r}$

So $E\left[\text{time}\right] = E\left[\text{pseudotime}\right] = \varphi_x(r)$

$\Rightarrow E\left[\min\{G_x, r\}\right] = \varphi_x(r) = s_x(r) + r h_x(r)$

## Exercise: Write $E\left[\min\{G_{x_1}, G_{x_2}\}\right]$ in two ways: using $h$ and using $\varphi$, $G$

Hint: $E[R] = \int_0^\infty P[R > r] \, dr$ for non-negative rv R

$= \int_0^\infty h_{x_1}(r) h_{x_2}(r) \, dr = E\left[\varphi_{x_1}(G_{x_2})\right] = E\left[\varphi_{x_2}(G_{x_1})\right]$

Generalizing to any # of balls:

$$\Phi(\vec{x}) = \int_0^\infty \prod_j h_{x_j}(r)\, dr = E\left[\varphi_{x_i}\left(\min_{j \neq i} G_{x_j}\right)\right]$$

Time to show $\Phi$ solves Bellman eq

# Def:

- $\underbrace{\Delta_i f(x_i)}_{f(\ldots, x_i, \ldots)} = E[f(X_{next}) \,|\, \text{start at } x_i] - f(x_i)$

- $\Gamma f(r) = -\frac{1}{r} f(r) = \frac{1}{r} 0 + \frac{r-1}{r} f(r) - f(r)$

(changing notation: $\varphi_x(\cdot) \leadsto \varphi(x,r)$)

Bellman eq for pay-to-win golf:

$$\underbrace{\text{"} \min\{\Delta_i, \Gamma\} \varphi(x_i, r) \text{"}}_{\text{change in cost-to-go}} = \underbrace{-1}_{\substack{\text{cost paid} \\ \text{this round}}}$$

Specifically:

$$\Delta_i \varphi(x_i, r) \geq -1,$$
with $=$ iff $x_i \in Y^*(\cdot) \Leftrightarrow g_{x_i} \leq r$

Bellman eq for multi-golf:

$$\min_i \Delta_i \Phi(\vec{x}) = -1$$

Want to show

$$\Delta_i \Phi(\vec{x}) \geq -1,$$
with $=$ iff $g_{x_i} \leq g_{x_j} \quad \forall j \neq i$

$$\Delta_i \Phi(\vec{x}) = E\left[\Delta_i \varphi_{x_i}\left(\min_{j \neq i} G_{x_j}\right)\right] \geq -1$$

because $g_{x_j}$ is min possible value of $G_{x_j}$

$$\text{with } = \text{ iff } g_{x_i} \leq G_{x_j} \text{ w.p. } 1 \Leftrightarrow g_{x_i} \leq g_{x_j}$$

# Extension: Sometimes goes through when multiple actions possible per ball

("choice of club")

[Whittle '80, Doval '18]

# Gittins optimality proof 3

👓 Queueing theory glasses

Change of objective: minimize __sum__ of completion times ⟺ minimize avg # incomplete jobs

**Question:** Can we adapt previous proofs to this setting? Yes!
What if jobs arrive over time? Harder...

Recall: $\frac{\partial}{\partial r} \varphi_x(r) = h_x(r)$ ⟹ $\frac{\partial}{\partial r} \frac{\varphi_x(r)}{r} = \frac{-S_x(r)}{r^2}$

$$\Rightarrow \int_0^\infty \frac{S_x(r)}{r^2} dr = \frac{"\varphi_x(0)"}{0} - \frac{"\varphi_x(\infty)"}{\infty} = 1$$

$"d\frac{1}{r}"$

**Lem:** $\int_0^\infty E[S_x(r)] d\frac{1}{r} = 1$

Let $W_{\vec{x}}(r) = \sum_j S_{x_j}(r)$ be "r-work" in system

Lemma ⟹ $\int_0^\infty E[W_{\vec{x}}(r)] d\frac{1}{r} = \#\vec{x}$

Let $W(r)$ be r-work at "uniformly random" time

**Thm:** $\int_0^\infty E[W(r)] d\frac{1}{r} = E[\# \text{jobs}]$ ⟵ WINE: Work Integral Number Equality

**Question:** How to minimize $W(r)$?

**Thm:** Gittins minimizes $E[\# \text{jobs}]$

**Extension:** Still holds w/ Poisson arrivals
[SGH '20, SH '21]



W(r)

Gittins: best!

time



Not Gittins: not best