

# A Graph Patrol Problem with Locally-Observable Random Attackers

Thomas Lowbridge and David Hodge

University Of Nottingham,UK

June 15, 2018

- Introduction to game
  - Game setup
  - Markov Decision Process(MDP) formulation
- Problem relaxation
  - Multi-node relaxation
  - Total-rate constraint
  - Lagrangian relaxation
- Single node problem
  - Developing a threshold policy
  - Fair costs
- Indices and heuristics

We wish to model a moving patroller with attackers who arrive on a network, with some type of suspicious attackers.

As an example, we could consider a security guard patrolling a museum, with attackers who are trying to steal parts of the exhibits in each room. We implement semi-intelligent thieves, if they arrive and see the guard they will wait until they try to steal.

A Patrolling game with random attackers and local-observations,  $G = G(Q, \mathbf{X}, \mathbf{b}, \boldsymbol{\lambda}, \mathbf{c})$  is made of 5 major components.

- A **Graph**,  $Q = (N, E)$ , made of nodes,  $N$  ( $|N| = n$ ), and a set of edges,  $E$  and an adjacency matrix,  $A$ .
- A vector of **attack time distributions**,  $\mathbf{X} = (X_1, \dots, X_n)$ .
- A vector of **observable capacities**,  $\mathbf{b} = (b_1, \dots, b_n)$ .
- A vector of **poisson arrival rates**,  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$ .
- A vector of **costs**,  $\mathbf{c} = (c_1, \dots, c_n)$ .

A Patrolling game with random attackers and local-observations,  $G = G(Q, \mathbf{X}, \mathbf{b}, \boldsymbol{\lambda}, \mathbf{c})$  is made of 5 major components.

- A **Graph**,  $Q = (N, E)$ , made of nodes,  $N$  ( $|N| = n$ ), and a set of edges,  $E$  and an adjacency matrix,  $A$ .
- A vector of **attack time distributions**,  $\mathbf{X} = (X_1, \dots, X_n)$ .
- A vector of **observable capacities**,  $\mathbf{b} = (b_1, \dots, b_n)$ .
- A vector of **poisson arrival rates**,  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$ .
- A vector of **costs**,  $\mathbf{c} = (c_1, \dots, c_n)$ .

With the **patroller moving instantaneously** between nodes and attackers who **arrive when the patroller is present waiting till the next time period** to begin their attack (who are observed as suspicious by the patroller).

# Example of Game

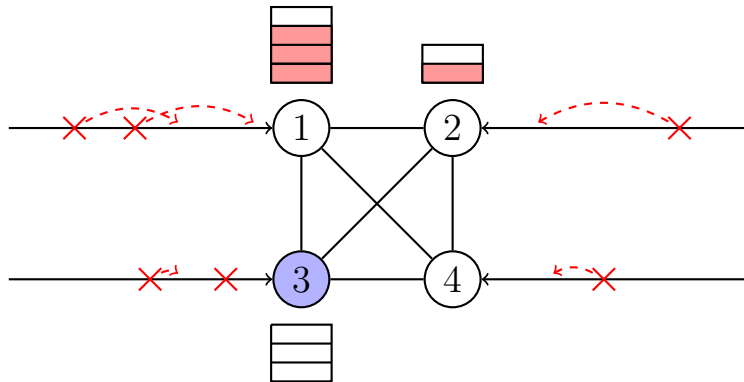


Figure: Example of  $G = (K_4, \mathbf{X}, \mathbf{b}, \boldsymbol{\lambda}, \mathbf{c})$  with patroller currently at node 3

# Example of timing

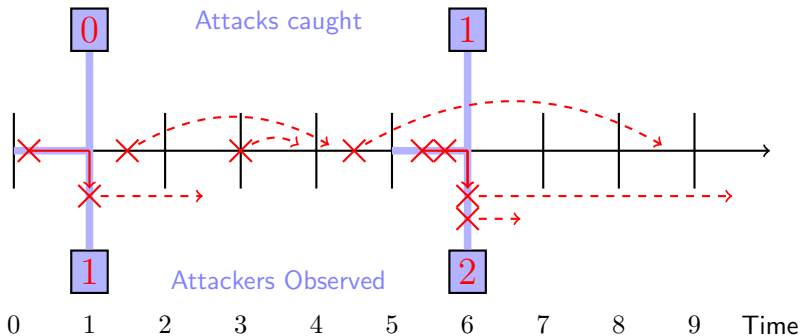


Figure: Example of timing at a node

# Markov Decision Process(MDP) formulation

This game is a MDP with states  $(s, v)$ , where  $s_i$  is the time since node  $i$  was last chosen to be visited and  $v_i$  is how many local-observations where observed when node  $i$  was last visited. With state space  $\Omega = \{(s, v) | s_i = 1, 2, \dots \text{ and } v_i = 0, 1, \dots, b_i \ \forall i = 1, \dots, n\}$ .



# Markov Decision Process(MDP) formulation

This game is a MDP with states  $(\mathbf{s}, \mathbf{v})$ , where  $s_i$  is the time since node  $i$  was last chosen to be visited and  $v_i$  is how many local-observations were observed when node  $i$  was last visited. With state space  $\Omega = \{(\mathbf{s}, \mathbf{v}) | s_i = 1, 2, \dots \text{ and } v_i = 0, 1, \dots, b_i \ \forall i = 1, \dots, n\}$ .

The current node can be identified by  $l(\mathbf{s}) = \operatorname{argmin}_i s_i$ . The current node will have  $s_i = 1$ . The available actions from any state is  $\mathcal{A}(\mathbf{s}) = \{j | A_{l(\mathbf{s}),j} = 1\}$ , that is nodes we can move to from the current node.

# Markov Decision Process(MDP) formulation

This game is a MDP with states  $(\mathbf{s}, \mathbf{v})$ , where  $s_i$  is the **time** since node  $i$  was last chosen to be visited and  $v_i$  is how many **local-observations** where observed when node  $i$  was last visited. With state space  $\Omega = \{(\mathbf{s}, \mathbf{v}) | s_i = 1, 2, \dots \text{ and } v_i = 0, 1, \dots, b_i \ \forall i = 1, \dots, n\}$ .

The current node can be identified by  $l(\mathbf{s}) = \operatorname{argmin}_i s_i$ . The current node will have  $s_i = 1$ . The available actions from any state is  $\mathcal{A}(\mathbf{s}) = \{j | A_{l(\mathbf{s}),j} = 1\}$ , that is nodes we can move to from the current node.

A transition from a state,  $(\mathbf{s}, \mathbf{v})$ , with a chosen action,  $i$ , is  $\phi(\mathbf{s}, \mathbf{v}, i) = (\tilde{\mathbf{s}}, \tilde{\mathbf{v}})$  where  $\tilde{\mathbf{s}}$  has;  $\tilde{s}_i = 1$  and  $\tilde{s}_j = s_j + 1 \ \forall j \neq i$  and  $\tilde{\mathbf{v}}$  has;  $\tilde{v}_i \sim TPO(\lambda_i, b_i)$  and  $\tilde{v}_j = v_j \ \forall j \neq i$ .

Where  $TPO(\lambda, b)$  is the Poisson distribution truncated at the value  $b$ . I.e

$$P(TPO(\lambda, b)) = \begin{cases} P(Po(\lambda)) & \text{if } i \neq b \\ P(Po(\lambda) \geq i) & \text{if } i = b \\ 0 & \text{Otherwise} \end{cases}$$

Note.  $TPO(\lambda, \infty) = Po(\lambda)$ .

To avoid certain problems, we would like to limit ourselves to a finite state space

As in Lin et al. [2013], we will bound the attack times to create a finite state space, and define  $B_j \equiv \min\{k \mid k \in \mathbb{Z}^+, P(X_j \leq k) = 1\}$ . We will also bound our observable capacities to finite integers  $b_i \in \mathbb{Z}_0^+$ . This restricts our state space to a finite size

$$\Omega = \{(\mathbf{s}, \mathbf{v}) \mid s_i = 1, 2, \dots, B_i + 1 \text{ and } v_i = 0, 1, \dots, b_i \forall i = 1, \dots, n\}.$$

To avoid certain problems, we would like to limit ourselves to a finite state space

As in Lin et al. [2013], we will bound the attack times to create a finite state space, and define  $B_j \equiv \min\{k \mid k \in \mathbb{Z}^+, P(X_j \leq k) = 1\}$ . We will also bound our observable capacities to finite integers  $b_i \in \mathbb{Z}_0^+$ . This restricts our state space to a finite size

$$\Omega = \{(\mathbf{s}, \mathbf{v}) \mid s_i = 1, 2, \dots, B_i + 1 \text{ and } v_i = 0, 1, \dots, b_i \forall i = 1, \dots, n\}.$$

Now however our transitions are limited and become  $\phi(\mathbf{s}, \mathbf{v}, i) = (\tilde{\mathbf{s}}, \tilde{\mathbf{v}})$  where  $\tilde{\mathbf{s}}$  has;  $\tilde{s}_i = 1$  and  $\tilde{s}_j = \min\{s_j + 1, B_j + 1\} \forall j \neq i$  and  $\tilde{\mathbf{v}}$  has;  $\tilde{v}_i \sim TPO(\lambda_i, b_i)$  and  $\tilde{v}_j = v_i \forall j \neq i$ .

# Example of state space

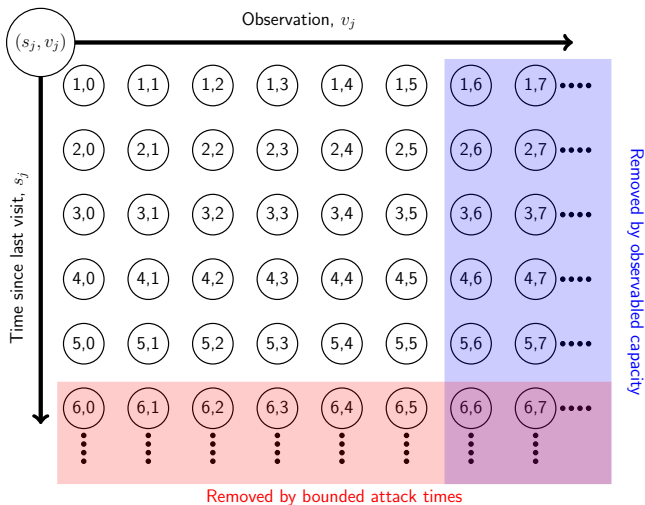


Figure: State space diagram, with  $b_j = 5$  and  $B_j = 4$  (e.g.  $X_j \leq 3.7$ )

# Markov Decision Process(MDP) formulation

The cost at a node is zero if that node is chosen, otherwise it is the cost of arrivals finishing in the next time period, plus the cost of local-observations finishing in the next time period.

$$C_j(\mathbf{s}, \mathbf{v}, i) = \begin{cases} \underbrace{c_j \lambda_j \int_{s_j-1}^{s_j} P(X_j \leq t) dt}_{\text{Arrivals finishing}} + \underbrace{c_j v_j P(s_j - 1 < X_j \leq s_j)}_{\text{Local-observed finishing}} & \text{if } j \neq i \\ 0 & \text{if } j = i \end{cases}$$

# Markov Decision Process(MDP) formulation

The cost at a node is zero if that node is chosen, otherwise it is the cost of arrivals finishing in the next time period, plus the cost of local-observations finishing in the next time period.

$$C_j(\mathbf{s}, \mathbf{v}, i) = \begin{cases} \underbrace{c_j \lambda_j \int_{s_j-1}^{s_j} P(X_j \leq t) dt}_{\text{Arrivals finishing}} + \underbrace{c_j v_j P(s_j - 1 < X_j \leq s_j)}_{\text{Local-observed finishing}} & \text{if } j \neq i \\ 0 & \text{if } j = i \end{cases}$$

We will restrict ourselves to deterministic attack times,  
 $P(X_j = x_j) = 1 \forall j$  so that

$$C_j(\mathbf{s}, \mathbf{v}, i) = \begin{cases} c_j \lambda_j R_j + c_j v_j & \text{for } s_j = B_j, i \neq j \\ c_j \lambda_j & \text{for } s_j = B_j + 1, i \neq j \\ 0 & \text{Otherwise} \end{cases}$$

Where  $R_j = B_j - x_j$ . We see that we only worry about incurring costs in states when  $s_i = B_i$  or  $B_i + 1$  for some  $i$ .

# Objective

Due to the finiteness of the state space, we can just focus on **stationary, deterministic** policies,  $\pi : \Omega \rightarrow \mathcal{A} \in \Pi$ . The patroller wants to choose a policy such that the long-run average cost is minimized, that is find the minimum cost (and policy) defined by,

$$C^{\text{OPT}}(s_0, v_0) \equiv \min_{\pi \in \Pi} \sum_{i=1}^n V_i(\pi, s_0, v_0)$$

Where  $V_i(\pi, s_0, v_0)$  is the **long-run average cost** incurred at node  $i$  under the policy,  $\pi$ , starting from state,  $(s_0, v_0)$  defined by ,

$$V_i(\pi, s_0, v_0) \equiv \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} C_i(\phi_{\pi}^k(s_0, v_0), \pi(\phi_{\pi}^k(s_0, v_0)))$$

Where  $\phi_{\pi}^k(s_0, v_0)$  is the state after  $k$  transitions starting from  $(s_0, v_0)$  under the policy  $\pi$ .



# Multiple node relaxation

We now relax the problem to that of a patroller who can visit multiple nodes in the next time period.

We will call this Problem the **multiple-node(MN)** problem. We extend the class of policies to

$$\Pi^{\text{MN}} = \{\pi \mid \pi : \Omega \rightarrow \{\alpha \mid \alpha_i \in \{0, 1\} \text{ for } i = 1, \dots, n\}\}$$

Where  $\alpha_i = 1$  if the patroller will visit node  $i$  in the next time period and  $\alpha_i = 0$  if the patroller will not visit node  $i$  in the next time period. Note.  $\Pi \subset \Pi^{\text{MN}}$ . Similar to  $C^{\text{OPT}}$  we define  $C^{\text{MN}}$  just with this larger class of policies.

## Definition (Long-run visit rate)

Define the long-run rate of visits to node  $i$ , starting at  $(s_0, v_0)$  under policy  $\pi$  by

$$\mu_i(\pi, s_0, v_0) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \alpha_{\phi_{\pi}^k(s_0, v_0), i}$$

## Definition (Long-run visit rate)

Define the long-run rate of visits to node  $i$ , starting at  $(s_0, v_0)$  under policy  $\pi$  by

$$\mu_i(\pi, s_0, v_0) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \alpha_{\phi_{\pi}^k(s_0, v_0), i}$$

We now impose the **total-rate** constraint, that is the long-run overall visit rate to all nodes is no greater than 1. and restrict the MN problem by this to get the total-rate(TR) problem and its class of policies

$$\Pi^{\text{TR}} = \left\{ \pi \in \Pi^{\text{MN}} \mid \sum_{i=1}^n \mu_i(\pi, s_0, v_0) \leq 1 \quad \forall (s_0, v_0) \in \Omega \right\}$$

Again  $\Pi \subset \Pi^{\text{TR}}$  and we define  $C^{\text{TR}}$  similar to  $C^{\text{OPT}}$  just on this larger class of policies.

We now relax the problem again by incorporating the total rate constraint into the objective function, with a Lagrange multiplier,  $\omega \geq 0$ . This forms

$$\begin{aligned} C(\omega) &= \min_{\pi \in \Pi^{\text{MN}}} \left\{ \sum_{i=1}^n V_i(\pi) + \omega \left( \sum_{i=0}^n \mu_i(\pi) - 1 \right) \right\} \\ &= \min_{\pi \in \Pi^{\text{MN}}} \sum_{i=1}^n (V_i(\pi) + \omega \mu_i(\pi)) - \omega \end{aligned}$$

By incorporating the total-rate constraint as a Lagrange multiplier we can drop the constraint, so that the patroller can choose to visit any number of nodes in each time period (each costing  $\omega$ ).

We now relax the problem again by incorporating the total rate constraint into the objective function, with a Lagrange multiplier,  $\omega \geq 0$ . This forms

$$\begin{aligned} C(\omega) &= \min_{\pi \in \Pi^{\text{MN}}} \left\{ \sum_{i=1}^n V_i(\pi) + \omega \left( \sum_{i=0}^n \mu_i(\pi) - 1 \right) \right\} \\ &= \min_{\pi \in \Pi^{\text{MN}}} \sum_{i=1}^n (V_i(\pi) + \omega \mu_i(\pi)) - \omega \end{aligned}$$

By incorporating the total-rate constraint as a Lagrange multiplier we can drop the constraint, so that the patroller can choose to visit any number of nodes in each time period (each costing  $\omega$ ).

We can show that,  $C^{\text{OPT}} \geq C^{\text{TR}} \geq C(\omega)$ .

# Single node problem

We now want to find  $C(\omega) = \min_{\pi \in \Pi^{\text{MN}}} \sum_{i=1}^n (V_i(\pi) + \omega \mu_i(\pi)) - \omega$ . This problem, as we are in the multi-node class, can be rewrote as

$$C(\omega) = \sum_{i=1}^n \left( \min_{\pi \in \Pi^{\text{MN}}} (V_i(\pi) + \omega \mu_i(\pi)) \right) - \omega$$

# Single node problem

We now want to find  $C(\omega) = \min_{\pi \in \Pi^{MN}} \sum_{i=1}^n (V_i(\pi) + \omega \mu_i(\pi)) - \omega$ . This problem, as we are in the multi-node class, can be rewritten as

$$C(\omega) = \sum_{i=1}^n \left( \min_{\pi \in \Pi^{MN}} (V_i(\pi) + \omega \mu_i(\pi)) \right) - \omega$$

That is for every node,  $i$ , try to minimize  $V_i(\pi) + \omega \mu_i(\pi)$ , where  $\omega$  can be interpreted as the service charge for visiting the node. For now we drop the node subscript and just try to find the following,

$$g = \min_{\pi \in \Pi^{MN}} V(\pi) + \omega \mu(\pi)$$

# Single node problem

We now want to find  $C(\omega) = \min_{\pi \in \Pi^{MN}} \sum_{i=1}^n (V_i(\pi) + \omega \mu_i(\pi)) - \omega$ . This problem, as we are in the multi-node class, can be rewritten as

$$C(\omega) = \sum_{i=1}^n \left( \min_{\pi \in \Pi^{MN}} (V_i(\pi) + \omega \mu_i(\pi)) \right) - \omega$$

That is for every node,  $i$ , try to minimize  $V_i(\pi) + \omega \mu_i(\pi)$ , where  $\omega$  can be interpreted as the service charge for visiting the node. For now we drop the node subscript and just try to find the following,

$$g = \min_{\pi \in \Pi^{MN}} V(\pi) + \omega \mu(\pi)$$

We can solve this problem to find bounds on  $g$  suggesting a certain type of threshold policy



## Definition (Threshold Policy)

Policy,  $\pi_{\text{Th}}(v_{\text{crit}})$ , is the policy which in states:

- $(s, v)$ ,  $s < B$  waits until  $(B, v)$
- $(B, v)$  renews now if  $v \geq v_{\text{crit}}$  and waits until  $(B + 1, v)$  if  $v < v_{\text{crit}}$
- $(B + 1, v)$  renews now.

# Developing a Threshold policy

## Definition (Threshold Policy)

Policy,  $\pi_{\text{Th}}(v_{\text{crit}})$ , is the policy which in states:

- $(s, v)$ ,  $s < B$  waits until  $(B, v)$
- $(B, v)$  renews now if  $v \geq v_{\text{crit}}$  and waits until  $(B + 1, v)$  if  $v < v_{\text{crit}}$
- $(B + 1, v)$  renews now.

Because of one of the bounds there is a maximum for  $v_{\text{crit}}$ ,

$v_{\text{max}} \equiv \max\{v \in \{0, 1, \dots, b\} \mid v \leq \lambda(1 - R)\}$ , with

$v_{\text{crit}} \in \{0, 1, \dots, v_{\text{max}} + 1\}$ .

## Definition (Threshold Policy)

Policy,  $\pi_{\text{Th}}(v_{\text{crit}})$ , is the policy which in states:

- $(s, v)$ ,  $s < B$  waits until  $(B, v)$
- $(B, v)$  renews now if  $v \geq v_{\text{crit}}$  and waits until  $(B + 1, v)$  if  $v < v_{\text{crit}}$
- $(B + 1, v)$  renews now.

Because of one of the bounds there is a maximum for  $v_{\text{crit}}$ ,

$v_{\text{max}} \equiv \max\{v \in \{0, 1, \dots, b\} \mid v \leq \lambda(1 - R)\}$ , with

$v_{\text{crit}} \in \{0, 1, \dots, v_{\text{max}} + 1\}$ . We get bounds on the long-run average cost from our solution of

- $g \leq c\lambda R$  if  $v_{\text{crit}} = 0$ .
- $c\lambda R + c(k - 1) < g \leq c\lambda R + kc$  if  $v_{\text{crit}} \neq 0, v_{\text{max}} + 1$
- $c\lambda R + cv_{\text{max}} < g \leq c\lambda$  if  $v_{\text{crit}} = v_{\text{max}} + 1$

# Example of a threshold policy

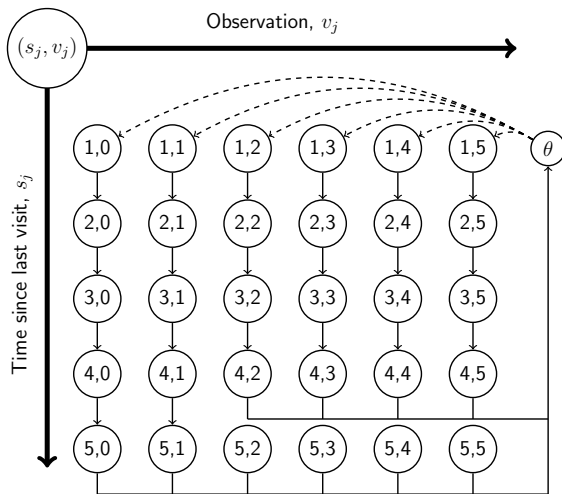


Figure 6.1: Threshold policy,  $\pi_{\text{Th}}(2)$ , with  $b_j = 5$  and  $B_j = 4$  (e.g.  $X_j \leq 3.7$ )

# Converting bounds

By using such a strategy of  $v_{\text{crit}}$  threshold, we get a long-run average cost of

$$g^{\pi_{\text{Th}}(v_{\text{crit}})}(\omega) = \frac{\text{Expected cost per renewal}}{\text{Expected renewal length}}$$
$$= \frac{\omega + c\lambda R(1 - P(TPo(\lambda, b) \geq v_{\text{crit}})) + c \sum_{i=0}^{v_{\text{crit}}-1} iP(TPo(\lambda, b) = i)}{B + 1 - P(TPo(\lambda, b) \geq v_{\text{crit}})}$$

Converting the bounds on  $g$  to bounds on  $\omega$  gives us the fair price

- $0 \leq \omega \leq c\lambda R(B + 1) \equiv \Delta(0)$  if  $v_{\text{crit}} = 0$ .
- $\Delta(v_{\text{crit}} - 1) < \omega \leq \Delta(v_{\text{crit}})$  if  $v_{\text{crit}} \neq 0, v_{\text{max}}$
- $\Delta(v_{\text{max}}) < \omega \leq \tilde{\Delta}$  if  $v_{\text{crit}} = v_{\text{max}} + 1$

Where  $\Delta, \tilde{\Delta}$  are appropriately defined by rearrangement. These now give us a fair cost to renew in  $(B, v)$  depending on  $v$ .

We can use these bounds on the visiting cost,  $\omega$ , to develop an index by reinserting the node subscript,  $i$ , to create

## Definition (Node Index)

We define the index for node  $i$ , when the system is in state,  $(s, v)$ , to be

$$W_i(s_i, v_i) = \begin{cases} 0 & \text{If } s_i < B_i, \\ \Delta_i(v_i) & \text{If } s_i = B_i, v_i < v_{i,\max}, \\ \widetilde{\Delta}_i & \text{If } s_i = B_i, v_i \geq v_{i,\max}, \\ \widetilde{\Delta}_i & \text{If } s_i = B_i + 1. \end{cases}$$

# Example of node index

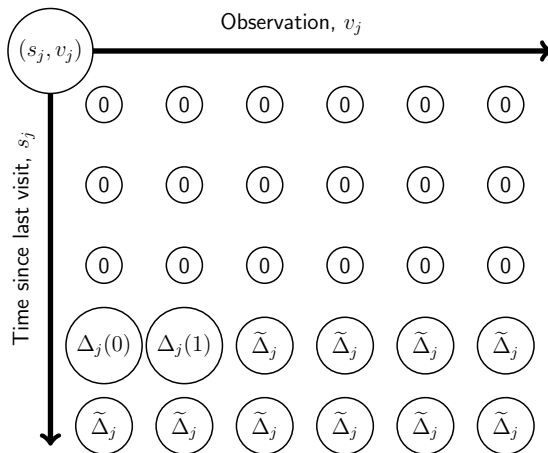


Figure 7.1: Index value for states

# Heuristic types

We now use the index to develop two types of heuristics, using the indices for states:



We now use the index to develop two types of heuristics, using the indices for states:

## Definition (Benefit Heuristic(BH))

The benefit heuristic(BH) has a patroller currently in state  $(s, v)$  look at all paths of length,  $l$ . They choose the path with the largest aggregate of indices collected along the path. The action is the first node in the path.

# Heuristic types

We now use the index to develop two types of heuristics, using the indices for states:

## Definition (Benefit Heuristic(BH))

The benefit heuristic(BH) has a patroller currently in state  $(s, v)$  look at all paths of length,  $l$ . They choose the path with the largest aggregate of indices collected along the path. The action is the first node in the path.

## Definition (Penalty Heuristic(PH))

To develop a penalty heuristic(PH) based on the index, have a patroller currently in state  $(s, v)$  look at all paths of length,  $l$ . They choose the path with the smallest aggregate of indices not collected along the path. The action is the first node in the path.

If we are looking at all paths of length  $l$ , we might as well consider applying the heuristic to all paths of lengths,  $1, \dots, l$ . We call this the heuristic depth,  $d = l$ .

Once we apply our heuristic choice we are left with  $d$  paths and therefore  $d$  suggestions for the first action to take. To decide on which action to take, we look at the average cost the path incurs and from the state at the end of the path we look at the average cost to decay to the neglected state  $(B + 1, v)$ .

We will call the depth heuristics,  $BH(d)$  and  $PH(d)$ . We now perform numerical experiments, for 500 scenarios on  $K_4$ , using the  $PH(d)$ , as suggested by Lin et al. [2013] instead of  $BH(d)$ , to illustrate a problem with the current index.

# Numerical Experiments

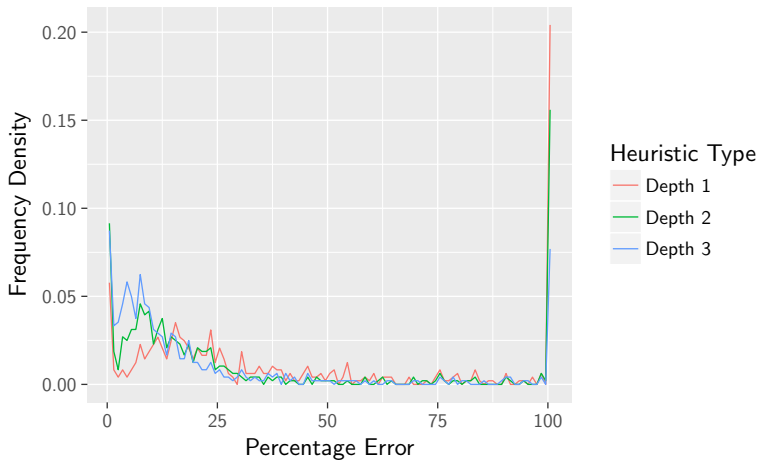


Figure: Percentage Error frequency for PH of depths 1, 2, 3

# An alternative index

We suggest a slight alteration, to uncap the index for,  $s_i = B_i, v_i \geq v_{i,\max}$ . This means that the index prioritizes observations that are about to finish.

## Definition (Alternative Plain Node Index)

We define the index for node  $i$ , when the system is in state,  $(s, v)$ , to be

$$W_i(s_i, v_i) = \begin{cases} 0 & \text{If } s_i < B_i, \\ \Delta_i(v_i) & \text{If } s_i = B_i, \\ \Delta_i(v_i + 1) & \text{If } s_i = B_i + 1, v_i < v_{i,\max}, \\ \widetilde{\Delta}_i & \text{If } s_i = B_i + 1, v_i \geq v_{i,\max}. \end{cases}$$

# Suggestion of node index fix

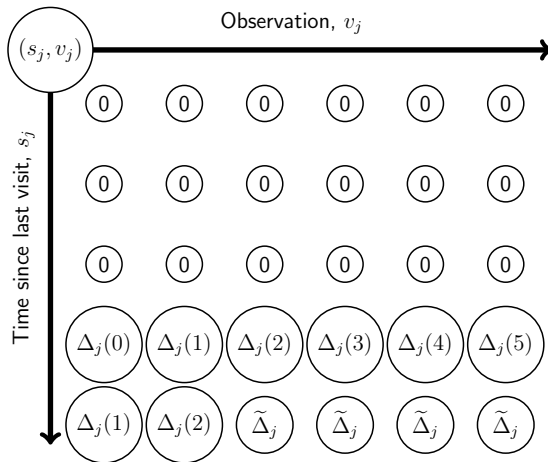
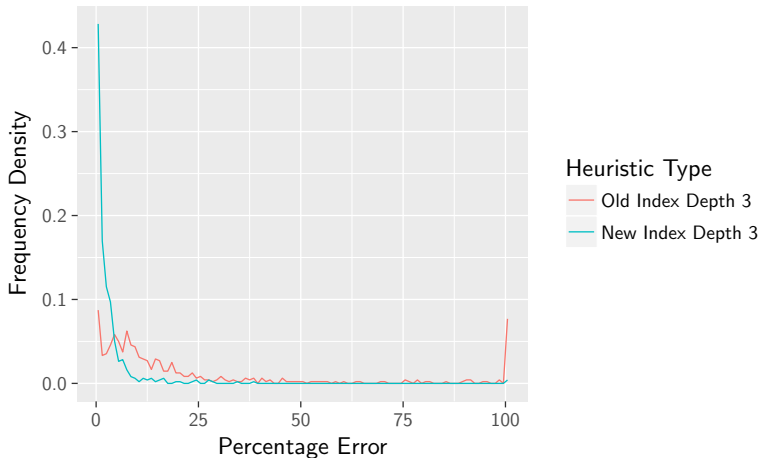


Figure: New suggested index value for states

# Comparing indices



**Figure:** Percentage Error frequency for PH of depth 3 for old and new index

- Investigate whether different types of graphs are more susceptible to different heuristic/Index types.
- Look at improving the heuristic, by working with cyclic costs rather than proxy average costs and decay costs.
- Improve theory to deal with generic distributions.
- Improve theory to deal with no waiting capacities.



Kyle Y Lin, Michael P Atkinson, Timothy H Chung, and Kevin D Glazebrook. A Graph Patrol Problem with Random Attack Times. *Operations Research*, 61(3):694–710, 2013. doi: 10.1287/opre.1120.1149. URL <http://pubsonline.informs.org><http://dx.doi.org/10.1287/opre.1120.1149>.

# Appendix: Developing the optimal policy

# Developing an optimal policy

We aim to develop an index on our state space, that is a fair price for visiting the node for a given state. We shall first show there is no fair price,  $\omega$  for visiting in any state,  $(s, v)$  with  $s < B$ .

From this position consider the policy  $\pi_k$  which waits  $k$  time periods and then renews and follows the optimal policy,  $\sigma$ , with  $k = 0, \dots, B - s$ .

Using such a policy will get us that

$$V_n^{\pi_k}(x, v) = \omega + E[V_{n-k-1}^{\sigma}(\theta)] \quad (1)$$

where  $\theta$  is the state upon a visit (i.e it is the state  $(1, V) \sim (1, TPO(\lambda))$ ).

Now we will pick policy  $\pi_{k+1}$  over  $\pi_k$  (or be indifferent) if

$$\begin{aligned} \lim_{n \rightarrow \infty} V_n^{\pi_k}(x, v) - V_n^{\pi_{k+1}}(x, v) &\geq 0 \\ \iff \lim_{n \rightarrow \infty} E[V_{n-k}^{\sigma}(\theta) - V_{n-k-1}^{\sigma}(\theta)] &\geq 0 \\ \iff g &\geq 0 \end{aligned}$$

Where  $g$  is the long-run average cost and so we know  $g \geq 0$  hence we will always wait instead of renewing.

# Developing and optimal policy

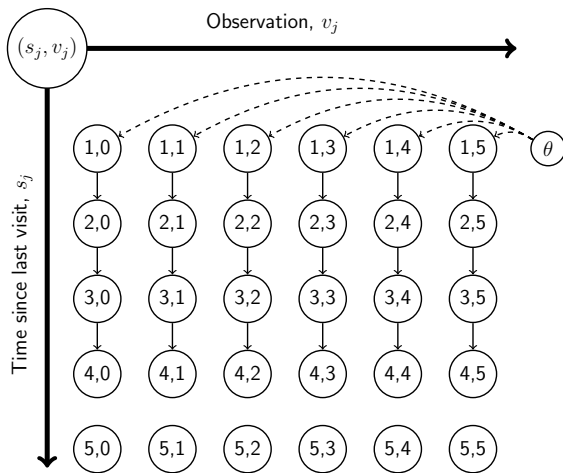


Figure: Threshold policy,  $\pi_{\text{Th}}(2)$ , with  $b_j = 5$  and  $B_j = 4$  (e.g.  $X_j \leq 3.7$ )

# Developing an optimal policy

Now we will skip to the state  $(B + 1, v)$  and suggest again a policy  $\pi_k$  which waits  $k$  time periods before renewing and then follows some optimal policy,  $\sigma$ .

Using such a policy will get us

$$V_n^{\pi_k}(x, v) = \omega + c\lambda k + E[V_{n-k-1}^{\sigma}(\theta)] \quad (2)$$

And again we will pick a policy  $\pi_{k+1}$  over  $\pi_k$  (or be indifferent) if

$$\begin{aligned} \lim_{n \rightarrow \infty} V_n^{\pi_k}(\lfloor B \rfloor + 2, 0) - V_n^{\pi_{k+1}}(\lfloor B \rfloor + 2, 0) &\geq 0 \\ \iff \lim_{n \rightarrow \infty} -c\lambda + E[V_{n-k}^{\sigma}(\theta) - V_{n-k-1}^{\sigma}(\theta)] &\geq 0 \\ \iff g &\geq c\lambda \end{aligned}$$

So hence if  $g \geq c\lambda$  we will wait forever, as this has no dependence on  $k$ . We will now argue that  $g_{\max} = c\lambda$ , that is at worst the long-run average cost is  $c\lambda$ . This can be seen by the strategy  $\pi_{\text{neg}}$ , a strategy which never renews no matter what state we are in. I.e the patroller neglects the node. As for large  $n$  we just pay  $c\lambda$  every time step. Hence in these states if we renew we will renew immediately.

# Developing an optimal policy

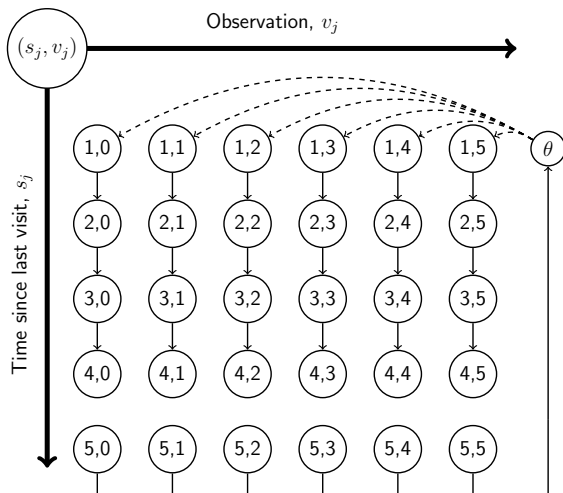


Figure: Threshold policy,  $\pi_{\text{Th}}(2)$ , with  $b_j = 5$  and  $B_j = 4$  (e.g.  $X_j \leq 3.7$ )

# Developing an optimal policy

Our neglecting strategy shows that  $g \leq c\lambda$  and hence if we are in state  $(B, v)$ , we can decide to renew now or renew in one time period (at  $(B + 1, v)$ ) then follow the optimal,  $\sigma$  we will choice to renew now over waiting if

$$\begin{aligned} \lim_{n \rightarrow \infty} V_n^{\pi_1}(B, v) - V_n^{\pi_0}(B, v) &\geq 0 \\ \iff \lim_{n \rightarrow \infty} c\lambda R + vc + E[V_{n-1}^\sigma(B + 1, 0)] - (\omega + E[V_{n-1}^\sigma(\theta)]) &\geq 0 \\ \iff g \leq c(\lambda R + v) \end{aligned}$$

So we renew now if  $g \leq c(\lambda R + v)$ , as we are guaranteed that  $g \leq c\lambda$  we are clearly in this region if  $c(\lambda R + v) \leq c\lambda \iff v \leq \lambda(1 - R)$ . We also note that the bound is increasing in  $v$  so if we will renew in  $v$  we definitely renew in  $v + 1, v + 2, \dots, b$ .