A Gittins Policy for Optimizing Tail Latency

Amit HarlevCornell CAM

Joint work with

George Yu Ziv Scully Cornell ORIE Cornell ORIE

1



We know how to minimize tail latency in queues!

Main takeaway





We know how to minimize tail latency in queues!

Main takeaway

for most information models





If you have a problem where this is needed, please come talk to us!

































queue





server

queue





server

queue





server

















queue





captures many information models in a single result



Examples of information models

Known Size (S = 6)





Examples of information models

Known Size (S = 6)

Unknown Size $(S \sim \text{Unif}\{2,4,6\})$







Examples of information models

Known Size (S = 6)

Unknown Size
($S \sim \text{Unif}\{2,4,6\}$)

Partial information (*S* ~ Unif{4,6})







queue

light-tailed M/G random → arrivals

Goal: scheduling policy π that minimizes *asymptotic behavior* of the *tail of response time*, $\mathbf{P}[T_{\pi} > x]$







queue

light-tailed M/G random → arrivals

Goal: scheduling policy π that minimizes *asymptotic behavior* of the *tail of response time*, $\mathbf{P}[T_{\pi} > x]$



 T_{π} = response time under policy π ("total time in system")













when *S* is light-tailed





$$\gamma_{\pi} = deconstraints$$
 $C_{\pi} = tail$

when *S* is light-tailed

ay rate of π

constant of π





$$\gamma_{\pi} = deconstraints$$
 $C_{\pi} = tail$

when *S* is light-tailed



threshold *x*

 $C_{\pi}e^{-\gamma_{\pi}x}$ (roughly)

ay rate of π

constant of π





 $\gamma_{\pi} = decay \ rate \ of \ \pi$ Weak optimality: optimal γ_{π} $C_{\pi} = tail constant of \pi$

when S is light-tailed



threshold *x*

 $C_{\pi}e^{-\gamma_{\pi}x}$ (roughly)





 $\gamma_{\pi} = decay \ rate \ of \ \pi$ Weak optimality: optimal γ_{π} $C_{\pi} = tail constant of \pi$

when *S* is light-tailed

Heavy-tailed S: PS & LAS have optimal asymptotic behavior Wierman & Zwart 2012, Yu & Scully 2024]

threshold *x*

asymptotic behavior

 $C_{\pi}e^{-\gamma_{\pi}x}$ (roughly)

Strong optimality: optimal γ_{π} and C_{π}







scheduling for response time: short jobs before long jobs



scheduling for response time: short jobs before long jobs

...leads to poor tail performance



scheduling for response time: short jobs before long jobs



...leads to poor tail performance

don't starve long jobs



Managing the tradeoff





Does any policy in the literature balance this?


Managing the tradeoff





Most policies: priority based on job's size or attained service

Does any policy in the literature balance this?



Managing the tradeoff





Most policies: priority based on job's size or attained service

Does any policy in the literature balance this?

not flexible enough!



Managing the tradeoff



Most policies: priority based on job's size or attained service

not flexible enough!



Existing work































FCFS weakly optimal [Boxma & Zwart]

strongly optimal!

 γ -Gittins [Harlev, Yu & Scully]



New policy: *γ*-Gittins!



Theorem

 γ -Gittins is strongly optimal, $C_{\gamma\text{-Gittins}} = \inf_{\pi} C_{\pi}$



New policy: *γ*-Gittins!





New policy: *γ*-Gittins!



Practice: what are the design takeaways?



New policy: *γ*-Gittins!



Practice: what are the design takeaways?

Theory: why was this hard to discover?



New policy: *γ*-Gittins!

Theory: why was this hard to discover?



Practice: what are the design takeaways?



What is *y*-Gittins?

11



11

What is γ-Gittins?







What is a boost policy?



What is γ-Gittins?





What is a boost policy?







What is a boost policy?

boost policy: serve jobs in order of **boosted** arrival time





boost policy: serve jobs in order of **boosted** arrival time



What is a boost policy?

boost(size) in [Yu & Scully, 2024]































0





































What does priority under γ-Gittins look like?



What does priority under γ-Gittins look like?

(Classical) Unknown Size





What does priority under γ -Gittins look like?

(Classical) Unknown Size







What does priority under γ -Gittins look like?

(Classical) Unknown Size







What does priority under γ -Gittins look like?

(Classical) Unknown Size






What does priority under γ -Gittins look like?

Partial Information S ∼ Unif{10, 60, 140} w/ size known after 10 units service



What does priority under γ -Gittins look like? 130 states 1/3 10 states 50 states

Partial Information

S ~ **Unif**{10, 60, 140}

w/ size known after 10 units service





Partial Information

S ~ **Unif**{10, 60, 140} w/ size known after 10 units service





S ~ **Unif**{10, 60, 140}



























If you want to minimize the tail of response time:



If you want to minimize the tail of response time:

1. priority should depends on both arrival time and job state



If you want to minimize the tail of response time:

- 1.
- 2.

priority should depends on both arrival time and job state

priority \approx big initial boost, then decrease based on state



If you want to minimize the tail of response time:

- priority should depends on both arrival time and job state 1.
- priority \approx big initial boost, then decrease based on state 2.
- priority can also increase based on state, but times where it 3. decreases are more important



If you want to minimize the tail of response time:

- priority should depends on both arrival time and job state 1.
- priority \approx big initial boost, then decrease based on state 2.
- priority can also increase based on state, but times where it 3. decreases are more important
- can use any information model, 4. but more information \implies better performance



New policy: *γ*-Gittins!



Practice: what are the design takeaways?

Theorem

 γ -Gittins is strongly optimal,

 $C_{\gamma ext{-Gittins}} = \inf_{\pi} C_{\pi}$



New policy: *γ*-Gittins!

Theorem γ -Gittins is strongly optimal, $C_{\gamma ext{-Gittins}} = \inf_{\pi} C_{\pi}$

Practice: what are the design takeaways?

Theory: why was this hard to discover?



Practice: what are the design takeaways?

Theory: why was this hard to discover?



Theorem

```
\gamma-Gittins is strongly optimal,
```

```
C_{\gamma-\text{Gittins}} = \inf C_{\pi}
```

boost = $\frac{1}{\pi} \log(\text{Gittins index})$







Why are we only realizing that Gittins is good for tails now?

Practice: what are the design takeaways?



Theorem

```
\gamma-Gittins is strongly optimal,
```

```
C_{\gamma-\text{Gittins}} = \inf C_{\pi}
```

boost = $\frac{1}{2}\log(\text{Gittins index})$

Theory: why was this hard to discover?







arm states can change when out of service





Can view $\mathbf{P}[T_{\pi} > x]$ as:

E[cost of job] when cost is $\mathbf{1}(T_{\pi} > x)$

arm states can change when out of service





Can view $\mathbf{P}[T_{\pi} > x]$ as:

arm states can change when out of service

E[cost of job] when cost is $\mathbf{1}(T_{\pi} > x)$ restless process





restless bandits are hard



E[cost of job] when cost is $\mathbf{1}(T_{\pi} > x)$ restless process









Gittins was designed specifically for *non-restless* bandits





Gittins was designed specifically for *non-restless* bandits



for optimizing asymptotic behavior of $\mathbf{P}[T_{\pi} > x]$, correct choice for "cost of job" is $e^{\gamma T_{\pi}}$







Gittins was designed specifically for *non-restless* bandits

or optimizing asymptotic behavior of $\mathbf{P}[T_{\pi} > x]$, correct choice for "cost of job" is $e^{\gamma T_{\pi}}$





Arm 1:



Gittins is optimal for multi-armed bandits with arrivals













Gittins is optimal for multi-armed bandits with arrivals

...but only if arrivals are time-homogenous











Goal: minimize $\mathbf{E}[e^{\gamma T_{\pi}}]$





- Gittins is optimal for multi-armed bandits with arrivals
 - ...but only if arrivals are time-homogenous





Gittins is optimal for multi-armed bandits with arrivals

...but only if arrivals are time-homogenous

Goal: minimize **E**

A: arri D_{π} : depa





$$[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$$
ival time
arture time





Gittins is optimal for multi-armed bandits with arrivals

...but only if arrivals are time-homogenous











Gittins is optimal for multi-armed bandits with arrivals

...but only if arrivals are time-homogenous

job-specific costs depends on arrival time \rightarrow time-inhomogeneous arrivals













Gittins is optimal for multi-armed bandits with arrivals

...but only if arrivals are time-homogenous



job-specific costs depends on arrival time \rightarrow *time-inhomogeneous arrivals*









job-specific costs depends on arrival time \rightarrow *time-inhomogeneous arrivals*




$$4 \quad S = 1$$















Goal: minimize $\mathbf{E}[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$

batch problem:







Goal: minimize $\mathbf{E}[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$

multi-armed bandit problem

batch problem:









Goal: minimize]

multi-armed bandit problem

$$\mathbf{E}[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$$

Solution: serve job with greatest Gittins(cost, state)







Goal: minimize

multi-armed bandit problem

batch problem:

Solution: serve job with greatest Gittins(cost, state)

$$\mathbf{E}[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$$

Gittins(cost, state) = cost·Gittins(1, state)







Goal: minimize]

multi-armed bandit problem

batch problem:

Solution: serve job with greatest Gittins(cost, state)

 $-\frac{1}{\gamma}\log(\operatorname{Gittins}(e^{-\gamma A},\operatorname{state}))$

$$\mathbf{E}[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$$

Gittins(cost, state) = cost·Gittins(1, state)

e)) =
$$A - \frac{1}{\gamma} \log(\text{Gittins}(1, \text{state}))$$







Goal: minimize]

multi-armed bandit problem

batch problem:

Solution: serve job with greatest Gittins(cost, state)

Gittins(cost, state) = cost·Gittins(1, state)

$$\mathbf{E}[e^{\gamma T_{\pi}}] = \mathbf{E}[e^{-\gamma A}e^{\gamma D_{\pi}}]$$

boost =
$$\frac{1}{\gamma} \log(\text{Gittins})$$

 $-\frac{1}{\nu}\log(\operatorname{Gittins}(e^{-\gamma A}, \operatorname{state})) = A - \frac{1}{\nu}\log(\operatorname{Gittins}(1, \operatorname{state}))$







multi-armed bandit problem

batch problem:









Ignoring arrivals in the γ -Gittins analysis?







Ignoring arrivals in the γ -Gittins analysis?

job sizes may be correlated when sampled as busy period







Ignoring arrivals in the γ -Gittins analysis?







why is it surprising that Gittins is the right choice?

Theory: why was this hard to discover?





why is it surprising that Gittins is the right choice?

Theory: why was this hard to discover?

• initially looks like restless bandit: hard, needs different tools





• initially looks like restless bandit: hard, needs different tools

• MAB with arrivals: have time-inhomogeneous arrivals

why is it surprising that Gittins is the right choice?

Theory: why was this hard to discover?





why is it surprising that Gittins is the right choice?

Theory: why was this hard to discover?

• initially looks like restless bandit: hard, needs different tools

• MAB with arrivals: have time-inhomogeneous arrivals

• ignoring arrivals: busy period \rightarrow batch problem leads to correlation





quantitative economic argument!

why is it surprising that Gittins is the right choice?

Theory: why was this hard to discover?

• initially looks like restless bandit: hard, needs different tools

• MAB with arrivals: have time-inhomogeneous arrivals

• ignoring arrivals: busy period \rightarrow batch problem leads to correlation







If you want to minimize the tail of response time:

- 1. priority should depends on both arrival time and job state
- 2. priority \approx big initial boost, then decrease based on state
- 3. priority can also increase based on state, but times where it decreases are more important
- 4. can use any information model, but more information \implies better performance



bonus: why might there be job size correlations?

server idle



$$S_3$$

Jnif{1, ε }

$$a_3 = 1$$



bonus: why might there be job size correlations?



Let $\varepsilon \ll 1$ and $S = \text{Unif}\{1, \varepsilon\}$

3 jobs:
$$a_1 = 0$$
, $a_2 = \varepsilon^2$, $a_3 = 1$

If
$$S_1 = \varepsilon$$
 then $S_2 = 1$

$$S_3$$

Jnif{1, ε }

$$a_3 = 1$$



26

bonus: why might there be job size correlations?





26