#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$
- The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- ightarrow a set of observations  $Z=O=\{z_1,\ldots,z_{|Z|}\}$
- Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t|s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $\mathit{O}(s_t, a_{t-1}, z_t) = \mathit{P}(z_t | s_t, a_{t-1})$
- Reward(=utility) R(s, a) for each state (and action).
- (discount factor  $\gamma \in <$  0, 1 >

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} \gamma^{t} R(s_{t}, a_{t}) \right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- ightarrow a set of observations  $Z=O=\{z_1,\ldots,z_{|Z|}\}$
- Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t|s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $\mathit{O}(s_t, a_{t-1}, z_t) = \mathit{P}(z_t | s_t, a_{t-1})$
- Reward(=utility) R(s, a) for each state (and action).
- (discount factor  $\gamma \in <$  0, 1 >

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} \gamma^{t} R(s_{t}, a_{t}) \right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
  - ightarrow a set of observations  $Z = O = \{z_1, \dots, z_{|Z|}\}$
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t|s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $\mathit{O}(s_t, a_{t-1}, z_t) = \mathit{P}(z_t|s_t, a_{t-1})$
- Reward(=utility) R(s, a) for each state (and action).
- (discount factor  $\gamma \in <$  0, 1 >

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} \gamma^{t} R(s_{t}, a_{t}) \right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- $\rightarrow\,$  a set of observations  $Z={\it O}=\{z_1,\ldots,z_{|Z|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
  - ightarrow observation matrix  $\mathit{O}(s_t, a_{t-1}, z_t) = \mathit{P}(z_t | s_t, a_{t-1})$
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <$  0, 1>

We maximize the expected cumulative reward  $max_{\pi}\mathbb{E}_{\pi}\left[\sum_{t=1}^{\infty}\gamma^{t}R(s_{t},a_{t})\right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- $\rightarrow\,$  a set of observations  $Z={\it O}=\{z_1,\ldots,z_{|Z|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
  - ightarrow observation matrix  $O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <$  0, 1>

We maximize the expected cumulative reward  $max_{\pi}\mathbb{E}_{\pi}\left[\sum_{t=1}^{\infty}\gamma^{t}R(s_{t},a_{t})\right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- $\rightarrow\,$  a set of observations  $Z=\mathit{O}=\{z_1,\ldots,z_{|\mathcal{Z}|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$ 
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <$  0, 1 >)

We maximize the expected cumulative reward  $max_{\pi}\mathbb{E}_{\pi}\left[\sum_{t=1}^{\infty}\gamma^{t}R(s_{t},a_{t})\right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- ightarrow a set of observations  $Z=O=\{z_1,\ldots,z_{|Z|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $\mathit{O}(\mathit{s}_t, \mathit{a}_{t-1}, \mathit{z}_t) = \mathit{P}(\mathit{z}_t | \mathit{s}_t, \mathit{a}_{t-1})$ 
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <$  0, 1 >

We maximize the expected cumulative reward  $max_{\pi}\mathbb{E}_{\pi}\left[\sum_{t=1}^{\infty}\gamma^{t}R(s_{t},a_{t})\right]$ .

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- ightarrow a set of observations  $Z=O=\{z_1,\ldots,z_{|Z|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$ 
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <0, 1>$ ).

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} [\sum_{t=1}^{\infty} \gamma^t R(s_t, a_t)].$ 

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- $\rightarrow\,$  a set of observations  ${\it Z}={\it O}=\{{\it z}_1,\ldots,{\it z}_{|{\it Z}|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$ 
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <0, 1>$ ).

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} [\sum_{t=1}^{\infty} \gamma^t R(s_t, a_t)].$ 

#### MDP The policy is a function of the state $\pi(s)$

MDP The policy is a function of the history  $\pi(a_{t-1}, z_{t-1}, \dots, z_1, a_0, b_0)$ 

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- ightarrow a set of observations  $Z=O=\{z_1,\ldots,z_{|Z|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$ 
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <0, 1>$ ).

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} [\sum_{t=1}^{\infty} \gamma^t R(s_t, a_t)].$ 

#### MDP The policy is a function of the state $\pi(s)$ MDP The policy is a function of the history $\pi(a_{t-1}, z_{t-1}, \dots, z_1, a_0, b_0)$ • or a function of the belief: $b: S \to (0, 1), \pi(b)$

#### Definition (Partially Observed Markov Decision Processes POMDP)

Partially Observed Markov Decision Processes is defined by:

- Finite set of states S,  $S_i = S$  for any time  $t \in \mathbb{N}_0$ ,
- $\rightarrow$  Initial **belief**  $b_0(s) = P(S_0)$ 
  - The set of possible actions (decisions) at any time  $A = \{a_1, \ldots, a_{|A|}\}$
- ightarrow a set of observations  $Z=O=\{z_1,\ldots,z_{|Z|}\}$ 
  - Transition matrix  $T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$
- ightarrow observation matrix  $O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$ 
  - Reward(=utility) R(s, a) for each state (and action).
  - (discount factor  $\gamma \in <0, 1>$ ).

We maximize the expected cumulative reward  $\max_{\pi} \mathbb{E}_{\pi} [\sum_{t=1}^{\infty} \gamma^t R(s_t, a_t)].$ 

#### MDP The policy is a function of the state $\pi(s)$ MDP The policy is a function of the history $\pi(a_{t-1}, z_{t-1}, \dots, z_1, a_0, b_0)$ • or a function of the **belief**: $b : S \to (0, 1)$ , $\pi(b)$

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$ .
- The reward R is a function of the state and the action
  - $\sim U(gold, l/r) = 10, U(tiger, l/r) = -100, U(s, listen) = -1, that is s$

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right *S* = {*left*, *right*}
- We may open any door or listen A = {left, right, listen},
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$ .
- The reward R is a function of the state and the action
  - $\sim U(gold, l/r) = 10, U(tiger, l/r) = -100, U(*, listen) = -1, that is$

The example is a variant of the Monty Hall problem.

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right S = {left, right}
- We may open any door or listen A = {left, right, listen},
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:

 $\sim$  the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$ 

- The reward R is a function of the state and the action
  - $\sim U(\mathit{gold}, l/r) = 10, U(\mathit{tiger}, l/r) = -100, U(*, \mathit{listen}) = -1,$  that is

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - $\Rightarrow$  the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward R is a function of the state and the action
  - $\sim U(gold, l/r) = 10, U(tiger, l/r) = -100, U(s, listen) = -1, that is s$

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
- The reward *R* is a function of the state and the action
  - $\sim U(gold, l/r) = 10, U(tiger, l/r) = -100, U(s, listen) = -1, that is$

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
- The reward *R* is a function of the state and the action
  - $\sim U(gold, l/r) = 10, U(tiger, l/r) = -100, U(s, listen) = -1, that is s$

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward R is a function of the state and the action
  - $\sim U(gold, l/r) = 10, U(tiger, l/r) = -100, U(*, listen) = -1, that is$

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we **listen** we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward *R* is a function of the state and the action
  - U(gold, l/r) = 10, U(tiger, l/r) = -100, U(\*, lister) = -1, that is

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward *R* is a function of the state and the action
  - U(gold, l/r) = 10, U(tiger, l/r) = -100, U(\*, listen) = -1, that is

The example is a variant of the Monty Hall problem.

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward R is a function of the state and the action

• U(gold, l/r) = 10, U(tiger, l/r) = -100, U(\*, listen) = -1, that is

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward R is a function of the state and the action
  - U(gold, l/r) = 10, U(tiger, l/r) = -100, U(\*, listen) = -1, that is

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 \rangle$
- The reward R is a function of the state and the action
  - U(gold, l/r) = 10, U(tiger, l/r) = -100, U(\*, listen) = -1, that is

- We face two doors.
  - There is a tiger behind one door,
  - there is a gold brick behind the other.
- The Tiger is left or right  $S = \{left, right\}$
- We may open any door or listen  $A = \{left, right, listen\},\$
- we search optimal policy for given observation and reward tables.
- We observe Z only if we listen we listen the tiger left TL or right TR
- we reset the world at the beginning and after opening any door:
  - the initial belief  $P(S_0) = \langle 0.5, 0.5 
    angle$
- The reward R is a function of the state and the action
  - U(gold, l/r) = 10, U(tiger, l/r) = -100, U(\*, listen) = -1, that is

Tiger Action	left	right	Z   S=?, A=listen	left	right
Listen	-1	-1	TL	0.85	0.15
left	-100	10	TR	0.15	0.85
right	10	-100	NoInfo	0	0

Finite horizon POMPD *t*,  $\gamma = 1$ :

• 
$$t = 1$$
  
 $EU_{t=1}(A = left/right) = \frac{-100+10}{2} = -45$   
 $EU_{t=1}(A = listen) = -1$ 

• horizon t = 2

$$T(s_{t-1}, a_{t-1}, s_t) = P(s_t | s_{t-1}, a_{t-1})$$
$$O(s_t, a_{t-1}, z_t) = P(z_t | s_t, a_{t-1})$$
• and  $t = 4$ 

### POMPD

Finite horizon POMPD *t*,  $\gamma = 1$ :

• 
$$t = 1$$
  
 $EU_{t=1}(A = left/right) = \frac{-100+10}{2} = -45$   
 $EU_{t=1}(A = listen) = -1$ 

• horizon 
$$t = 2$$



and t = 4

### POMPD

Finite horizon POMPD *t*,  $\gamma = 1$ :

• 
$$t = 1$$
  
 $EU_{t=1}(A = left/right) = \frac{-100+10}{2} = -45$   
 $EU_{t=1}(A = listen) = -1$   
• horizon  $t = 2$ 

$$\begin{array}{lll} T(s_{t-1},a_{t-1},s_t) & = & P(s_t|s_{t-1},a_{t-1}) \\ O(s_t,a_{t-1},z_t) & = & P(z_t|s_t,a_{t-1}) \end{array}$$





#### • $\gamma = 0.75$

- we iterate until convergence
- Then, we create a graph by joining two successive time slices together.
- We may omit nodes that are not reachable from the initial belief  $b_0(s) = 0.5$ .

- $\gamma = 0.75$
- we iterate until convergence
- Then, we create a graph by joining two successive time slices together.
- We may omit nodes that are not reachable from the initial belief  $b_0(s) = 0.5$ .



- $\gamma = 0.75$
- we iterate until convergence
- Then, we create a graph by joining two successive time slices together.
- We may omit nodes that are not reachable from the initial belief  $b_0(s) = 0.5$ .





Figure 16: Policy graph for tiger example

•  $\gamma = 0.75$ 

- we iterate until convergence
- Then, we create a graph by joining two successive time slices together.
- We may omit nodes that are not reachable from the initial belief  $b_0(s) = 0.5$ .





Figure 16: Policy graph for tiger example



Figure 17: Trimmed policy graph for tiger example

#### • The history is aggregated in the probability distribution over states

- history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
- **belief**  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0),$
- initial belief  $b_0(s) = P(S_0)$ .
- In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - a transition we eliminate unobserved spectral
  - $\sim$  an observation we condition by  $z_{t^{-}}$
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{|}) \\ &= \frac{\sum_s O(s^{|}, a_{t-1}, z_t) T(s, a_{t-1}, s^{|}) b_{t-1}(s)}{P(z_t | b_{t-1}, a_{t-1})} \end{aligned}$$

#### • The history is aggregated in the probability distribution over states

- history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
- **belief**  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0),$
- initial belief  $b_0(s) = P(S_0)$ .
- In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - $\sim$  a transition we eliminate unobserved  $s_{
    m c}$
  - an observation we condition by  $z_{t^{-}}$
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\dagger}) \\ &= \frac{\sum_s O(s^{\dagger}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\dagger}) b_{t-1}(s)}{P(z_t | b_{t-1}, a_{t-1})} \end{aligned}$$

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - belief  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - $\sim$  a transition we eliminate unobserved  $s_{
    m constraint}$
  - $\sim$  an observation we condition by  $z_{t^{-}}$
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{|}) \\ &= \frac{\sum_s O(s^{|}, a_{t-1}, z_t) T(s, a_{t-1}, s^{|}) b_{t-1}(s)}{P(z_t | b_{t-1}, a_{t-1})} \end{aligned}$$

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - belief  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - $\sim$  a transition we eliminate unobserved  $s_{\rm f}$
  - an observation we condition by z<sub>i</sub>.
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\dagger}) \\ &= \frac{\sum_s O(s^{\dagger}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\dagger}) b_{t-1}(s)}{P(z_t | b_{t-1}, a_{t-1})} \end{aligned}$$

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - belief  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - a transition we eliminate unobserved  $s_{t-1}$
  - an observation we condition by z<sub>t</sub>.
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{|}) \\ &= \frac{\sum_s O(s^{|}, a_{t-1}, z_t) T(s, a_{t-1}, s^{|}) b_{t-1}(s)}{P(z_t | b_{t-1}, a_{t-1})} \end{aligned}$$
- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - **belief**  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).

#### • We update belief after any iteration. The update consists of:

- a transition we eliminate unobserved s<sub>t-1</sub>
- an observation we condition by  $z_t$ .
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\mid}) \\ &= \frac{\sum_s O(s^{\mid}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\mid}) b_{t-1}(s)}{P(z_t \mid b_{t-1}, a_{t-1})} \end{aligned}$$

• Markov with respect to b since au does not depend on time.

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - **belief**  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - a transition we eliminate unobserved s<sub>t-1</sub>
  - an observation we condition by  $z_t$ .
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\mid}) \\ &= \frac{\sum_s O(s^{\mid}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\mid}) b_{t-1}(s)}{P(z_t \mid b_{t-1}, a_{t-1})} \end{aligned}$$

Markov with respect to b since τ does not depend on time.

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - **belief**  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - a transition we eliminate unobserved  $s_{t-1}$
  - an observation we condition by  $z_t$ .

belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\mid}) \\ &= \frac{\sum_s O(s^{\mid}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\mid}) b_{t-1}(s)}{P(z_t \mid b_{t-1}, a_{t-1})} \end{aligned}$$

• Markov with respect to b since  $\tau$  does not depend on time.

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - belief  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - a transition we eliminate unobserved  $s_{t-1}$
  - an observation we condition by  $z_t$ .
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\mid}) \\ &= \frac{\sum_s O(s^{\mid}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\mid}) b_{t-1}(s)}{P(z_t \mid b_{t-1}, a_{t-1})} \end{aligned}$$

19. prosince 2024 14 / 1 - 112

- The history is aggregated in the probability distribution over states
  - history  $h_t = \{a_0, z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t\}$
  - **belief**  $b_t(s) = P(S = s | z_t, a_{t-1}, \dots, a_0, b_0)$ ,
  - initial belief  $b_0(s) = P(S_0)$ .
  - In the tiger example a single number b(left), since the other probability is 1 b(left).
- We update belief after any iteration. The update consists of:
  - a transition we eliminate unobserved  $s_{t-1}$
  - an observation we condition by  $z_t$ .
- belief update

$$\begin{aligned} \tau(b_{t-1}, a_{t-1}, z_t) &= b_t(s^{\mid}) \\ &= \frac{\sum_s O(s^{\mid}, a_{t-1}, z_t) T(s, a_{t-1}, s^{\mid}) b_{t-1}(s)}{P(z_t \mid b_{t-1}, a_{t-1})} \end{aligned}$$

• Markov with respect to b since  $\tau$  does not depend on time.

#### • Strategy (policy) is a function $\pi(b) \rightarrow a$ ,

optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = argmax_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

#### value function

- initial  $V_0(b) = max_p \sum_{s \in S} R(s,s)b(s)$ .
- recursively

 $V_{1}(b) = \max_{s} \left[ \sum_{s \in S} R(s, s) b(s) + \gamma \sum_{s \in S} P(z|s, b) V_{t-1}(\tau(b, s, z)) \right],$ 

• optimal strategy for the horizon t:

 $\pi_t^*(b) = \operatorname{argmax}_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right].$ 

- Strategy (policy) is a function  $\pi(b) \rightarrow a$ ,
- optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

- value function
  - initial  $V_0(b) = max_s \sum_{s \in S} R(s, s) b(s)$
  - recursively
    - $V_i(b) = max_i \sum_{s \in S} R(s, s)b(s) + \gamma \sum_{s \in S} P(z|s, b)V_{i-1}(r(b, s, z)) | i$
- optimal strategy for the horizon t:
  - $\pi_t^*(b) = \operatorname{argmax}_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right].$

- Strategy (policy) is a function  $\pi(b) \rightarrow a$ ,
- optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

#### • value function

- initial  $V_0(b) = max_a \sum_{s \in S} R(s, a)b(s)$
- recursively

 $V_t(b) = \max_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right],$ 

• optimal strategy for the horizon t:

 $\pi_t^*(b) = \operatorname{argmax}_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right].$ 

- Strategy (policy) is a function  $\pi(b) \rightarrow a$ ,
- optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

#### • value function

- initial  $V_0(b) = max_a \sum_{s \in S} R(s, a)b(s)$
- recursively

 $V_t(b) = \max_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right],$ 

optimal strategy for the horizon t

 $\pi_*^*(b) = \operatorname{argmax}_a\left[\sum_{s \in S} R(s, a)b(s) + \gamma \sum_{z \in Z} P(z|a, b)V_{t-1}(\tau(b, a, z))\right].$ 



- Strategy (policy) is a function  $\pi(b) \rightarrow a$ ,
- optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

#### • value function

• initial 
$$V_0(b) = max_a \sum_{s \in S} R(s, a)b(s)$$

• recursively  $V_t(b) = max_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right],$ 

• optimal strategy for the horizon t

 $\pi_*^*(b) = \operatorname{argmax}_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right].$ 



• Strategy (policy) is a function  $\pi(b) \rightarrow a$ ,

optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

#### value function

- initial  $V_0(b) = max_a \sum_{s \in S} R(s, a)b(s)$
- recursively

$$V_t(b) = \max_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right],$$

• optimal strategy for the horizon t:

 $\pi_t^*(b) = \operatorname{argmax}_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right].$ 



- Strategy (policy) is a function  $\pi(b) \rightarrow a$ ,
- optimal strategy maximizes the expected discounted cumulative reward

$$\pi^*(b_0) = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (\gamma^t \cdot r_t) | b_0 
ight]$$

#### • value function

- initial  $V_0(b) = max_a \sum_{s \in S} R(s, a)b(s)$
- recursively  $V_t(b) = max_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right],$
- optimal strategy for the horizon *t*:

$$\pi_t^*(b) = \operatorname{argmax}_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right].$$





### $|\Gamma_t| = O(|A| \cdot |\Gamma_{t-1}|^{|Z|})$

• value function  $V_t(b)$  can be represented by a finite number of hyperplanes

• each hyperplane is represented as a vector  $\alpha$ :

$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}$$

• initial:  $\Gamma_0(b) = \{ \langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle \}_{a \in A}$ 

• From

 $= V_1(b) = max_{\theta} \left[ \sum_{n \in \mathcal{N}} R(a, n)b(n) + \gamma \sum_{n \in \mathcal{N}} P(a|a, b) V_{\theta-1}(n(b, a, a)) \right]_{\theta}$ =  $r(b_0, a_0, c_{0,1}) = \frac{2\pi i \Omega(a_0, a_0, a_1) \Omega(a_0, a_2)}{\Omega(a_0, a_1) \Omega(a_0, a_2)}$ 



 $V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}.$ 

• initial:  $\Gamma_0(b) = \{\langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle\}_{a \in A}$ 

• at the time t:  $V_t(b) = \max_{\alpha \in \Gamma_t} \sum_{s \in S} \alpha(s) b(s)$ .

From

 $= \mathcal{N}_{0}(b) = \max_{i} \left[ \sum_{a \in \mathcal{A}} \mathcal{R}(a, a)b(a) + \gamma \sum_{a \in \mathcal{A}} \mathcal{P}(a|a, b)\mathcal{N}_{i-1}(x(b, a, a)) \right]$  $= \mathcal{R}(b_{i}, a_{i}, a_{i}) = \frac{2i\partial(a_{i}, a_{i}, a_{i})\partial(a_{i}, a_{i}, b)}{\partial(a_{i}, a_{i}, a_{i})} = \frac{2i\partial(a_{i}, a_{i}, a_{i})\partial(a_{i}, a_{i}, a_{i})}{\partial(a_{i}, a_{i}, a_{i})}$ 



$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}$$

- initial:  $\Gamma_0(b) = \{ \langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle \}_{a \in A}$
- at the time t:  $V_t(b) = \max_{\alpha \in \Gamma_t} \sum_{s \in S} \alpha(s)b(s)$

From

 $= \mathcal{O}(b) = \max_{i \in \mathcal{O}} \mathbb{E}[\max_{i \in \mathcal{O}} \mathcal{P}(a, a)b(a) + \gamma \sum_{i \in \mathcal{O}} \mathcal{P}(a, a, b)V_{i-1}(\tau(b, a, c))] = \frac{2\pi i \mathcal{O}(\tau_{i}, a_{i}, a_{i})}{m_{i}} \frac{1}{m_{i}} \frac{1}{m_{$ 



$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}.$$

- initial:  $\Gamma_0(b) = \{\langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle\}_{a \in A}$
- at the time t:  $V_t(b) = \max_{\alpha \in \Gamma_t} \sum_{s \in S} \alpha(s) b(s)$ .

From

 $= V_t(b) = \max_{\theta} \left[ \sum_{s \in S} R(s, a)b(s) + \gamma \sum_{s \in Z} P(z|a, b)V_{t-1}(\tau(b, a, z)) \right]$   $= \tau(b_0, \sigma_0, \varepsilon_{0,1}) = \frac{\sum_{s \in S} R(s, a)b(s)}{P(a, 1)b(a)}$ 



$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}$$

- initial:  $\Gamma_0(b) = \{\langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle\}_{a \in A}$
- at the time t:  $V_t(b) = \max_{\alpha \in \Gamma_t} \sum_{s \in S} \alpha(s) b(s)$ .

From

• 
$$V_t(b) = \max_a \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{z \in Z} P(z|a, b) V_{t-1}(\tau(b, a, z)) \right]$$
  
•  $\tau(b_t, a_t, z_{t+1}) = \frac{\sum_s O(s^{\parallel}, a_t, z_{t+1}) T(s, a_t, s^{\parallel}) b_t(s)}{Pr(z_{t+1}|b_t, a_t)}$ 



$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}$$

- initial:  $\Gamma_0(b) = \{ \langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle \}_{a \in A}$ • at the time  $t: V_t(b) = max_{\alpha \in \Gamma_t} \sum_{e \in S} \alpha(s)b(s).$
- From

• 
$$V_t(b) = \max_a \left[ \sum_{s \in S} R(s, a)b(s) + \gamma \sum_{z \in Z} P(z|a, b)V_{t-1}(\tau(b, a, z)) \right]:$$
  
• 
$$\tau(b_t, a_t, z_{t+1}) = \frac{\sum_s O(s^{t}, a_t, z_{t+1}) T(s, a_t, s^{t})b_t(s)}{P(z_{t+1}|b_t, a_t)}$$



$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \ldots, \alpha_m\} = \{a; v_0, v_1, \ldots, v_m\}.$$

- initial:  $\Gamma_0(b) = \{ \langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle \}_{a \in A}$ • at the time  $t: V_t(b) = max_{\alpha \in \Gamma_t} \sum_{e \in S} \alpha(s)b(s).$
- From

• 
$$V_t(b) = max_a \left[ \sum_{s \in S} R(s, a)b(s) + \gamma \sum_{z \in Z} P(z|a, b)V_{t-1}(\tau(b, a, z)) \right]$$
  
•  $\tau(b_t, a_t, z_{t+1}) = \frac{\sum_s O(s^{|}, a_t, z_{t+1})T(s, a_t, s^{|})b_t(s)}{Pr(z_{t+1}|b_t, a_t)}$ 

#### $\alpha$ vectors



$$|\Gamma_t| = O(|A| \cdot |\Gamma_{t-1}|^{|Z|})$$

value function V<sub>t</sub>(b) can be represented by a finite number of hyperplanes
 each hyperplane is represented as a vector α:

$$V_t(b) \Leftrightarrow \Gamma_t = \{a; \alpha_0, \alpha_1, \dots, \alpha_m\} = \{a; v_0, v_1, \dots, v_m\}$$
  
• initial:  $\Gamma_0(b) = \{\langle a; R(s_1, a), R(s_2, a), \dots, R(s_{|S|}, a) \rangle\}_{a \in A}$   
• at the time  $t$ :  $V_t(b) = max_{\alpha \in \Gamma_t} \sum_{a \in S} \alpha(s)b(s)$ .

From

• 
$$V_t(b) = max_a \left[ \sum_{s \in S} R(s, a)b(s) + \gamma \sum_{z \in Z} P(z|a, b)V_{t-1}(\tau(b, a, z)) \right]$$
:  
•  $\tau(b_t, a_t, z_{t+1}) = \frac{\sum_s O(s^{|}, a_t, z_{t+1})T(s, a_t, s^{|})b_t(s)}{Pr(z_{t+1}|b_t, a_t)}$ 

$$V_t(b) = \max_a \left| \sum_{s \in S} R(s, a) b(s) \right|$$

+ 
$$\gamma \sum_{z \in Z} \max_{\alpha \in \Gamma_{t-1}} \sum_{s' \in S} \sum_{s \in S} T(s, a, s') O(s', a, z) \alpha(s') b(s)$$

٦

• temporal sets  $\forall \alpha_i \in \Gamma_{t-1}$  ( $j \in S$  is the  $\alpha$  dimension):

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(j) = R(j,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(j) = \gamma \sum_{s' \in S} T(j,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• The utility for the action a summed over possible observation results  $z_k$ :

$$\Gamma_t^a = \Gamma_t^{a,+} + \Gamma_t^{a,z_1} \oplus \Gamma_t^{a,z_2} \oplus \ldots \oplus \Gamma_t^{a,z_m}$$

• the new value function for the time t:  $\Gamma_t \leftarrow \bigcup_{a \in A} \Gamma_t^a$ .

• We remove all  $\alpha$  that are dominated by others

> or to avoid to generate many of them at all  $|\mathbb{I}_t| = O(|A| \cdot |\mathbb{I}_{t-1}|^{|Z|})$ 

• temporal sets  $\forall \alpha_i \in \Gamma_{t-1}$  ( $j \in S$  is the  $\alpha$  dimension):

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(j) = R(j,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(j) = \gamma \sum_{s' \in S} T(j,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• The utility for the action a summed over possible observation results  $z_k$ :

$$\Gamma_t^{a} = \Gamma_t^{a,+} + \Gamma_t^{a,z_1} \oplus \Gamma_t^{a,z_2} \oplus \ldots \oplus \Gamma_t^{a,z_m}$$

• the new value function for the time  $t: \Gamma_t \leftarrow \bigcup_{a \in A} \Gamma_t^a$ .

• We remove all  $\alpha$  that are dominated by others

• temporal sets  $\forall \alpha_i \in \Gamma_{t-1}$  ( $j \in S$  is the  $\alpha$  dimension):

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(j) = R(j,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(j) = \gamma \sum_{s' \in S} T(j,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• The utility for the action a summed over possible observation results  $z_k$ :

$$\Gamma_t^{a} = \Gamma_t^{a,+} + \Gamma_t^{a,z_1} \oplus \Gamma_t^{a,z_2} \oplus \ldots \oplus \Gamma_t^{a,z_m}$$

the new value function for the time t: Γ<sub>t</sub> ← ⋃<sub>a∈A</sub> Γ<sup>a</sup><sub>t</sub>.
We remove all α that are dominated by others

• temporal sets  $\forall \alpha_i \in \Gamma_{t-1}$  ( $j \in S$  is the  $\alpha$  dimension):

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(j) = R(j,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(j) = \gamma \sum_{s' \in S} T(j,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• The utility for the action a summed over possible observation results  $z_k$ :

$$\Gamma_t^{a} = \Gamma_t^{a,+} + \Gamma_t^{a,z_1} \oplus \Gamma_t^{a,z_2} \oplus \ldots \oplus \Gamma_t^{a,z_m}$$

- the new value function for the time  $t: \Gamma_t \leftarrow \bigcup_{a \in A} \Gamma_t^a$ .
- We remove all  $\alpha$  that are dominated by others
  - there are strategies to remove them earlier
  - or to avoid to generate many of them at all  $|\Gamma_t| = O(|A| \cdot |\Gamma_{t-1}|^{|Z|})$

• temporal sets  $\forall \alpha_i \in \Gamma_{t-1}$  ( $j \in S$  is the  $\alpha$  dimension):

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(j) = R(j,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(j) = \gamma \sum_{s' \in S} T(j,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• The utility for the action a summed over possible observation results  $z_k$ :

$$\Gamma_t^{a} = \Gamma_t^{a,+} + \Gamma_t^{a,z_1} \oplus \Gamma_t^{a,z_2} \oplus \ldots \oplus \Gamma_t^{a,z_m}$$

- the new value function for the time  $t: \Gamma_t \leftarrow \bigcup_{a \in A} \Gamma_t^a$ .
- We remove all  $\alpha$  that are dominated by others
  - there are strategies to remove them earlier
  - or to avoid to generate many of them at all  $|\Gamma_t| = O(|A| \cdot |\Gamma_{t-1}|^{|Z|})$

• temporal sets  $\forall \alpha_i \in \Gamma_{t-1}$  ( $j \in S$  is the  $\alpha$  dimension):

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(j) = R(j,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(j) = \gamma \sum_{s' \in S} T(j,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• The utility for the action a summed over possible observation results  $z_k$ :

$$\Gamma_t^{a} = \Gamma_t^{a,+} + \Gamma_t^{a,z_1} \oplus \Gamma_t^{a,z_2} \oplus \ldots \oplus \Gamma_t^{a,z_m}$$

- the new value function for the time  $t: \Gamma_t \leftarrow \bigcup_{a \in A} \Gamma_t^a$ .
- We remove all  $\alpha$  that are dominated by others
  - there are strategies to remove them earlier
  - or to avoid to generate many of them at all  $|\Gamma_t| = O(|A| \cdot |\Gamma_{t-1}|^{|Z|})$ .

#### POMDP!

1: procedure POLICY POMDP(T)  $\Gamma = \{[None; 0, ..., 0]\}$ 2: for t = 1 to T do 3:  $\Gamma' = \emptyset$ 4: for all  $(a', v_1^k, \ldots, v_N^k) \in \Gamma$ , all  $a \in A$ , all  $z \in O$  do 5: for i = 1 to N do  $\triangleright$  for all states  $s_i$ 6:  $v_{\mu,z,i}^{k} = \gamma \sum_{i=1}^{N} v_{i}^{k} p(z|s_{i}) p(s_{i}|a,s_{i})$ 7: end for 8: end for ▷ next: free choice  $k(z) \in \Gamma$  for every z 9: for all *u*, all k(1), ..., k(M) = (1, ..., 1) to  $(|\Gamma|, ..., |\Gamma|)$  do 10: for i = 1 to N do 11:  $v'_i = \left[ r(s_j, a) + \sum_z v^{k(z)}_{a, z, j} \right]$ ho aggregate over  $z \in O$ 12: end for 13. add  $(a; v'_1, \ldots, v'_N)$  to  $\Gamma'$  $\triangleright$  one  $\Gamma_{t}^{a}$  element 14: end for 15: 16: (optimally prune  $\Gamma'$ )  $\Gamma = \Gamma'$ 17: end for 18: 19: return  $\Gamma$ 20. end procedure

# POMPD policy!

• From the set of  $\alpha$  vectors  $\Gamma$  we get the policy:

#### policy

- 1: procedure POLICY\_POMDP( $\Gamma, b = [b_1, \dots, b_N]$ )
- 2:  $\hat{a} = \arg \max_{(a;v_1^k,\ldots,v_N^k) \in \Gamma} \sum_{i=1}^N v_i^k b_i$
- 3: return â
- 4: end procedure



Figure 15.7 The benefit of point-based value iteration over general value iteration: Shown in (a) is the exact value function at horizon T = 30 for a different example, which consists of 120 constraints, after pruning. On the right is the result of the PBVI algorithm retaining only 11 linear functions. Both functions yield virtually indistinguisable results when applied to control.



Figure 17: Trimmed policy graph for tiger example

- $\bullet~\times$  for the rest of today lecture
- Pineau & all.: Anytime Point-Based Approximations for Large POMDPs, JAIR 2006
- Pearl the Nursebot
- Find a person





(a) t=1

(b) t=7







(b) Reconstruction

c) t=29

d) t=17

c) t=12

3

#### • We evaluate the *belief* only in a finite number of points

• only one  $\alpha$  vector for each point

$$\begin{aligned} &\Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(s) = R(s,a) \\ &\Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• max for FINITE number of  $b \in B$ 

$$\alpha_{b} = \operatorname{argmax}_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \sum_{z \in Z} \operatorname{argmax}_{\alpha \in \Gamma_{t}^{a, z}} \sum_{s \in S} \alpha(s) b(s) \right]$$
  
$$\Gamma_{t} = \bigcup_{b \in B} \{\alpha_{b}\}$$

- We evaluate the *belief* only in a finite number of points
- $\bullet\,$  only one  $\alpha\,$  vector for each point

$$\begin{split} & \Gamma_t^{a,+} & \leftarrow & \alpha^{a,+}(s) = R(s,a) \\ & \Gamma_t^{a,z} & \leftarrow & \alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s,a,s') O(s',a,z) \alpha(s'), \end{split}$$

• max for FINITE number of  $b \in B$ 

$$\alpha_{b} = \operatorname{argmax}_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \sum_{z \in Z} \operatorname{argmax}_{\alpha \in \Gamma_{t}^{a, z}} \sum_{s \in S} \alpha(s) b(s) \right]$$
  
$$\Gamma_{t} = \bigcup_{b \in B} \{\alpha_{b}\}$$

- We evaluate the *belief* only in a finite number of points
- $\bullet\,$  only one  $\alpha\,$  vector for each point

$$\begin{aligned} \Gamma^{a,+}_t &\leftarrow \alpha^{a,+}(s) = R(s,a) \\ \Gamma^{a,z}_t &\leftarrow \alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s,a,s') O(s',a,z) \alpha(s'), \end{aligned}$$

• max for FINITE number of  $b \in B$ 

$$\alpha_{b} = \operatorname{argmax}_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \sum_{z \in Z} \operatorname{argmax}_{\alpha \in \Gamma_{t}^{s, z}} \sum_{s \in S} \alpha(s) b(s) \right]$$
  
$$\Gamma_{t} = \bigcup_{b \in B} \{\alpha_{b}\}$$

- We evaluate the *belief* only in a finite number of points
- $\bullet\,$  only one  $\alpha\,$  vector for each point

$$\begin{split} & \Gamma_t^{a,+} & \leftarrow & \alpha^{a,+}(s) = R(s,a) \\ & \Gamma_t^{a,z} & \leftarrow & \alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s,a,s') O(s',a,z) \alpha(s'), \end{split}$$

• max for FINITE number of  $b \in B$ 

$$\alpha_{b} = \operatorname{argmax}_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \sum_{z \in Z} \operatorname{argmax}_{\alpha \in \Gamma_{t}^{s, z}} \sum_{s \in S} \alpha(s) b(s) \right]$$
  
$$\Gamma_{t} = \bigcup_{b \in B} \{\alpha_{b}\}$$

## POMDP Evaluation for the Fixed Number of B Points

#### 1: procedure BACKUP( $B, \Gamma_{t-1}$ ) **for** each action $a \in A$ **do** 2: **for** each observation $z \in Z$ **do** 3: **for** each solution vector $\alpha_i \in \Gamma_{t-1}$ **do** 4: $\alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s, a, s') O(s', a, z) \alpha(s'), \forall s \in S$ 5: end for 6: $\Gamma_t^{a,z} = \bigcup_i \alpha^{a,z}(s)$ 7: end for 8: end for 9: $\Gamma_t = \emptyset$ 10: for each belief point $b \in B$ do 11: 12. $\alpha_{b} = \operatorname{argmax}_{a} \left[ \sum_{s \in S} R(s, a) b(s) + \sum_{z \in Z} \operatorname{argmax}_{\alpha \in \Gamma^{a, z}} \sum_{s \in S} \alpha(s) b(s) \right]$ if $\alpha_b \notin \Gamma_t$ then 13: $\Gamma_t = \Gamma_t \cup \alpha_h$ 14. end if 15 end for 16· return $\Gamma_t$ 17. 18: end procedure Probabilistic Graphical Models MDP. POMDP 12 19. prosince 2024

24 22 / 1 - 112

- 1: procedure PBVI-MAIN(  $B_{Init}, \Gamma_0, N, T$  )
- 2:  $B = B_{Init}$
- 3:  $\Gamma = \Gamma_{Init}$
- 4: for *N* expansions do
- 5: for T iterations do
- 6:  $\Gamma = BACKUP(B, \Gamma)$
- 7: end for
- 8:  $B_{new} = EXPAND(B, \Gamma)$
- 9: end for
- 10: return Γ
- 11: end procedure

T either a horizon or we select a error bound  $\gamma^t ||V_0^* - V^*||$ .

# Expand: New Points Selection

1) at random



2) greedy maximal error improvement

- b' a new candidate
- $\circ$  the upper error bound in b'
  - $\int_{a}^{b} \frac{d^{2}}{ds} = a(s)(b'(s) b(s)) \qquad b'(s) \ge b(s) \\ \int_{a}^{b} \frac{d^{2}}{ds} = a(s)(b'(s) b(s)) \qquad b'(s) \le b(s)$
- $\sim b$  on the fringe, the error weighted by the probability of observations:

$$= \max_{x \in \mathcal{X}} \sum_{x \in \mathcal{X}} o(x, a, x) o(x(x, a, x)) = \max_{x \in \mathcal{X}} \sum_{x \in \mathcal{X}} o(x, a, x) o(x(x, a, x)) o(x(x,$$
1) at random



#### 2) greedy maximal error improvement

- b' a new candidate
- the upper error bound in *k*

$$f(b') \le \min_{b \in B} \sum_{s \in S} \begin{cases} (\frac{R_{max}}{1-\gamma} - \alpha(s))(b'(s) - b(s)) & b'(s) \ge b(s) \\ (\frac{R_{min}}{1-\gamma} - \alpha(s))(b'(s) - b(s)) & b'(s) < b(s) \end{cases}$$

• b on the fringe, the error weighted by the probability of observations:

$$\begin{aligned} \epsilon(b) &= \max_{a \in A} \sum_{z \in Z} O(b, a, z) \epsilon(\tau(b, a, z)) \\ &= \max_{a \in A} \sum_{z \in Z} \left[ \sum_{s' \in S} \sum_{s \in S} T(s, a, s') O(s', a, z) b(s) \right] \epsilon(\tau(b, a, z)). \end{aligned}$$

3

1) at random



- 2) greedy maximal error improvement
  - b' a new candidate
  - the upper error bound in *l*

$$\epsilon(b') \le \min_{b \in B} \sum_{s \in S} \begin{cases} \left(\frac{R_{max}}{1-\gamma} - \alpha(s)\right)(b'(s) - b(s)) & b'(s) \ge b(s) \\ \left(\frac{R_{min}}{1-\gamma} - \alpha(s)\right)(b'(s) - b(s)) & b'(s) < b(s) \end{cases}$$

• b on the fringe, the error weighted by the probability of observations:

$$\epsilon(b) = \max_{a \in A} \sum_{z \in Z} O(b, a, z) \epsilon(\tau(b, a, z))$$
$$= \max_{a \in A} \sum_{z \in Z} \left[ \sum_{s' \in S} \sum_{s \in S} T(s, a, s') O(s', a, z) b(s) \right] \epsilon(\tau(b, a, z)).$$

1) at random



- 2) greedy maximal error improvement
  - b' a new candidate
  - the upper error bound in b'

$$\epsilon(b') \leq \min_{b \in B} \sum_{s \in S} \begin{cases} (\frac{R_{max}}{1-\gamma} - \alpha(s))(b'(s) - b(s)) & b'(s) \geq b(s) \\ (\frac{R_{min}}{1-\gamma} - \alpha(s))(b'(s) - b(s)) & b'(s) < b(s) \end{cases}$$

• b on the fringe, the error weighted by the probability of observations:

$$\begin{aligned} \epsilon(b) &= \max_{a \in A} \sum_{z \in Z} O(b, a, z) \epsilon(\tau(b, a, z)) \\ &= \max_{a \in A} \sum_{z \in Z} \left[ \sum_{s' \in S} \sum_{s \in S} T(s, a, s') O(s', a, z) b(s) \right] \epsilon(\tau(b, a, z)). \end{aligned}$$

1) at random



- 2) greedy maximal error improvement
  - b' a new candidate
  - the upper error bound in b'

$$\epsilon(b') \le \min_{b \in B} \sum_{s \in S} \begin{cases} \left(\frac{R_{max}}{1-\gamma} - \alpha(s)\right) (b'(s) - b(s)) & b'(s) \ge b(s) \\ \left(\frac{R_{min}}{1-\gamma} - \alpha(s)\right) (b'(s) - b(s)) & b'(s) < b(s) \end{cases}$$

b on the fringe, the error weighted by the probability of observations:

$$\epsilon(b) = \max_{a \in A} \sum_{z \in Z} O(b, a, z) \epsilon(\tau(b, a, z))$$
$$= \max_{a \in A} \sum_{z \in Z} \left[ \sum_{s' \in S} \sum_{s \in S} T(s, a, s') O(s', a, z) b(s) \right] \epsilon(\tau(b, a, z)).$$

# Augmented MDP Application

#### a) and c) plans that ignore robot's perceptual uncertainty.

- b) and d) was computed by AMDP algorithm. It avoids regions where the robot is more likely to get lost.
- Useful for sensors with maximal range below 4m





# Augmented MDP Application

- a) and c) plans that ignore robot's perceptual uncertainty.
- b) and d) was computed by AMDP algorithm. It avoids regions where the robot is more likely to get lost.
- Useful for sensors with maximal range below 4m





# Augmented MDP Application

- a) and c) plans that ignore robot's perceptual uncertainty.
- b) and d) was computed by AMDP algorithm. It avoids regions where the robot is more likely to get lost.
- Useful for sensors with maximal range below 4m.





#### QMPD

• Each belief point adds a single linear 'Bellman' equation

$$\hat{V}(b) = \mathbb{E}_b[\hat{V}(s)] = \sum_{i=1}^N b(s_i)\hat{V}(s_i)$$

- Line 4: The value function is calculated in the original space X and controls u
- Line 6: The generalization to belief space.
- $\hat{V}$  can be used as an input to the POMDP evaluation.
- we may continue with T iterations, even small values of T help.

#### QMDP

1: procedure QMDP( b )  $\hat{V} = MDP\_discrete\_value\_iteration()$ 2: **for** each action  $a \in A$  do 3: **for** each state  $s \in S$  **do** 4.  $Q(s, a) = R(s, a) + \gamma \sum_{s' \in S} \hat{V}(s) p(s'|a, s)$ 5 end for 6. end for 7. **return** arg max<sub>a</sub>  $\sum_{s \in S} b(s)Q(s, a)$ 8. 9: end procedure

- Pengfei Zhu, Xin Li, Pascal Poupart, Guanghui Miao: On Improving Deep Reinforcement Learning for POMDPs, https://arxiv.org/abs/1704.07978v6
- Maximilian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, Shimon Whiteson: Deep Variational Reinforcement Learning for POMDPs, Proceedings of the 35 th International Conference on Machine

- Pengfei Zhu, Xin Li, Pascal Poupart, Guanghui Miao: On Improving Deep Reinforcement Learning for POMDPs, https://arxiv.org/abs/1704.07978v6
- Maximilian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, Shimon Whiteson: Deep Variational Reinforcement Learning for POMDPs, Proceedings of the 35 th International Conference on Machine

#### • Simple screens, logic

- 18 actions
- used to train an agent by Reinforcement learning with no prior information.
- MDP if fully observable



(e) MsPacman

• flickering for POMPD training (50% of figures are black)



(a) ChopperCommand

- Simple screens, logic
- 18 actions
- used to train an agent by Reinforcement learning with no prior information.
- MDP if fully observable



(e) MsPacman

• flickering for POMPD training (50% of figures are black)



(a) ChopperCommand

- Simple screens, logic
- 18 actions
- used to train an agent by Reinforcement learning with no prior information.
- MDP if fully observable



(e) MsPacman

• flickering for POMPD training (50% of figures are black)



(a) ChopperCommand

- Simple screens, logic
- 18 actions
- used to train an agent by Reinforcement learning with no prior information.
- MDP if fully observable



(e) MsPacman

• flickering for POMPD training (50% of figures are black)



(a) ChopperCommand

- Simple screens, logic
- 18 actions
- used to train an agent by Reinforcement learning with no prior information.
- MDP if fully observable



(e) MsPacman

• flickering for POMPD training (50% of figures are black)



standard reinforcement learning update

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

• NN  $\theta$  is trained to minimize the loss function

 $L(\theta_i) = \mathbb{E}_{s,a,r,s'}[(r + \gamma \max_{a'} Q(s', a'; \theta_i^{old}) - Q(s', a'; \theta_i))^2].$ 

DRQN uses recurrent LSTM network



• standard reinforcement learning update

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

• NN  $\boldsymbol{\theta}$  is trained to minimize the loss function

$$L(\theta_i) = \mathbb{E}_{s,a,r,s'}[(r + \gamma \max_{a'} Q(s',a';\theta_i^{old}) - Q(s',a';\theta_i))^2].$$

DRQN uses recurrent LSTM network



hidden state keeps the track in flickered images

• standard reinforcement learning update

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

• NN  $\boldsymbol{\theta}$  is trained to minimize the loss function

$$L(\theta_i) = \mathbb{E}_{s,a,r,s'}[(r + \gamma \max_{a'} Q(s',a';\theta_i^{old}) - Q(s',a';\theta_i))^2].$$

DRQN uses recurrent LSTM network



hidden state keeps the track in flickered images

• standard reinforcement learning update

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

• NN  $\boldsymbol{\theta}$  is trained to minimize the loss function

$$L(\theta_i) = \mathbb{E}_{s,a,r,s'}[(r + \gamma \max_{a'} Q(s',a';\theta_i^{old}) - Q(s',a';\theta_i))^2].$$

DRQN uses recurrent LSTM network



- hidden state keeps the track in flickered images
- ADRQN Action-specific Deep Recurent Q-network

• standard reinforcement learning update

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

• NN  $\boldsymbol{\theta}$  is trained to minimize the loss function

$$L(\theta_i) = \mathbb{E}_{s,a,r,s'}[(r + \gamma \max_{a'} Q(s',a';\theta_i^{old}) - Q(s',a';\theta_i))^2].$$

DRQN uses recurrent LSTM network



• hidden state keeps the track in flickered images

- The observation (image) is aggregated by three convolutional layers to 3136 dimensions
- the action is extended from 18 dimensional indicator vector to 512 dimensional vector
- these are concatenated and fed to LSTM
- result is reduced to the number of action dimension, an action is selected and passed to simulator
- the simulator provides the observation o and a new pass begins.



- The observation (image) is aggregated by three convolutional layers to 3136 dimensions
- the action is extended from 18 dimensional indicator vector to 512 dimensional vector
- these are concatenated and fed to LSTM
- result is reduced to the number of action dimension, an action is selected and passed to simulator
- the simulator provides the observation o and a new pass begins.



- The observation (image) is aggregated by three convolutional layers to 3136 dimensions
- the action is extended from 18 dimensional indicator vector to 512 dimensional vector
- these are concatenated and fed to LSTM
- result is reduced to the number of action dimension, an action is selected and passed to simulator
- the simulator provides the observation o and a new pass begins.



- The observation (image) is aggregated by three convolutional layers to 3136 dimensions
- the action is extended from 18 dimensional indicator vector to 512 dimensional vector
- these are concatenated and fed to LSTM
- result is reduced to the number of action dimension, an action is selected and passed to simulator
- the simulator provides the observation *o* and a new pass begins.



- The observation (image) is aggregated by three convolutional layers to 3136 dimensions
- the action is extended from 18 dimensional indicator vector to 512 dimensional vector
- these are concatenated and fed to LSTM
- result is reduced to the number of action dimension, an action is selected and passed to simulator
- the simulator provides the observation o and a new pass begins.



#### • Simulate *M* episodes

- actual reward r<sub>j</sub> is used
- future reward from the old network  $\theta^{old}$
- adjust weights  $\theta$  by the gradient of

$$[(r_{j} + \gamma max_{a}Q(h_{j}, a_{j}, o_{j+1}, a; \theta_{i}^{old}) - Q(h_{i-1}, a_{i-1}, o_{i}, a_{i}; \theta))^{2}]$$

#### Algorithm 1 Action-specific Deep Recurrent Q-Network

- 1: Initialize the replay memory D, # of iterations M
- 2: Initialize Q-Network and Target-Network with  $\theta$  and  $\theta^-$  respectively
- 3: for episode = 1 to M do

7.

- 4: Initialize the first action  $a_0 = no \ operation$ ,  $h_0 = \mathbf{0}$
- 5: Init.the first obs.o1 with the preprocessed first screen
- 6: while  $o_t \neq terminal$  do
  - Select a random action  $a_t$  with the probability  $\epsilon$
- 8: Else select  $a_t = \operatorname{argmax}_a Q(h_{t-1}, a_{t-1}, o_t, a; \theta)$
- 9: Execute action  $a_t$
- 10: Obtain reward  $r_t$  and resulting observation  $o_{t+1}$
- 11: Store transition  $\langle \{a_{t-1}, o_t\}, a_t, r_t, o_{t+1} \rangle$  as one record of the current episode in *D*
- 12: Randomly sample a minibatch of transition sequences  $\langle a_{j-1}, o_j, a_j, r_j, o_{j+1} \rangle$  from D
- 13: Compute Q-value of target network

$$y_{j} = \begin{cases} r_{j}, & o_{j+1} = terminal \\ r_{j} + \gamma \max_{a} Q(h_{j}, a_{j}, o_{j+1}, a; \theta^{-}), & o_{j+1} \neq terminal \end{cases}$$

14: Compute the gradient of 
$$(y_j - Q(h_{j-1}, a_{j-1}, o_j, a_j; \theta))^2$$
 to update  $\theta$   
15: end while  
16: end for

Mountain hike

#### POMPD

- $S = O = A = \mathbb{R}^2$
- state, observed state, the desire step.
- $\circ S_{i+1} = S_i + \tilde{S}_i + S_i$
- c ~ N(0, 0.25 ≤ l)
- $|\delta_t| \le 0.5$ , with the length capped  $|\delta_t| \le 0.5$ .
- $o_t = s_t + \epsilon_{o_t t}$
- $\bullet : \epsilon \sim N(0, m \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
- $R_t = r(x_t, y_t) = 0.01[3_t], r$  the figure.
- $h \sim h([-8.5, -8.5]^2, l),$
- T=75 steps.



• Red: Recurrent Neural Network

- Black: Deep Variational Reinforcement Learning
- dots: observations.

Mountain hike

#### POMPD

- $S = O = A = \mathbb{R}^2$
- state, observed state, the desired step.
- $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
- $\epsilon \sim N(0, 0.25 \cdot I)$
- $\tilde{a}_t$  is  $a_t$  with the length capped to  $\cdot 10$  $|\tilde{a}_t| \leq 0.5$ .
- $o_t = s_t + \epsilon_{o,t}$
- $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
- *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
- $b_0 \sim N([-8.5, -8.5]^T, I),$
- T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - $\epsilon \sim N(0, 0.25 \cdot I)$
  - $\tilde{a}_t$  is  $a_t$  with the length capped to  $\cdot 10$  $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|*ã<sub>t</sub>*|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I),$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_s$ ,
  - $\epsilon \sim N(0, 0.25 \cdot I)$
  - $\tilde{a}_t$  is  $a_t$  with the length capped to  $\cdot 10$  $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|*ã<sub>t</sub>*|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I),$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - $\epsilon \sim N(0, 0.25 \cdot I)$
  - $\tilde{a}_t$  is  $a_t$  with the length capped to 10  $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I),$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to 10  $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I),$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to 10 $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I),$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to |10| $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I), \sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to 10 $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim \textit{N}(0, \sigma_0 \cdot \textit{I}), \ \sigma_0 \in \{0, 1.5, 3\}$
  - $R_t = r(x_t, y_t) 0.01 |\tilde{a}_t|$ , r the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I),$
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to 10 $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I)$ ,  $\sigma_0 \in \{0, 1.5, 3\}$
  - R<sub>t</sub> = r(x<sub>t</sub>, y<sub>t</sub>) 0.01|ã<sub>t</sub>|, r the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$ ,
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.
- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to 10 $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I)$ ,  $\sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$ ,
  - T=75 steps



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to |10| $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I)$ ,  $\sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$ ,
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to |10| $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I)$ ,  $\sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$ ,
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to |10| $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I)$ ,  $\sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$ ,
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

- Mountain hike
- POMPD
  - $S = O = A = \mathbb{R}^2$
  - state, observed state, the desired step.
  - $s_{t+1} = s_t + \tilde{a}_t + \epsilon_{s,t}$
  - *ϵ* ∼ *N*(0, 0.25 · *I*)
  - $\tilde{a}_t$  is  $a_t$  with the length capped to |10| $|\tilde{a}_t| \leq 0.5$ .
  - $o_t = s_t + \epsilon_{o,t}$
  - $\epsilon \sim N(0, \sigma_0 \cdot I)$ ,  $\sigma_0 \in \{0, 1.5, 3\}$
  - *R<sub>t</sub>* = *r*(*x<sub>t</sub>*, *y<sub>t</sub>*) − 0.01|ã<sub>t</sub>|, *r* the color in the figure.
  - $b_0 \sim N([-8.5, -8.5]^T, I)$ ,
  - T=75 steps.



- Red: Recurrent Neural Network
- Black: Deep Variational Reinforcement Learning
- dots: observations.

### Recurrent Neural Network

### • *h<sub>t</sub>* latent state

- the current update *n<sub>s</sub>* steps
- computational graphs ng steps
- n<sub>g</sub> > n<sub>s</sub> greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

#### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs n<sub>g</sub> steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs ng steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs ng steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs ng steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs ng steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs ng steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

### Recurrent Neural Network

- *h<sub>t</sub>* latent state
- the current update n<sub>s</sub> steps
- computational graphs ng steps
- $n_g > n_s$  greatly improves the performance.

### DVRL

- K = 30 particles  $(h_t^k, z_t^k, w_t^k)$
- $h_t^k$  the latent state of an RNN
- $z_t^k$  an additional state to learn the transition
- $w_t^k$  an particle importance weight.



(a) **RNN-based approach.** The RNN acts as an encoder for the action-observation history, on which actor and critic are conditioned. The networks are updated end-to-end with an RL loss.



(b) DVRL. The agent learns a generative model which is used to update a belief distribution. Actor and critic now condition on the belief. The generative model is learned to optimise both the ELBO and the RL loss.

• resample particles

• RNN latent state update

$$u_{t-1}^k \sim \textit{Discrete}\left(rac{w_{t-1}^k}{\sum_{j=1}^K w_{t-1}^k}
ight)$$

$$h_t^k = \psi_{\theta}^{RNN}(h_{t-1}^{u_{t-1}^k}, z_t^k, a_{t-1}, o_t)$$

the weight: probability o<sub>t</sub> in our model

model the transition

$$z_t^k \sim q_\phi(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)$$

$$y_{t}^{k} = \frac{p_{\theta}(z_{t}^{k}|h_{t-1}^{u_{t-1}^{k}}, a_{t-1})p_{\theta}(o_{t}|h_{t-1}^{u_{t-1}^{k}}, z_{t}^{k}, a_{t-1})}{q_{\phi}(z_{t}^{k}|h_{t-1}^{u_{t-1}^{k}}, a_{t-1}, o_{t})}$$



э

• resample particles

• RNN latent state update

$$u_{t-1}^k \sim \textit{Discrete}\left(rac{w_{t-1}^k}{\sum_{j=1}^{K} w_{t-1}^k}
ight)$$

$$h_t^k = \psi_{\theta}^{RNN}(h_{t-1}^{u_{t-1}^{*}}, z_t^k, a_{t-1}, o_t)$$

• the weight: probability o<sub>t</sub> in our model

model the transition

$$z_t^k \sim q_{\phi}(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)$$

$$u_t^k = rac{p_{ heta}(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}) p_{ heta}(o_t | h_{t-1}^{u_{t-1}^k}, z_t^k, a_{t-1})}{q_{\phi}(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)}$$



- resample particles
  - $u_{t-1}^k \sim \textit{Discrete}\left(rac{w_{t-1}^k}{\sum_{i=1}^K w_{t-1}^k}
    ight)$
- RNN latent state update

$$h_{t}^{k} = \psi_{\theta}^{RNN}(h_{t-1}^{u_{t-1}^{k}}, z_{t}^{k}, a_{t-1}, o_{t})$$

• the weight: probability o<sub>t</sub> in our model

model the transition

$$z_t^k \sim q_\phi(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)$$

$$y_t^k = rac{p_{ heta}(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}) p_{ heta}(o_t | h_{t-1}^{u_{t-1}^k}, z_t^k, a_{t-1})}{q_{\phi}(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)}$$



resample particles

• RNN latent state update

$$u_{t-1}^k \sim \textit{Discrete}\left(rac{w_{t-1}^k}{\sum_{j=1}^K w_{t-1}^k}
ight)$$

$$h_{t}^{k} = \psi_{\theta}^{RNN}(h_{t-1}^{u_{t-1}^{k}}, z_{t}^{k}, a_{t-1}, o_{t})$$

• the weight: probability o<sub>t</sub> in our model

k

model the transition

$$z_t^k \sim q_\phi(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t) \qquad w_t^k = \frac{p_\theta(z_t^k | h_{t-1}^{u_{t-1}}, a_{t-1}) p_\theta(o_t | h_{t-1}^{u_{t-1}}, z_t^k, a_{t-1})}{q_\phi(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)}$$



.k

$$A_{\eta}^{t,i}(a_{t+1}, s_{t+1}) = \left(\sum_{j=0}^{n_s - i - 1} \gamma^j r_{r+i+j} + \gamma^{n_s - i} V_{\eta}^{old}(s_{t+n_s}) - V_{\eta}(s_{t+i})\right)$$

#### Loss function

• policy  $\rho$ :  $\mathcal{L}_{t}^{A}(\rho) = -\frac{1}{n_{s}n_{s}} \sum_{l=0}^{n_{s}} \log \pi_{\rho}(a_{t+1}|s_{t+1}) \mathcal{A}_{\eta}^{t,l,old}(a_{t+1},s_{t+1})$ • value  $\mathcal{L}_{t}^{V}(\eta) = \frac{1}{n_{s}n_{s}} \sum_{l=0}^{n_{s}} \mathcal{A}_{\eta}^{n_{s}-1} \mathcal{A}_{\eta}^{t,l}(a_{t+1},s_{t+1})^{2}$ • exploration  $\mathcal{L}_{t}^{P}(\eta) = \frac{1}{n_{s}n_{s}} \sum_{l=0}^{n_{s}} \mathcal{A}_{\eta}^{n_{s}-1} \mathcal{A}_{\eta}^{n_{s}}(a_{t+1},s_{t+1})^{2}$ • exploration  $\mathcal{L}_{t}^{P}(\eta) = \frac{1}{n_{s}n_{s}} \sum_{l=0}^{n_{s}} \mathcal{A}_{\eta}^{n_{s}-1} \mathcal{A}_{\eta}^{n_{s}}(a_{t+1},s_{t+1})^{2}$ •  $\rho_{s}\eta_{s}$  for  $\mathcal{L}_{t}^{P}(\eta) = \frac{1}{n_{s}n_{s}} \sum_{l=0}^{n_{s}} \mathcal{L}_{t}^{n_{s}}(a_{t},a_{t+1},a_{t+1})^{2}$ 

$$A_{\eta}^{t,i}(a_{t+1}, s_{t+1}) = \left(\sum_{j=0}^{n_s - i - 1} \gamma^j r_{r+i+j} + \gamma^{n_s - i} V_{\eta}^{old}(s_{t+n_s}) - V_{\eta}(s_{t+i})\right)$$

#### Loss function

• policy 
$$\rho$$
:  $\mathcal{L}_{t}^{\mathcal{A}}(\rho) = -\frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{e}} \sum_{i=0}^{n_{s}-1} \log \pi_{\rho}(a_{t+1}|s_{t+1}) \mathcal{A}_{\eta}^{t,i,old}(a_{t+1},s_{t+1})$ 

- value  $\mathcal{L}_t^V(\eta) = rac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} A_\eta^{t,i}(a_{t+1},s_{t+1})^2$
- exploration  $\mathcal{L}_t^H(\eta) = \frac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} H(\pi_{\rho}(.|s_{t+1}))$
- p, q, fit ELBO  $\mathcal{L}_t^{ELBO}(\theta, \phi) = -\frac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} \log \left(\frac{1}{K} \sum_{i=1}^K w_{t+i}^k\right)$

$$A_{\eta}^{t,i}(a_{t+1}, s_{t+1}) = \left(\sum_{j=0}^{n_s - i - 1} \gamma^j r_{r+i+j} + \gamma^{n_s - i} V_{\eta}^{old}(s_{t+n_s}) - V_{\eta}(s_{t+i})\right)$$

### Loss function

• policy 
$$\rho$$
:  $\mathcal{L}_{t}^{\mathcal{A}}(\rho) = -\frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{e}} \sum_{i=0}^{n_{s}-1} \log \pi_{\rho}(a_{t+1}|s_{t+1}) \mathcal{A}_{\eta}^{t,i,old}(a_{t+1},s_{t+1})$ 

• value 
$$\mathcal{L}_t^V(\eta) = rac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} \mathcal{A}_\eta^{t,i}(a_{t+1},s_{t+1})^2$$

• exploration 
$$\mathcal{L}_t^H(\eta) = \frac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} H(\pi_{\rho}(.|s_{t+1}))$$

• 
$$p, q$$
, fit - ELBO  $\mathcal{L}_t^{ELBO}(\theta, \phi) = -\frac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} \log \left( \frac{1}{K} \sum_{i=1}^K w_{t+i}^k \right)$ 

æ

$$A_{\eta}^{t,i}(a_{t+1}, s_{t+1}) = \left(\sum_{j=0}^{n_s - i - 1} \gamma^j r_{r+i+j} + \gamma^{n_s - i} V_{\eta}^{old}(s_{t+n_s}) - V_{\eta}(s_{t+i})\right)$$

### Loss function

• policy 
$$\rho$$
:  $\mathcal{L}_{t}^{A}(\rho) = -\frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{e}} \sum_{i=0}^{n_{s}-1} \log \pi_{\rho}(a_{t+1}|s_{t+1}) A_{\eta}^{t,i,old}(a_{t+1},s_{t+1})$ 

• value 
$$\mathcal{L}_t^V(\eta) = rac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} \mathcal{A}_\eta^{t,i}(a_{t+1},s_{t+1})^2$$

• exploration 
$$\mathcal{L}_t^H(\eta) = \frac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} H(\pi_{\rho}(.|s_{t+1}))$$

• 
$$p, q$$
, fit - ELBO  $\mathcal{L}_t^{ELBO}(\theta, \phi) = -\frac{1}{n_e n_s} \sum_{envs}^{n_e} \sum_{i=0}^{n_s-1} \log \left(\frac{1}{K} \sum_{i=1}^K w_{t+i}^k\right)$ 

æ

$$A_{\eta}^{t,i}(a_{t+1}, s_{t+1}) = \left(\sum_{j=0}^{n_s - i - 1} \gamma^j r_{r+i+j} + \gamma^{n_s - i} V_{\eta}^{old}(s_{t+n_s}) - V_{\eta}(s_{t+i})\right)$$

### Loss function

• policy 
$$\rho: \mathcal{L}_{t}^{A}(\rho) = -\frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{e}} \sum_{i=0}^{n_{s}-1} \log \pi_{\rho}(a_{t+1}|s_{t+1}) A_{\eta}^{t,i,old}(a_{t+1},s_{t+1})$$
  
• value  $\mathcal{L}_{t}^{V}(\eta) = \frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{s}-1} A_{\eta}^{t,i}(a_{t+1},s_{t+1})^{2}$   
• exploration  $\mathcal{L}_{t}^{H}(\eta) = \frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{e}} \sum_{i=0}^{n_{s}-1} H(\pi_{\rho}(.|s_{t+1}))$   
•  $p, q, \text{ fit} - \text{ELBO } \mathcal{L}_{t}^{ELBO}(\theta, \phi) = -\frac{1}{n_{e}n_{s}} \sum_{envs}^{n_{e}} \sum_{i=0}^{n_{s}-1} \log \left(\frac{1}{K} \sum_{i=1}^{K} w_{t+i}^{k}\right)$ 

- The average return is slightly better for DVRI





(a) Influence of the particle number on per- (b) Performance of the full DVRL algorithm formance for DVRL. Only using one particle compared to setting  $\lambda^E = 0$  ("No ELBO") tion in the latent state.

is not sufficient to encode enough informa- or not backpropagating the policy gradients through the encoder ("No joint optim").

Frames  $\times 10^7$ (c) Influence of the maximum backpropagation length  $n_{\sigma}$  on performance. Note that RNN suffers most from very short lengths. This is consistent with our conjecture that RNN relies mostly on memory, not inference.

--- RNN

4000





19. prosince 2024 36 / 113 - 245

- The average return is slightly better for DVRI
- One particle is not enough





(a) Influence of the particle number on per- (b) Performance of the full DVRL algorithm formance for DVRL. Only using one particle compared to setting  $\lambda^E = 0$  ("No ELBO") tion in the latent state.

is not sufficient to encode enough informa- or not backpropagating the policy gradients through the encoder ("No joint optim").

(c) Influence of the maximum backpropagation length  $n_{\sigma}$  on performance. Note that RNN suffers most from very short lengths. This is consistent with our conjecture that RNN relies mostly on memory, not inference.

Frames  $\times 10^7$ 

--- RNN

4000



19. prosince 2024 36 / 113 - 245

- The average return is slightly better for DVRI
- One particle is not enough
- Backpropagation is necessary, ELBO improves the result







(a) Influence of the particle number on per- (b) Performance of the full DVRL algorithm formance for DVRL. Only using one particle compared to setting  $\lambda^E = 0$  ("No ELBO") is not sufficient to encode enough informa- or not backpropagating the policy gradients tion in the latent state.

through the encoder ("No joint optim").

(c) Influence of the maximum backpropagation length  $n_{\sigma}$  on performance. Note that RNN suffers most from very short lengths. This is consistent with our conjecture that RNN relies mostly on memory, not inference.



- The average return is slightly better for DVRI
- One particle is not enough
- Backpropagation is necessary, ELBO improves the result
- memory length improves the result





(a) Influence of the particle number on per- (b) Performance of the full DVRL algorithm formance for DVRL. Only using one particle compared to setting  $\lambda^E = 0$  ("No ELBO") tion in the latent state.

is not sufficient to encode enough informa- or not backpropagating the policy gradients through the encoder ("No joint optim").

--- RNN 4000 Frames  $\times 10^7$ 

(c) Influence of the maximum backpropagation length  $n_{\sigma}$  on performance. Note that RNN suffers most from very short lengths. This is consistent with our conjecture that RNN relies mostly on memory, not inference.



- The average return is slightly better for DVRI
- One particle is not enough
- Backpropagation is necessary, ELBO improves the result
- memory length improves the result
- with minimal memory  $n_g = 5$  DVRL works, RNN does not.



tion in the latent state.



(a) Influence of the particle number on per- (b) Performance of the full DVRL algorithm formance for DVRL. Only using one particle compared to setting  $\lambda^{E} = 0$  ("No ELBO") is not sufficient to encode enough informa- or not backpropagating the policy gradients through the encoder ("No joint optim").

Frames (c) Influence of the maximum backpropagation length  $n_{\sigma}$  on performance. Note that RNN suffers most from very short lengths. This is consistent with our conjecture that RNN relies mostly on memory, not inference.

 $\times 10^{2}$ 

--- RNN

4000



- Bring you computer (if possible)
- Check moodle for libraries to install
  - Python 3
  - pgmpy, numpy, matplotlib
  - networkx
  - (graphviz, sklearn)
  - may be others.

# Have a nice holiday.

- Bring you computer (if possible)
- Check moodle for libraries to install
  - Python 3
  - pgmpy, numpy, matplotlib
  - networkx
  - (graphviz, sklearn)
  - may be others.

# Have a nice holiday.

- Karkus, Hsu, Lee: QMDP-Net: Deep Learning for Planning under Partial Observability https://proceedings.neurips.cc/paper/2017/file/e9412ee564384b987d086df32d4 Paper.pdf
- Eric Mueller and Mykel J. Kochenderfer :**Multi-Rotor Aircraft Collision Avoidance using Partially Observable Markov Decision Processes**, American Institute of Aeronautics and Astronautics

• Karkus, Hsu, Lee: QMDP-Net: Deep Learning for Planning under Partial Observability

https://proceedings.neurips.cc/paper/2017/file/e9412ee564384b987d086df32d4 Paper.pdf

• Eric Mueller and Mykel J. Kochenderfer :**Multi-Rotor Aircraft Collision Avoidance using Partially Observable Markov Decision Processes**, American Institute of Aeronautics and Astronautics

- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations.
- state 2D, (3D)
  - relative range states n<sub>1</sub>, n<sub>2</sub>, (n<sub>2</sub>)
  - velocities for the ownship  $x_{m}, x_{m'}$ ( $y_m$ )
  - velocities for the intruder v<sub>in</sub>, v<sub>y</sub>,
     (m)
  - absolute displacement from the desired trajectory  $d_{ij}$ ,  $d_{ji}$ ,  $(d_i)$
  - the desired trajectory is normalized to unit velocity in costs and zero velocity in the



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations.
- state 2D, (3D)
  - relative range states 6, 19, (6)
  - velocities for the ownship  $x_{m}, x_{m'}$ ( $y_m$ )
  - velocities for the intruder  $v_{ee}$ ,  $v_{ge}$ ( $v_{e}$ )
  - absolute displacement from the desired trajectory  $d_{ij}$ ,  $d_{ji}$ ,  $(d_i)$
  - the desired trajectory is normalized to unit velocity in provident of the velocity in the



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - $\sim$  relative range states  $r_{e_1} r_{e_2}$  ( $r_e$ )
  - velocities for the ownship  $x_{\alpha}, x_{\alpha'}$ ( $y_{\alpha}$ )
  - velocities for the intruder  $v_{ee}$ ,  $v_{ge}$ ( $v_{e}$ )
  - absolute displacement from the desired trajectory  $d_{\gamma}$ ,  $d_{\gamma}$ ,  $\{d_{\gamma}\}$
  - the desired trajectory is normalized to unit velocity in provident or velocity in the



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D
  - relative range states 6, 7, (6)
  - velocities for the ownship  $x_{av}, x_{av}$  ( $Y_{av}$ )
  - velocities for the intruder  $v_{60}$ ,  $v_{7}$ ,  $(v_{6})$
  - absolute displacement from the desired trajectory  $d_{\gamma\gamma}$   $d_{\gamma\gamma}$   $(d_{\gamma})$
  - the desired trajectory is normalized to unit velocity in zonis and zero velocity in the



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D
  - relative range states r<sub>x</sub>, r<sub>y</sub>, (r<sub>2</sub>)
     velocities for the ownship v<sub>ot</sub>, v<sub>ot</sub>
    - $(\mathbf{v}_{ac})$
  - $(v_{\mathcal{C}})$
  - absolute displacement from the desired trajectory  $d_{r_1}$   $d_{r_2}$
  - the desired trajectory is normalized to unit velocity in zonis and zero velocity in the


- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - relative range states  $r_x$ ,  $r_y$ ,  $(r_z)$
  - velocities for the ownship v<sub>ox</sub>, v<sub>oy</sub>, (v<sub>oz</sub>)
  - velocities for the intruder v<sub>ix</sub>, v<sub>iy</sub>, (v<sub>iz</sub>)
  - absolute displacement from the desired trajectory  $d_x$ ,  $d_y$ ,  $(d_z)$
  - the desired trajectory is normalized to unit velocity in the xaxis and zero velocity in the y axis.



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - relative range states  $r_x$ ,  $r_y$ ,  $(r_z)$
  - velocities for the ownship v<sub>ox</sub>, v<sub>oy</sub> (v<sub>oz</sub>)
  - velocities for the intruder v<sub>ix</sub>, v<sub>iy</sub>, (v<sub>iz</sub>)
  - absolute displacement from the desired trajectory  $d_x$ ,  $d_y$ ,  $(d_z)$
  - the desired trajectory is normalized to unit velocity in the xaxis and zero velocity in the y axis.



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - relative range states  $r_x$ ,  $r_y$ ,  $(r_z)$
  - velocities for the ownship  $v_{ox}$ ,  $v_{oy}$ ,  $(v_{oz})$
  - velocities for the intruder v<sub>ix</sub>, v<sub>iy</sub>, (v<sub>iz</sub>)
  - absolute displacement from the desired trajectory  $d_x$ ,  $d_y$ ,  $(d_z)$
  - the desired trajectory is normalized to unit velocity in the xaxis and zero velocity in the y axis.



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - relative range states  $r_x$ ,  $r_y$ ,  $(r_z)$
  - velocities for the ownship v<sub>ox</sub>, v<sub>oy</sub>, (v<sub>oz</sub>)
  - velocities for the intruder  $v_{ix}$ ,  $v_{iy}$ ,  $(v_{iz})$
  - absolute displacement from the desired trajectory  $d_x$ ,  $d_y$ ,  $(d_z)$
  - the desired trajectory is normalized to unit velocity in the xaxis and zero velocity in the y axis.



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - relative range states  $r_x$ ,  $r_y$ ,  $(r_z)$
  - velocities for the ownship v<sub>ox</sub>, v<sub>oy</sub>, (v<sub>oz</sub>)
  - velocities for the intruder  $v_{ix}$ ,  $v_{iy}$ ,  $(v_{iz})$
  - absolute displacement from the desired trajectory d<sub>x</sub>, d<sub>y</sub>, (d<sub>z</sub>)
  - the desired trajectory is normalized to unit velocity in the xaxis and zero velocity in the y axis.



- the algorithms designed for fixed-wing aircraft analyze
  - turns
  - vertical meneuvers
- multirotor aircraft (drones) and helicopters can also
  - horizontal plane accelerations
- state 2D, (3D)
  - relative range states  $r_x$ ,  $r_y$ ,  $(r_z)$
  - velocities for the ownship  $v_{ox}$ ,  $v_{oy}$ ,  $(v_{oz})$
  - velocities for the intruder  $v_{ix}$ ,  $v_{iy}$ ,  $(v_{iz})$
  - absolute displacement from the desired trajectory d<sub>x</sub>, d<sub>y</sub>, (d<sub>z</sub>)
  - the desired trajectory is normalized to unit velocity in the xaxis and zero velocity in the y axis.



### the prediction horizon is very short

- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- $N_*$  noise to the ownship, intruder, x and y axis
- Bellman update

transition from s with acceleration a to  $s^{\dagger}$ 

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s,a) \max_{a^{\mid}} Q[s^{\mid},a^{\mid}].$$

r<sub>x</sub>  $V_{ix} - V_{ox}$ r<sub>v</sub>  $V_{iy} - V_{oy}$  $a_x + N_{ox}$  $\dot{v}_{ox}$ =  $a_v + N_{ov}$ *v₀*<sub>v</sub> =. Vix Nix =Niv ċ<sub>iv</sub> = d<sub>×</sub> = $V_{tx} - V_{ox}$  $\dot{d}_{v}$  $V_{tv} - V_{ov}$ 

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- N<sub>\*</sub> noise to the ownship, intruder, x and y axis
- Bellman update

transition from s with acceleration a to  $s^{\dagger}$ 

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s,a) \max_{a^{\mid}} Q[s^{\mid},a^{\mid}].$$

$$\begin{array}{rcl} \dot{r}_{x} & = & v_{ix} - v_{ox} \\ \dot{r}_{y} & = & v_{iy} - v_{oy} \\ \dot{v}_{ox} & = & a_{x} + N_{ox} \\ \dot{v}_{oy} & = & a_{y} + N_{oy} \\ \dot{v}_{ix} & = & N_{ix} \\ \dot{v}_{iy} & = & N_{iy} \\ \dot{d}_{x} & = & v_{tx} - v_{ox} \\ \dot{d}_{y} & = & v_{ty} - v_{oy} \end{array}$$

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- *a<sub>x</sub>*, *a<sub>y</sub>* acceleration by the ownship
- N<sub>\*</sub> noise to the ownship, intruder, x and y axis
- Bellman update

transition from s with acceleration a to s

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s,a) \max_{a^{\mid}} Q[s^{\mid},a^{\mid}].$$

r<sub>x</sub>  $V_{ix} - V_{ox}$ r<sub>v</sub>  $= v_{iv} - v_{ov}$  $= a_x + N_{ox}$  $\dot{v}_{ox}$  $a_v + N_{ov}$ *v₀*<sub>v</sub> = . Vi× Nix =Niv ċ<sub>iv</sub> =d<sub>×</sub>  $= V_{tx} - V_{ox}$  $V_{tv} - V_{ov}$  $\dot{d}_{v}$ 

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- N<sub>\*</sub> noise to the ownship, intruder, x and y axis
- Bellman update transition from s with acceleration

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s|} T(s^{|}|s,a) max_{a|} Q[s^{|},a^{|}].$$

r<sub>x</sub>  $V_{ix} - V_{ox}$ r<sub>v</sub>  $= v_{iy} - v_{oy}$ *v₀*x  $= a_x + N_{ox}$  $= a_v + N_{ov}$ *v₀*<sub>v</sub> . Vi× Nix =  $\dot{v}_{iy}$ Niv =d<sub>×</sub>  $= v_{tx} - v_{ox}$  $V_{tv} - V_{ov}$  $\dot{d}_{v}$ 

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- N<sub>\*</sub> noise to the ownship, intruder, x and y axis
   N<sub>0</sub>(μ = 0, 0.30s<sup>-2</sup>), N<sub>1</sub>(μ = 0, 0.45s<sup>-2</sup>),
- Bellman update transition from s with acceleration a to s<sup>1</sup>

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s^{|}} T(s^{|}|s,a) \max_{a^{|}} Q[s^{|},a^{|}].$$

$$\dot{r}_{x} = v_{ix} - v_{ox}$$
  
 $\dot{r}_{y} = v_{iy} - v_{oy}$   
 $\dot{v}_{ox} = a_{x} + N_{ox}$   
 $\dot{v}_{oy} = a_{y} + N_{oy}$   
 $\dot{v}_{ix} = N_{ix}$   
 $\dot{v}_{iy} = N_{iy}$   
 $\dot{d}_{x} = v_{tx} - v_{ox}$   
 $\dot{d}_{y} = v_{ty} - v_{oy}$ 

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- $N_*$  noise to the ownship, intruder, x and y axis

•  $N_o(\mu=0,0.30s^{-2}),~N_i(\mu=0,0.45s^{-2})$ 

Bellman update

transition from *s* with acceleration *a* to *s* 

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s,a) \max_{a^{\mid}} Q[s^{\mid},a^{\mid}].$$

 $\begin{array}{rcl} \dot{r}_{x} & = & v_{ix} - v_{ox} \\ \dot{r}_{y} & = & v_{iy} - v_{oy} \\ \dot{v}_{ox} & = & a_{x} + N_{ox} \\ \dot{v}_{oy} & = & a_{y} + N_{oy} \\ \dot{v}_{ix} & = & N_{ix} \\ \dot{v}_{iy} & = & N_{iy} \\ \dot{d}_{x} & = & v_{tx} - v_{ox} \\ \dot{d}_{y} & = & v_{ty} - v_{oy} \end{array}$ 

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- $N_*$  noise to the ownship, intruder, x and y axis
  - $N_o(\mu = 0, 0.30s^{-2}), N_i(\mu = 0, 0.45s^{-2}),$

Bellman update

transition from s with acceleration a to s

$$Q[s,a] \leftarrow R(s,a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s,a) \max_{a^{\mid}} Q[s^{\mid},a^{\mid}].$$

$$\begin{array}{rcl} \dot{r}_{x} & = & v_{ix} - v_{ox} \\ \dot{r}_{y} & = & v_{iy} - v_{oy} \\ \dot{v}_{ox} & = & a_{x} + N_{ox} \\ \dot{v}_{oy} & = & a_{y} + N_{oy} \\ \dot{v}_{ix} & = & N_{ix} \\ \dot{v}_{iy} & = & N_{iy} \\ \dot{d}_{x} & = & v_{tx} - v_{ox} \\ \dot{d}_{y} & = & v_{ty} - v_{oy} \end{array}$$

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- $N_*$  noise to the ownship, intruder, x and y axis
  - $N_o(\mu = 0, 0.30s^{-2}), N_i(\mu = 0, 0.45s^{-2}),$
- Bellman update

transition from *s* with acceleration *a* to *s* 

$$Q[s, a] \leftarrow R(s, a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s, a) \max_{a^{\mid}} Q[s^{\mid}, a^{\mid}].$$

$$\begin{array}{rcl} \dot{r}_{x} & = & v_{ix} - v_{ox} \\ \dot{r}_{y} & = & v_{iy} - v_{oy} \\ \dot{v}_{ox} & = & a_{x} + N_{ox} \\ \dot{v}_{oy} & = & a_{y} + N_{oy} \\ \dot{v}_{ix} & = & N_{ix} \\ \dot{v}_{iy} & = & N_{iy} \\ \dot{d}_{x} & = & v_{tx} - v_{ox} \\ \dot{d}_{y} & = & v_{ty} - v_{oy} \end{array}$$

- the prediction horizon is very short
- updates done every 0.1 to 1 seconds
- simple update equation are sufficient
- not a benefit to using more complex dynamic equations.
- $a_x$ ,  $a_y$  acceleration by the ownship
- $N_*$  noise to the ownship, intruder, x and y axis

• 
$$N_o(\mu=0,0.30s^{-2})$$
,  $N_i(\mu=0,0.45s^{-2})$ ,

Bellman update

transition from s with acceleration a to  $s^{|}$ 

$$Q[s, a] \leftarrow R(s, a) + \gamma \sum_{s^{\mid}} T(s^{\mid}|s, a) max_{a^{\mid}} Q[s^{\mid}, a^{\mid}].$$

$$\begin{array}{rcl} \dot{r}_{x} & = & v_{ix} - v_{ox} \\ \dot{r}_{y} & = & v_{iy} - v_{oy} \\ \dot{v}_{ox} & = & a_{x} + N_{ox} \\ \dot{v}_{oy} & = & a_{y} + N_{oy} \\ \dot{v}_{ix} & = & N_{ix} \\ \dot{v}_{iy} & = & N_{iy} \\ \dot{d}_{x} & = & v_{tx} - v_{ox} \\ \dot{d}_{y} & = & v_{ty} - v_{oy} \end{array}$$

- collision
- physically impossible states
- keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

- collision
- physically impossible states
- keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

- collision
- physically impossible states
- keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

- collision
- physically impossible states
- keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

- collision
- physically impossible states
- keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

- Minimum reward R<sub>min</sub>
  - collision
  - physically impossible states
  - keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

- Minimum reward R<sub>min</sub>
  - collision
  - physically impossible states
  - keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory

•  $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

19. prosince 2024 41 / 113 - 245

э

- Minimum reward R<sub>min</sub>
  - collision
  - physically impossible states
  - keeps the sum finite
- we prefer no acceleration
- we prefer long distance to the intruder
- we prefer short distance to the desired trajectory
- $K_s, K_T, R_{min}$  weights was learned, k weights was = 1.

$$R(s,a) = max \left[ R_{min}, -(k_{ax}|a_x| + k_{ay}|a_y|) - K_s \frac{1}{k_{rx}r_x^2 + k_{ry}r_y^2} - K_T(k_{dx}d_x^2 + k_{dy}d_y^2) \right]$$

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

State variable	State Description	Discretization
$r_x, r_y$	Intruder range components	-15, [-7, -3], -1, 0, 1, [3, 7], 15
$v_{ox}, v_{oy}$	ownship velocity components	$-5, -3, -1, 0, 1, 3, 5 \ s^{-1}$
$v_{ix}, v_{iy}$	intruder velocity components	$-5, [-3], -1, 0, 1, [3], 5 \ s^{-1}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

### offline optimization

#### a few hours for coarse discretization, 1 PC

- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

State variable	State Description	Discretization
$r_x, r_y$ $v_{ox}, v_{oy}$ $v_{ix}, v_{iy}$	Intruder range components ownship velocity components intruder velocity components	$\begin{array}{c} -15, [-7, -3], -1, 0, 1, [3, 7], 15 \\ -5, -3, -1, 0, 1, 3, 5 \ s^{-1} \\ -5, [-3], -1, 0, 1, [3], 5 \ s^{-1} \end{array}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

State variable	State Description	Discretization
$r_x, r_y$ $v_{ox}, v_{oy}$ $v_{ix}, v_{iy}$	Intruder range components ownship velocity components intruder velocity components	$\begin{array}{c} -15, [-7, -3], -1, 0, 1, [3, 7], 15\\ -5, -3, -1, 0, 1, 3, 5 \ s^{-1}\\ -5, [-3], -1, 0, 1, [3], 5 \ s^{-1} \end{array}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

#### all values normalized

State variable	State Description	Discretization
$r_x, r_y$ $v_{ox}, v_{oy}$ $v_{ix}, v_{iy}$	Intruder range components ownship velocity components intruder velocity components	$\begin{array}{c} -15, [-7, -3], -1, 0, 1, [3, 7], 15\\ -5, -3, -1, 0, 1, 3, 5 \ s^{-1}\\ -5, [-3], -1, 0, 1, [3], 5 \ s^{-1} \end{array}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

э

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

#### all values normalized

the coarse set contained a total of 765,625 discrete states

State variable	State Description	Discretization
$r_x, r_y$	Intruder range components	-15, [-7, -3], -1, 0, 1, [3, 7], 15
$v_{ox}, v_{oy}$	ownship velocity components	$-5, -3, -1, 0, 1, 3, 5 \ s^{-1}$
$v_{ix}, v_{iy}$	intruder velocity components	$-5, [-3], -1, 0, 1, [3], 5 \ s^{-1}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

- the coarse set contained a total of 765,625 discrete states
- the finely discretized version contained 9,529,569 states.

State variable	State Description	Discretization
$r_x, r_y$	Intruder range components	-15, [-7, -3], -1, 0, 1, [3, 7], 15
$v_{ox}, v_{oy}$	ownship velocity components	$-5, -3, -1, 0, 1, 3, 5 \ s^{-1}$
$v_{ix}, v_{iy}$	intruder velocity components	$-5, [-3], -1, 0, 1, [3], 5 \ s^{-1}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.
- all values normalized
  - the coarse set contained a total of 765,625 discrete states
  - the finely discretized version contained 9,529,569 states.

State variable	State Description	Discretization
$r_x, r_y$	Intruder range components	-15, [-7, -3], -1, 0, 1, [3, 7], 15
$v_{ox}, v_{oy}$	ownship velocity components	$-5, -3, -1, 0, 1, 3, 5 \ s^{-1}$
$v_{ix}, v_{iy}$	intruder velocity components	$-5, [-3], -1, 0, 1, [3], 5 \ s^{-1}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

### offline optimization

- a few hours for coarse discretization, 1 PC
- initially stationary intruders
- intruders moving at uniform velocity with a variety of relative headings angles
- intruders state and dynamic uncertainty were added to the encounters.

#### all values normalized

- the coarse set contained a total of 765,625 discrete states
- the finely discretized version contained 9,529,569 states.

State variable	State Description	Discretization
$r_x, r_y$	Intruder range components	-15, [-7, -3], -1, 0, 1, [3, 7], 15
$v_{ox}, v_{oy}$	ownship velocity components	$-5, -3, -1, 0, 1, 3, 5 \ s^{-1}$
$v_{ix}, v_{iy}$	intruder velocity components	$-5, [-3], -1, 0, 1, [3], 5 \ s^{-1}$
$d_x, d_y$	desired trajectory distance	-10, [-3], -1, 0, 1, [3], 10

э

- The primary goal is to remain safely separated from the intruder aircraft.
  - r<sub>5%CPA</sub> 'the closest point of approach', we allow 5% trajectories a little bit closer.

- Figure: required 1.5 units, never closer than 1.1 units.
- Mean deviation distance from the desired trajectory  $\mu_{dev}$ .



Figure 3: Separation metric used to evaluate the collision avoidance algorithm



- The primary goal is to remain safely separated from the intruder aircraft.
  - r<sub>5%CPA</sub> 'the closest point of approach', we allow 5% trajectories a little bit closer.

- Figure: required 1.5 units, never closer than 1.1 units.
- Mean deviation distance from the desired trajectory  $\mu_{dev}$ .



Figure 3: Separation metric used to evaluate the collision avoidance algorithm



- The primary goal is to remain safely separated from the intruder aircraft.
  - r<sub>5%CPA</sub> 'the closest point of approach', we allow 5% trajectories a little bit closer.

- Figure: required 1.5 units, never closer than 1.1 units.
- Mean deviation distance from the desired trajectory  $\mu_{dev}$ .



Figure 3: Separation metric used to evaluate the collision avoidance algorithm



- The primary goal is to remain safely separated from the intruder aircraft.
  - r<sub>5%CPA</sub> 'the closest point of approach', we allow 5% trajectories a little bit closer.

- Figure: required 1.5 units, never closer than 1.1 units.
- Mean deviation distance from the desired trajectory  $\mu_{dev}$ .



Figure 3: Separation metric used to evaluate the collision avoidance algorithm



# Reward Tuning – Bayesian Optimization

- We tune  $R_P = (K_T, K_s, R_{min})$
- β weights the two objective functions

$$F(R_P) = (\beta \times (r_{5\% CPA})^{-1} + (1 - \beta) \times \mu_{dev})$$

- Gaussian process models  $F(R_P)$ .
- We determine the point at which the objective function is expected to have the largest improvement, *E*[*I*(*F*(*R<sub>P</sub>*))] over that of the current minimum.
- This set of *R<sub>P</sub>* is passed to QMDP to evaluate.
- until convergence.



Figure 5: Process for tuning POMDP reward parameters
- We tune  $R_P = (K_T, K_s, R_{min})$
- $\beta$  weights the two objective functions

$$F(R_P) = (\beta \times (r_{5\% CPA})^{-1} + (1-\beta) \times \mu_{dev})$$

- Gaussian process models  $F(R_P)$ .
- We determine the point at which the objective function is expected to have the largest improvement, E[I(F(R<sub>P</sub>))] over that of the current minimum.
- This set of *R<sub>P</sub>* is passed to QMDP to evaluate.
- until convergence.



Figure 5: Process for tuning POMDP reward parameters

- We tune  $R_P = (K_T, K_s, R_{min})$
- $\beta$  weights the two objective functions

$$F(R_P) = (\beta \times (r_{5\% CPA})^{-1} + (1-\beta) \times \mu_{dev})$$

#### • Gaussian process models $F(R_P)$ .

- We determine the point at which the objective function is expected to have the largest improvement, *E*[*I*(*F*(*R<sub>P</sub>*))] over that of the current minimum.
- This set of *R<sub>P</sub>* is passed to QMDP to evaluate.
- until convergence.



Figure 5: Process for tuning POMDP reward parameters

- We tune  $R_P = (K_T, K_s, R_{min})$
- $\beta$  weights the two objective functions

$$F(R_P) = (\beta \times (r_{5\% CPA})^{-1} + (1-\beta) \times \mu_{dev})$$

- Gaussian process models  $F(R_P)$ .
- We determine the point at which the objective function is expected to have the largest improvement, *E*[*I*(*F*(*R<sub>P</sub>*))] over that of the current minimum.
  - This set of *R<sub>P</sub>* is passed to QMDP to evaluate.
- until convergence.



Figure 5: Process for tuning POMDP reward parameters

- We tune  $R_P = (K_T, K_s, R_{min})$
- $\beta$  weights the two objective functions

$$F(R_P) = (\beta \times (r_{5\% CPA})^{-1} + (1-\beta) \times \mu_{dev})$$

- Gaussian process models  $F(R_P)$ .
- We determine the point at which the objective function is expected to have the largest improvement, *E*[*I*(*F*(*R<sub>P</sub>*))] over that of the current minimum.
- This set of *R<sub>P</sub>* is passed to QMDP to evaluate.
- until convergence.



Figure 5: Process for tuning POMDP reward parameters

- We tune  $R_P = (K_T, K_s, R_{min})$
- $\beta$  weights the two objective functions

$$F(R_P) = (\beta \times (r_{5\% CPA})^{-1} + (1 - \beta) \times \mu_{dev})$$

- Gaussian process models  $F(R_P)$ .
- We determine the point at which the objective function is expected to have the largest improvement, *E*[*I*(*F*(*R<sub>P</sub>*))] over that of the current minimum.
- This set of *R<sub>P</sub>* is passed to QMDP to evaluate.
- until convergence.



Figure 5: Process for tuning POMDP reward parameters

# **Bayesian Optimization**

- we know QMDP and F values for one or more x = R<sub>P</sub> points
- we search the point  $x = R_P^*$  to observe
- we minimize y = F(R<sub>P</sub><sup>\*</sup>) and search the maximal probability of improvement
- 'the chance to improve' is expressed by the **Expected improvement** (*EI*)



Peter I. Frazier: A Tutorial on Bayesian Optimization, rXiv:1807.02811v1 [stat.ML] 8 Jul 2018

#### • 194 parameter sets evaluated

- $\beta$  between 0.01 and 0.99
- resulting in nine non-dominated, Pareto-optimal designs.



#### Pareto Optimal frontier

- 194 parameter sets evaluated
- $\bullet~\beta$  between 0.01 and 0.99 .

• resulting in nine non-dominated, Pareto-optimal designs.



- 194 parameter sets evaluated
- $\bullet~\beta$  between 0.01 and 0.99 .
- resulting in nine non-dominated, Pareto-optimal designs.



#### • Left: intruder starts at (0,0),

- random heading, fixed velocity of the intruder
- the ownship starts at the blue cross

- Right: The goal is hovering
- the intruder comes from the right with the unknown behaviour.



- Left: intruder starts at (0,0),
- random heading, fixed velocity of the intruder
- the ownship starts at the blue cross

- Right: The goal is hovering
- the intruder comes from the right with the unknown behaviour.



- Left: intruder starts at (0,0),
- random heading, fixed velocity of the intruder
- the ownship starts at the blue cross

- Right: The goal is hovering
- the intruder comes from the right with the unknown behaviour.



- Left: intruder starts at (0,0),
- random heading, fixed velocity of the intruder
- the ownship starts at the blue cross

#### • Right: The goal is hovering

 the intruder comes from the right with the unknown behaviour.



- Left: intruder starts at (0,0),
- random heading, fixed velocity of the intruder
- the ownship starts at the blue cross

- Right: The goal is hovering
- the intruder comes from the right with the unknown behaviour.





The fine discretization improves the results.

Figure 10: Cumulative distributions of encounter model metrics as a function of state discretization

æ

- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller γ did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.
- The intruder is stationary.



- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller γ did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.



æ

Probabilistic Graphical Models Applications 13

- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller γ did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.



#### The intruder is stationary.

- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller  $\gamma$  did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.



The intruder is stationary.

- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller  $\gamma$  did not ensure the return to the desired path.

#### • Figures: Owhship at the origin

- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.





- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.

Probabilistic Graphical Models

• smaller  $\gamma$  did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, d = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.





- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller  $\gamma$  did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.



- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller  $\gamma$  did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, d = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.



The intruder is stationary.

- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller  $\gamma$  did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.



#### • The intruder is stationary.

<sup>19.</sup> prosince 2024 49 / 113 - 245

- considerable improvement in convergence speed by initializing by the value of previously evaluated policy
- the value of maximum negative reward influenced the convergence speed
- $\gamma \leftarrow 0.99$  taking hundred iterations to converge.
- smaller  $\gamma$  did not ensure the return to the desired path.

- Figures: Owhship at the origin
- different intruder positions
- policies indicated by color: black=up, red=right
- left: both own and intruder velocities are zero, *d* = 0
- right: owhship is moving in the positive y-axis direction at 1 s<sup>-1</sup> with zero trajectory error and nominal trajectory matches the velocity.
- Provide a state of the state of
- The intruder is stationary.

#### • Uncertainty does not increase with time, QMPD is justifiable.

- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = \max_{a} \left| \sum_{s \in S_p} Q(s, a) b(s) \right|$$

- The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions
  - rectangular interpolation between 2<sup>th</sup> nearest neighbors
  - simplex interpolation between n-1 nearest neighbor
  - prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = \max_{a} \left| \sum_{s \in S_{p}} Q(s, a) b(s) \right|$$

- The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions
  - rectangular interpolation between 2<sup>th</sup> nearest neighbors
  - simplex interpolation between n +1 nearest neighbor
  - prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = \max_{a} \left| \sum_{s \in S_{p}} Q(s, a) b(s) \right|$$

- The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions
  - rectangular interpolation between 2' nearest neighbor
  - simplex interpolation between n+1 nearest neighbor
  - prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = \max_{a} \left[ \sum_{s \in S_p} Q(s, a) b(s) \right]$$

The value Q(s<sup>(k)</sup>, a)b<sup>(k)</sup> approximated from QMDP solutions

- rectangular interpolation between 2" nearest neighbor.
- simplex interpolation between n+1 nearest neighbor
- prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = max_{a}\left[\sum_{s \in S_{p}} Q(s, a)b(s)\right]$$

- The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions
  - rectangular interpolation between 2<sup>n</sup> nearest neighbor
  - simplex interpolation between n + 1 nearest neighbor
  - prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = max_a \left[ \sum_{s \in S_p} Q(s, a) b(s) \right]$$

- The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions
  - rectangular interpolation between 2<sup>n</sup> nearest neighbor
  - simplex interpolation between n + 1 nearest neighbor
  - prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = max_a \left[ \sum_{s \in S_p} Q(s, a) b(s) \right]$$

- The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions
  - rectangular interpolation between 2<sup>n</sup> nearest neighbor
  - simplex interpolation between n+1 nearest neighbor
  - prior work has found little benefit to using more sophisticated approaches.

- Uncertainty does not increase with time, QMPD is justifiable.
- State uncertainty is incorporated only when actions are selected
- a set of potential states is calculated from the observations received at each step.
- the probability distribution on potential states is the belief used to select an action.

$$\pi(b) = max_{a}\left[\sum_{s \in \mathcal{S}_{p}} Q(s, a)b(s)
ight]$$

• The value  $Q(s^{(k)}, a)b^{(k)}$  approximated from QMDP solutions

- rectangular interpolation between 2<sup>n</sup> nearest neighbor
- simplex interpolation between n+1 nearest neighbor
- prior work has found little benefit to using more sophisticated approaches.

# Table of Content

- Basics, Classifiers
- Variable Elimination Algorithm, Hidden Markov Models
- Markov Random Fields and Other Models
- Junction Tree Algorithm (Optimized Evaluation)
- 5 Approximate Evaluation
- 6 Structure Learning
- 🕜 Bayesian Learning, EM Algorithm
- 8 Gaussian Graphical Models
- Gaussian Processes
- 10 Variational Approximation
- Decision Trees, Decision Graphs
- 12 MDP, POMDP
  - 3 Applications

# Summary Links



<sup>19.</sup> prosince 2024 50 / 246 - 247