



## Patrol avoider

Résolution d'un jeu stochastique par des MDPs augmentés

CHÂTEL Romain

MOUADDIB Abdel-illah

27 novembre 2019

Jeux stochastiques et MDPs, késako?

Jeu stochastique et surveillance militaire

Modélisations

Premiers résultats

Conclusion

- ▶ Extension de la théorie des jeux à l'incertain
- ▶ Prend l'incertitude en compte (environnement, capteurs, effecteurs ...)
- ▶ Un ou plusieurs agents avec des buts
- ▶ Fonction de gains
- ▶ Maximiser le gain immédiat, ET le gain espéré futur
- ▶ Actions stochastiques, information incomplète/imparfaite

## Observabilité totale

- ▶ Jeux stochastiques à un seul agent
- ▶ Un espace d'états  $S$ , propriété markovienne
- ▶ Un espace d'actions  $A$
- ▶ Un modèle de transitions  $T : S \times A \times S \mapsto [0, 1] \Rightarrow Pr(s'|s, a)$
- ▶ Une fonction de gains/d'utilité  $U : S \times A \mapsto \mathbb{R}$   
 $\Rightarrow U(s, a) = R(s) - Cost(s, a)$
- ▶ Les agents décident en appliquant une politique  $\pi : S \mapsto A$
- ▶ Equation d'optimalité de Bellman :

$$v^{\pi^*}(s) = \max_{a \in A} [U(s, a) + \gamma \sum_{s' \in S} T(s, a, s') v^{\pi^*}(s')]$$

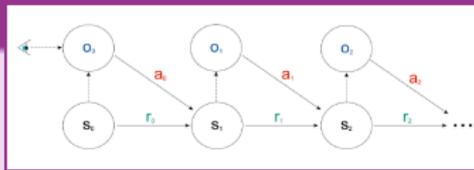
- ▶ Résolution par *Value Iteration*, *Policy Iteration*, *Linear programming*

...

## Observabilité partielle

Extension à l'incertitude sur l'état de l'agent

- ▶ Un espace d'observations  $\Omega$
- ▶ Un modèle d'observations  $TO : \Omega \times S \times A \mapsto [0, 1] \Rightarrow Pr(o|s', a)$
- ▶ Inférer l'état courant à partir des observations  
 $\Rightarrow$  Perte de la propriété markovienne
- ▶ Transformation en BMDP, espace de croyances sur les états  
 $B = [0, 1]^{|S|}$
- ▶ Fonction de mise à jour des croyances  
 $\zeta : B \times A \times B \Rightarrow Pr(b'|b, a) = \sum_{o \in \Omega} Pr(b'|b, a, o) Pr(o|b, a)$
- ▶ Résolution difficile



Jeux stochastiques et MDPs, késako ?

Jeu stochastique et surveillance militaire

Modélisations

Premiers résultats

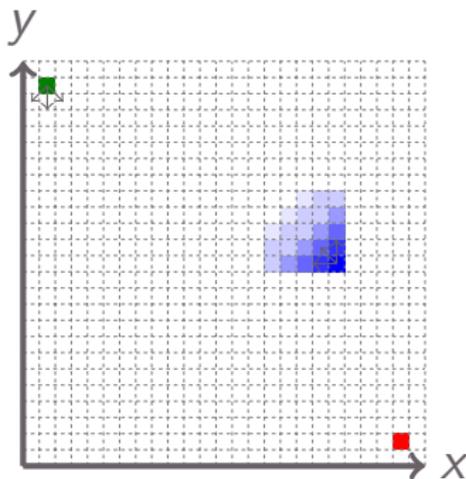
Conclusion

Jeu stochastique : (grille stochastique)

- ▶ Une **cible**
- ▶ Un **drone intrus**
- ▶ Un **drone de surveillance**
  - ▶ Un champ de perception
  - ▶ Plusieurs comportements possibles  
 $k_1 \rightarrow \dots \rightarrow k_2 \rightarrow \dots \rightarrow k_1 \rightarrow \dots$   
 $k_1$  : « parcours en zigzags »  
 $k_2$  : « couvrir la cible »

But :

- ▶ Compromis entre atteindre la cible et prendre le minimum de risque
- ▶ Graal : observabilité partielle sur l'état et le comportement de l'adversaire (satellite)

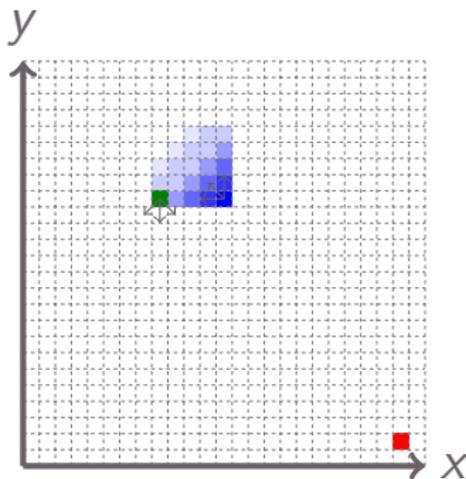


Jeu stochastique : (grille stochastique)

- ▶ Une **cible**
- ▶ Un **drone intrus**
- ▶ Un **drone de surveillance**
  - ▶ Un champ de perception
  - ▶ Plusieurs comportements possibles  
 $k_1 \rightarrow \dots \rightarrow k_2 \rightarrow \dots \rightarrow k_1 \rightarrow \dots$   
 $k_1$  : « parcours en zigzags »  
 $k_2$  : « couvrir la cible »

But :

- ▶ Compromis entre atteindre la cible et prendre le minimum de risque
- ▶ Graal : observabilité partielle sur l'état et le comportement de l'adversaire (satellite)

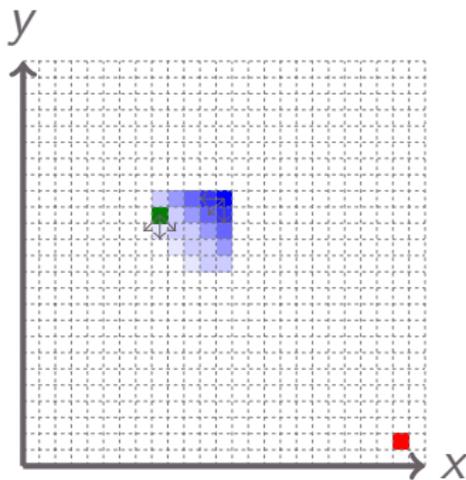


Jeu stochastique : (grille stochastique)

- ▶ Une **cible**
- ▶ Un **drone intrus**
- ▶ Un **drone de surveillance**
  - ▶ Un champ de perception
  - ▶ Plusieurs comportements possibles  
 $k_1 \rightarrow \dots \rightarrow k_2 \rightarrow \dots \rightarrow k_1 \rightarrow \dots$   
 $k_1$  : « parcours en zigzags »  
 $k_2$  : « couvrir la cible »

But :

- ▶ Compromis entre atteindre la cible et prendre le minimum de risque
- ▶ Graal : observabilité partielle sur l'état et le comportement de l'adversaire (satellite)

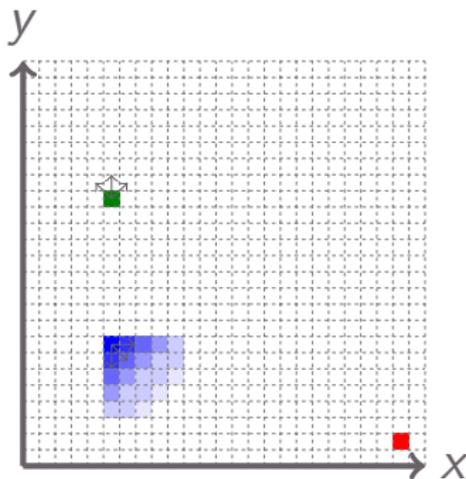


Jeu stochastique : (grille stochastique)

- ▶ Une **cible**
- ▶ Un **drone intrus**
- ▶ Un **drone de surveillance**
  - ▶ Un champ de perception
  - ▶ Plusieurs comportements possibles  
 $k_1 \rightarrow \dots \rightarrow k_2 \rightarrow \dots \rightarrow k_1 \rightarrow \dots$   
 $k_1$  : « parcours en zigzags »  
 $k_2$  : « couvrir la cible »

But :

- ▶ Compromis entre atteindre la cible et prendre le minimum de risque
- ▶ Graal : observabilité partielle sur l'état et le comportement de l'adversaire (satellite)



## Jeux de sécurité

- ▶ Une équipe en défense qui protège des cibles (allocation de ressources)
- ▶ Une équipe en attaque qui doit choisir quelle cible attaquer
- ▶ Résolution par jeu de Stackelberg
- ▶ Network Security Game (difficile)

## Intérêt

- ▶ Modélisation du point de vue de l'attaquant
- ▶ Restriction à un ensemble de comportements donnés pour la défense (non-optimale)
- ▶ Dynamique comportementale complexe

## MDP et Optimalité : Vorobeychik and Singh (2012)

*« Pour tout jeu de Stackelberg stochastique à somme générale et à facteur d'atténuation, si le meneur suit une politique markovienne stationnaire alors il existe une politique déterministe markovienne et stationnaire qui est une politique best-response pour le suiveur. »*

## Contraintes

- ▶ Utilisation de MDP  $\gamma$ -pondérés
- ▶ Le patrouilleur doit suivre une politique markovienne stationnaire

1. Observabilité totale
2. Observabilité partielle sur la politique du patrouilleur
3. Observabilité partielle sur l'état du patrouilleur

Jeux stochastiques et MDPs, késako ?

Jeu stochastique et surveillance militaire

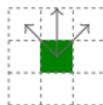
**Modélisations**

Premiers résultats

Conclusion

Symbole	Description	Domaine
$X$	Ensemble des lignes de la grille	$\{1, \dots, R\}$
$Y$	Ensemble des colonnes de la grille	$\{1, \dots, C\}$
$O$	Orientation des drones	$\{N, NE, E, SE, S, SW, W, NW\}$
$S^I$	Espace d'états de l'intrus	$X \times Y \times O$
$S^P$	Espace d'états du patrouilleur	$X \times Y \times O \times \dots$
$i$	État de l'intrus	$S^I$
$i _{XY}$	État de l'intrus projeté sur $XY$	$X \times Y$
$p$	État du patrouilleur	$S^P$
$\Pi^P$	Politiques possibles du patrouilleur	$\{k_1, \dots, k_N\}$
$k$	Politique courante du patrouilleur	$\Pi^P$
$\mathcal{T}^P$	Modèle de transition du patrouilleur	$S^P \times A \times S^P \mapsto [0, 1]$
$\tau$	Position de la cible	$X \times Y$

- ▶ Modèle commun  $A \Rightarrow$  simplicité
- ▶ Avancer d'une case dans une des 3 directions possibles :  
{GO-UP-LEFT, GO-STRAIGHT, GO-UP-RIGHT}



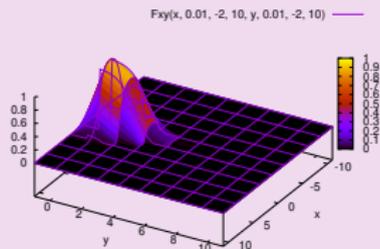
Équivalent à :

Orientation	Actions possibles
NORTH	{GO-NORTH, GO-NORTH-EAST, GO-NORTH-WEST}
NORTH-EAST	{GO-NORTH-EAST, GO-EAST, GO-NORTH}
EAST	{GO-EAST, GO-SOUTH-EAST, GO-NORTH-EAST}
SOUTH-EAST	{GO-SOUTH-EAST, GO-SOUTH, GO-EAST}
SOUTH	{GO-SOUTH, GO-SOUTH-WEST, GO-SOUTH-EAST}
SOUTH-WEST	{GO-SOUTH-WEST, GO-WEST, GO-SOUTH}
WEST	{GO-WEST, GO-NORTH-WEST, GO-SOUTH-WEST}
NORTH-WEST	{GO-NORTH-WEST, GO-NORTH, GO-WEST}

## Risque

Fonction  $D^P : S^I \times S^P \mapsto [0 \dots 1]$  :

$$D_{ip}^P = Pr(\text{DETECTED} = T|i, p) = \frac{\mathcal{F}_{pi|XY}^P}{\sum_i \mathcal{F}_{pi|XY}^P}$$



## Récompense

$$R_i^I = \begin{cases} 1 & \text{si } i|_{XY} == \tau, \\ 0 & \text{sinon.} \end{cases}$$

## Principe

1. Agréger les différentes politiques du patrouilleur en formant un modèle de transitions global  $T^P$
2. Gérer les transitions entre les comportements via un modèle de transitions inter-politiques  $\theta^P$

## Modèle de transition agrégé

$$T^P : S^P \times \Pi^P \times S^P \mapsto [0, 1]$$
$$T_{p'k_p}^P = \mathcal{T}_{p'k_p}^P = Pr(p'|p, k_p)$$

## Modèle de transition inter-politique

$$\theta^P : \Pi^P \times \Pi^P \times S^I \times S^P$$
$$\theta_{k'kip}^P = Pr(k'|k, i, p)$$

## Fonction d'utilité

Compromis entre prise de risque et atteindre la cible :

$$U_{ip}^I = \beta R_i^I + (1 - \beta)(1 - D_{ip}^P)$$

## MDP augmenté résultant

Un MDP augmenté  $\langle S, A, T^I, T^P, \theta^P, U^I \rangle$

- ▶  $S : S^I \times S^P \times \Pi^P$
- ▶  $T^I : S^I \times A \times S^I \mapsto [0, 1]$

$$V_{ipk}^* = U_{ip}^I + \gamma \max_a \sum_{i'p'k'} T_{i'ai}^I T_{p'kp}^P \theta_{k'kip}^P V_{i'p'k'}^*$$

## Programme linéaire primal

$$\text{Minimize } \sum_{i,p,k} \mu_{ipk} V_{ipk}$$

Subject to

$$V_{ipk} \geq U_{ip}^I + \gamma \sum_{i',p',k'} T_{iai'}^I T_{pkp'}^P \theta_{k'kip} V_{i'p'k'}, \forall (a, i, p, k) \in A \times S^I \times S^P \times \Pi^P$$

$$\text{With } \sum_{i,p,k} \mu_{ipk} = 1$$

## Programme linéaire dual

$$\text{Maximize } \sum_{i,p,k} \sum_a U_{ip}^I Q_{ipka}$$

Subject to

$$\sum_a Q_{i'p'k'a} - \gamma \sum_{ipk} \sum_a T_{iai'}^I T_{pkp'}^P \theta_{k'kip} Q_{ipka} = \mu_{i'p'k'}, \forall (i', p', k') \in S^I \times S^P \times \Pi^P$$

$$\text{With } \sum_{i,p,k} \mu_{ipk} = 1$$

## Principe

1. On reprend les 2 principes précédents.
2. On construit une fonction de valeur pour le patrouilleur en considérant que le but du patrouilleur est de mettre en danger l'intrus.  $\Rightarrow$  Surestimation du risque.
3. On retranche cette fonction de valeur à celle de l'intrus.  
 $\Rightarrow$  pseudo-minmax.

## Le patrouilleur

- ▶ Fonction de récompense  $R_{pi}^P = D_{ip}^P$
- ▶ Fonction de valeur :

$$V_{pik}^{P*} = R_{pi}^P + \gamma \min_{a \in A} \sum_{i'p'k'} T_{i'ai}^I, T_{p'kp}^P \theta_{k'kip}^P V_{p'i'k'}^{P*}$$

## MDP augmenté résultant

Un MDP augmenté  $\langle S, A, T^I, T^P, \theta^P, U^I, V^{P*} \rangle$

►  $U_{ip}^I = R_i$

$$V_{ipk}^* = \beta \left[ U_{ip}^I + \gamma \max_a \sum_{i'p'k'} T_{i'ai}^I T_{p'kp}^P \theta_{k'kip}^P V_{i'p'k'}^* \right] - (1 - \beta) V_{pik}^{P*}$$

## Programme Linéaire

Adaptation du précédent.

## Question

Nash-optimalité ?

## Principe

1. On maintient une croyance sur la politique suivie par la défense
2. On considère que les couples  $(i, p)$  sont des observations sur la politique utilisée par le patrouilleur.

## Belief MDP

Un BMDP augmenté  $\langle S, A, T^I, T^P, \theta^P, U^I, \zeta_\pi \rangle$

▶  $S : S^I \times S^P \times B_\pi^P$

▶  $b'_\pi = \zeta_\pi(i, p, b_\pi)$

$$b'_\pi(k') = Pr(k'|i, p, b_\pi) = \sum_k b_\pi(k) \theta_{k'kip}^P$$

$$V_{ipb_\pi}^* = U_{ip}^I + \gamma \max_a \sum_{i'p} T_{i'ai}^I \left[ \sum_k T_{p'kp}^P b_\pi(k) \right] V_{i'p'\zeta_\pi(i,p,b_\pi)}^*$$

## Principe

1. On maintient une seconde croyance sur l'état du patrouilleur
2. Un oracle (satellite?) donne des observations sur l'état du patrouilleur

## Belief MDP

Un BMDP augmenté  $\langle S, A, T^I, T^P, \theta^P, U^I, \Omega, TO, \zeta_\pi, \zeta_s \rangle$

- ▶  $S : S^I \times B_s^P \times B_\pi^P$
- ▶  $U : U_{ib_s}^I = \beta \left[ 1 - \sum_p b_s(p) D_{ip}^P \right] + (1 - \beta) R_i^I$
- ▶  $\Omega : X \times Y \times O$
- ▶  $TO : \Omega \times S^P \times \Pi^P \Rightarrow Pr(o|p', k)$
- ▶  $\zeta_s : \Omega \times B_s^P \times B_\pi^P \Rightarrow Pr(p'|o, b_s, b_\pi), \forall p'$

$$V_{ib_s b_\pi}^* = U_{ib_s}^I + \gamma \max_a \sum_{i'} T_{i' ai}^I \sum_{op'k} TO_{op'k} \sum_p T_{p'kp}^P b_\pi(k) b_s(p) V_{i' \zeta_s(o, b_s, b_\pi) \zeta_\pi(i, p, b_\pi)}^*$$

---

## Algorithme 1 : Adaptation de l'algorithme POMCP

---

**Données :**  $i, h, \mathcal{G}, \mathcal{T}, \pi_{\text{rollout}}, \gamma, \epsilon, c$

**Resultat :**  $\pi^* = \text{Search}(i, h)$

**definir**  $\text{Search}(i, h)$ :

```

répéter
  | si  $h = \emptyset$ :
  |   |  $p \sim I_s$ 
  |   |  $k \sim I_\pi$ 
  |   | sinon:
  |   |   |  $p \sim \mathcal{B}_s(h)$ 
  |   |   |  $k \sim \mathcal{B}_\pi(h)$ 
  |   |   |  $\text{Simulate}(i, p, k, h, 0)$ 
jusqu'à  $\text{Timeout}()$ 
retourner  $\underset{a}{\text{argmax}} V(ha)$ 
    
```

**definir**  $\text{Rollout}(i, p, k, h, \text{depth})$ :

```

si  $\gamma^{\text{depth}} < \epsilon$ :
  | retourner 0
  |  $a \sim \pi_{\text{rollout}}(h, \cdot)$ 
  |  $i', p', k', o, u \sim \mathcal{G}(i, p, k, a)$ 
retourner  $u +$ 
  |  $\gamma \cdot \text{Rollout}(i', p', k', haio, \text{depth} +$ 
  | 1)
    
```

**definir**  $\text{Simulate}(i, p, k, h, \text{depth})$ :

```

si  $\gamma^{\text{depth}} < \epsilon$ :
  | retourner 0
si  $h \notin \mathcal{T}$ :
  | pour  $a$  dans  $A'$ :
  |   |  $\mathcal{T}(ha) \leftarrow$ 
  |   |   |  $(N_{\text{init}}(ha), V_{\text{init}}(ha), \emptyset, \emptyset)$ 
  |   | retourner
  |   |   |  $\text{Rollout}(i, p, k, h, \text{depth})$ 
a =  $\underset{b \in A'}{\text{argmax}} V(hb) + c \sqrt{\frac{\log N(h)}{N(hb)}}$ 
 $i', p', k', o, u \sim \mathcal{G}(i, p, k, a)$ 
 $U \leftarrow u +$ 
 $\gamma \cdot \text{Simulate}(i', p', k', haio, \text{depth} +$ 
1)
 $\mathcal{B}_s(h) \leftarrow \mathcal{B}_s(h) \cup \{p\}$ 
 $\mathcal{B}_\pi(h) \leftarrow \mathcal{B}_\pi(h) \cup \{k\}$ 
 $N(h) \leftarrow N(h) + 1$ 
 $N(ha) \leftarrow N(ha) + 1$ 
 $V(ha) \leftarrow V(ha) + \frac{U - V(ha)}{N(ha)}$ 
retourner  $U$ 
    
```

Jeux stochastiques et MDPs, késako ?

Jeu stochastique et surveillance militaire

Modélisations

**Premiers résultats**

Conclusion

## Instances

- ▶ Environnement : (9/10, 5/100)
  - ▶ 2 comportements pour patrouilleur :
    1. Patrouille en zigzags (boustrophédon)
    2. Protéger la cible
    3. Transitions : probabilité de détection  $> 0.3$
- ⇒ Taille espace d'état :  $|S| = 272n^5$
- ⇒ Nombre de transitions :  $54|S|$

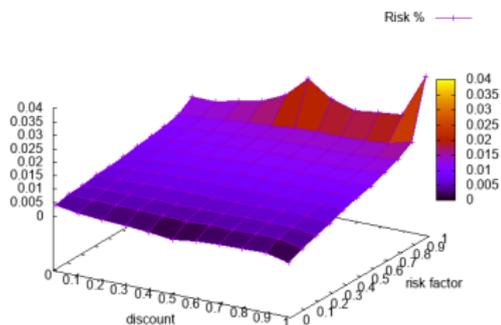
## Sélection

- ▶ *GridSearch* sur  $\gamma$  et  $\beta$
- ▶  $< \text{risque}$  et  $> \text{cible}$

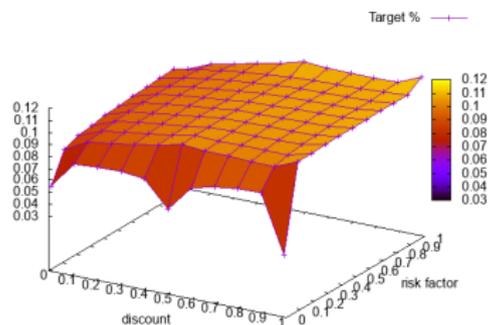
## Évaluation

- ▶ On simule pendant 30 min :
  1. %interceptions
  2. %risque
  3. %cible
  4. %utilité

## Grille 8 × 8 :

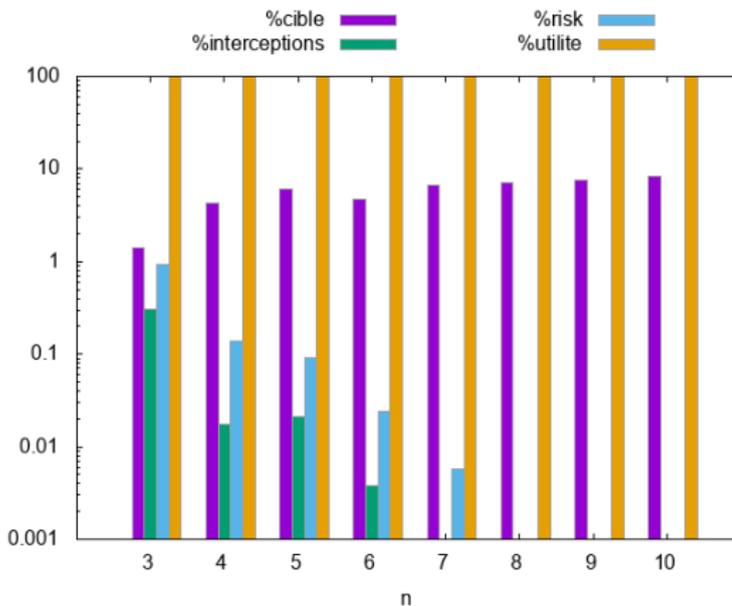


Risque(%) en fonction des paramètres libres

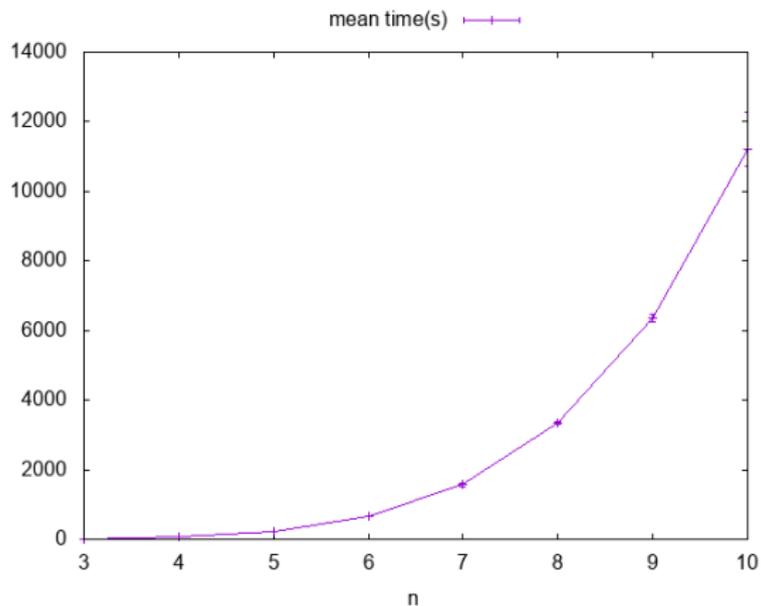


Cible (%) en fonction des paramètres libres

⇒ Compromis ( $\gamma = 0.99, \beta = 1e - 4$ ).



Métriques, modèle observabilité complète



Temps de planification, modèle observabilité complète

Jeux stochastiques et MDPs, késako ?

Jeu stochastique et surveillance militaire

Modélisations

Premiers résultats

Conclusion

- ▶ Première modélisation « simple » du problème
- ▶ Résultats encourageants pour l'observabilité totale
- ▶ Bémol : fonction de valeur non-convexe
- ▶ Propose un algorithme pour l'observabilité partielle
- ▶ Implémentation + expériences à compléter
- ▶ Décomposition structure  $\Rightarrow$  passage à l'échelle

# Questions ?

