

Artificial Neural Networks (Gerstner). Solutions for week 13

Intrinsically motivated exploration

Exercise 1. Information-gain, surprise, and the number of observations

Consider an environment with a finite set of states \mathcal{S} and a finite set of actions \mathcal{A} . At each time $t > 0$, we assume that the agent uses its past experiences (i.e., $s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t$) and estimates the environment transition probabilities as

$$\hat{p}^{(t)}(s'|s, a) = \frac{T_{s,a,s'}^{(t)} + \epsilon}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon},$$

where $T_{s,a}^{(t)}$ is the number of times that the agent has taken action a in state s until time t , $T_{s,a,s'}^{(t)}$ is the number of times that taking action a in state s took the agent to state s' , $|\mathcal{S}|$ is the total number of states, and $\epsilon > 0$ is a small constant to avoid division by zero.

Consider $s_t = s$ and $a_t = a$. In this exercise, we study the Information Gain (IG) of the transition $(s, a) \rightarrow s_{t+1}$ and its link to surprise and number of observations.

- Given the transition $(s, a) \rightarrow s_{t+1}$ at $t + 1$, what is the updated $\hat{p}^{(t+1)}(s'|s, a)$ for all $s' \in \mathcal{S}$? Write your answer as a function of $T_{s,a,s'}^{(t)}$, $T_{s,a}^{(t)}$, $|\mathcal{S}|$ and ϵ .
- Show that

$$\hat{p}^{(t+1)}(s'|s, a) - \hat{p}^{(t)}(s'|s, a) = \frac{1}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} \left(\delta_{s',s_{t+1}} - \hat{p}^{(t)}(s'|s, a) \right),$$

where δ is the Kronecker delta function.

- One approach to define information gain is the L1 norm of the difference between $\hat{p}^{(t+1)}(\cdot|s, a)$ and $\hat{p}^{(t)}(\cdot|s, a)$:

$$\text{IG}_{t+1} = \sum_{s' \in \mathcal{S}} \left| \hat{p}^{(t+1)}(s'|s, a) - \hat{p}^{(t)}(s'|s, a) \right|.$$

Find IG_{t+1} as a function of $T_{s,a}^{(t)}$, $|\mathcal{S}|$, ϵ , and $\hat{p}^{(t)}(s_{t+1}|s, a)$.

How does increasing the number of observation $T_{s,a}^{(t)}$ influence the information gain IG_{t+1} ?

- One of the many ways to define the surprise of the transition $(s, a) \rightarrow s_{t+1}$ is to use the notion of ‘State Prediction Error’ (see [Modirshanechi et al. 2022](#) for its link to other definitions of surprise):

$$\text{SPE}_{t+1} = 1 - \hat{p}^{(t)}(s_{t+1}|s, a).$$

Rewrite IG_{t+1} as a function of $T_{s,a}^{(t)}$, $|\mathcal{S}|$, ϵ , and SPE_{t+1} .

How does the information gain IG_{t+1} relate to the state prediction error SPE_{t+1} ?

- Assume that we know the true transition probabilities $p(\cdot|s, a)$ and that $\lim_{t \rightarrow \infty} T_{s,a}^{(t)} = \infty$ (i.e., agents choose each action infinitely many times). For a given next state $s_{t+1} = s'$, find the limits

$$\lim_{t \rightarrow \infty} \text{SPE}_{t+1} \quad \text{and} \quad \lim_{t \rightarrow \infty} \text{IG}_{t+1}.$$

What do these results imply about seeking SPE or IG as intrinsic rewards in the presence of stochasticity?

Which intrinsic reward is less prone to the noisy-TV problem?

Solution:

- Using the definition of the estimated transition probabilities, we have

$$\hat{p}^{(t+1)}(s'|s, a) = \frac{T_{s,a,s'}^{(t+1)} + \epsilon}{T_{s,a}^{(t+1)} + |\mathcal{S}|\epsilon} = \frac{T_{s,a,s'}^{(t)} + \epsilon + \delta_{s',s_{t+1}}}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1}.$$

b. Using part a, we can write

$$\begin{aligned}\hat{p}^{(t+1)}(s'|s, a) - \hat{p}^{(t)}(s'|s, a) &= \frac{T_{s,a,s'}^{(t)} + \epsilon + \delta_{s',s_{t+1}}}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} - \frac{T_{s,a,s'}^{(t)} + \epsilon}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon} \\ &= \frac{1}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} \left(T_{s,a,s'}^{(t)} + \epsilon + \delta_{s',s_{t+1}} - \left(T_{s,a,s'}^{(t)} + \epsilon \right) \underbrace{\frac{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon}}_{=1 + \frac{1}{\frac{T_{s,a}^{(t)}}{|\mathcal{S}|\epsilon + 1}}} \right).\end{aligned}$$

By simplifying the term in the parantheses and reusing the definition of $\hat{p}^{(t)}(s'|s, a)$, we have

$$\hat{p}^{(t+1)}(s'|s, a) - \hat{p}^{(t)}(s'|s, a) = \frac{1}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} \left(\delta_{s',s_{t+1}} - \hat{p}^{(t)}(s'|s, a) \right).$$

c. We can use the result of part b and write

$$\begin{aligned}\text{IG}_{t+1} &= \frac{1}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} \sum_{s' \in \mathcal{S}} \left| \delta_{s',s_{t+1}} - \hat{p}^{(t)}(s'|s, a) \right| \\ &= \frac{1}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} \left(1 - \hat{p}^{(t)}(s_{t+1}|s, a) + \underbrace{\sum_{s' \in \mathcal{S} - \{s_{t+1}\}} \hat{p}^{(t)}(s'|s, a)}_{=1 - \hat{p}^{(t)}(s_{t+1}|s, a)} \right) = \frac{2(1 - \hat{p}^{(t)}(s_{t+1}|s, a))}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1}.\end{aligned}$$

IG_{t+1} is a strictly decreasing function $T_{s,a}^{(t)}$.

d. We can use the result of part c and write

$$\text{IG}_{t+1} = \frac{2 \text{SPE}_{t+1}}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1}.$$

This means that the information-gain as defined in this exercise is equal to a scaled version of the state prediction error, while the scaling factor decreases over time. Hence, the same SPE will lead to less changes in the estimate of transition probabilities with accumulation of observations.

e. We have

$$\lim_{t \rightarrow \infty} \text{SPE}_{t+1} = 1 - \lim_{t \rightarrow \infty} \frac{T_{s,a,s'}^{(t)} + \epsilon}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon} = 1 - p(s'|s, a)$$

and

$$\lim_{t \rightarrow \infty} \text{IG}_{t+1} = \lim_{t \rightarrow \infty} \frac{2 \text{SPE}_{t+1}}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} = \lim_{t \rightarrow \infty} \frac{2 \lim_{t \rightarrow \infty} \text{SPE}_{t+1}}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} = \lim_{t \rightarrow \infty} \frac{2(1 - p(s'|s, a))}{T_{s,a}^{(t)} + |\mathcal{S}|\epsilon + 1} = 0.$$

This implies that the limit of IG is always equal to zero independently of the level of the stochasticity in the environment. However, SPE can have high values even in the limit of $t \rightarrow \infty$ if the environment is stochastic (if $0 < p(s'|s, a) < 1$). Hence, SPE is prone to the Noisy-TV problem but IG is not.

Exercise 2. Disagreement and information-gain

Consider an environment with a finite set of states \mathcal{S} and a finite set of actions \mathcal{A} . At each time $t > 0$, we assume that the agent uses its past experiences (i.e., $s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t$) and estimates the environment transition probabilities with K different and parallel models, i.e., for $k \in \{1, \dots, K\}$, we have

$$\hat{p}_k^{(t)}(s'|s, a) = \frac{T_{s,a,s'}^{(t)} + \hat{p}_k^{(0)}(s'|s, a)}{T_{s,a}^{(t)} + 1},$$

where $T_{s,a}^{(t)}$ is the number of times that the agent has taken action a in state s until time t , $T_{s,a,s'}^{(t)}$ is the number of times that taking action a in state s took the agent to state s' , and $\hat{p}_k^{(0)}(s'|s, a)$ is the random initialization of model $k \in \{1, \dots, K\}$. The agent uses

$$\hat{p}^{(t)}(s'|s, a) = \frac{1}{K} \sum_{k=1}^K \hat{p}_k^{(t)}(s'|s, a) = \frac{T_{s,a,s'}^{(t)} + \hat{p}^{(0)}(s'|s, a)}{T_{s,a}^{(t)} + 1},$$

as its final estimate.

Consider $s_t = s$ and $a_t = a$. In this exercise, we study how the disagreement of the different K models relate to the information-gain of the transition $(s, a) \rightarrow s_{t+1}$.

- a. Repeat what you did in [Exercise 1](#) to calculate the information gain defined as

$$\text{IG}_{t+1} = \sum_{s' \in \mathcal{S}} \left| \hat{p}^{(t+1)}(s'|s, a) - \hat{p}^{(t)}(s'|s, a) \right|$$

as a function of $T_{s,a}^{(t)}$ and $\text{SPE}_{t+1} = 1 - \hat{p}^{(t)}(s_{t+1}|s, a)$.

- b. We define the disagreement at time t as

$$D_t = \frac{1}{K} \sum_{k=1}^K \sum_{s' \in \mathcal{S}} \left(\hat{p}^{(t)}(s'|s, a) - \hat{p}_k^{(t)}(s'|s, a) \right)^2.$$

Find D_t as a function of $T_{s,a}^{(t)}$ and the initial disagreement D_0 at $t = 0$.

- c. We now compare three different intrinsic rewards: the State Prediction Error SPE, the Information Gain IG, and the Disagreement D. Let us assume that we know the true transition probabilities $p(\cdot|s, a)$ and that $\lim_{t \rightarrow \infty} T_{s,a}^{(t)} = \infty$ (i.e., agents choose each action infinitely many times). For a given next state $s_{t+1} = s'$, compare the limits

$$\lim_{t \rightarrow \infty} \text{SPE}_{t+1}, \quad \lim_{t \rightarrow \infty} \text{IG}_{t+1}, \quad \text{and} \quad \lim_{t \rightarrow \infty} D_t.$$

What do these results imply about seeking these different intrinsic rewards in the presence of stochasticity? Which intrinsic reward is less prone to the noisy-TV problem?

Solution:

- a. We can repeat the exact same procedure as in [Exercise 1](#) and find

$$\text{IG}_{t+1} = \frac{2 \text{SPE}_{t+1}}{T_{s,a}^{(t)} + 2}.$$

- b. Using the definition of D_t , we have

$$\begin{aligned} D_t &= \frac{1}{K} \sum_{k=1}^K \sum_{s' \in \mathcal{S}} \left(\frac{T_{s,a,s'}^{(t)} + \hat{p}^{(0)}(s'|s, a)}{T_{s,a}^{(t)} + 1} - \frac{T_{s,a,s'}^{(t)} + \hat{p}_k^{(0)}(s'|s, a)}{T_{s,a}^{(t)} + 1} \right)^2 \\ &= \frac{1}{(T_{s,a}^{(t)} + 1)^2} \frac{1}{K} \sum_{k=1}^K \sum_{s' \in \mathcal{S}} \left(\hat{p}^{(0)}(s'|s, a) - \hat{p}_k^{(0)}(s'|s, a) \right)^2 = \frac{D_0}{(T_{s,a}^{(t)} + 1)^2}. \end{aligned}$$

- c. We have (similar to [Exercise 1](#))

$$\lim_{t \rightarrow \infty} \text{SPE}_{t+1} = 1 - p(s'|s, a), \quad \lim_{t \rightarrow \infty} \text{IG}_{t+1} = 0, \quad \text{and} \quad \lim_{t \rightarrow \infty} D_t = 0.$$

This implies that the limits of disagreement and IG are always equal to zero independently of the level of the stochasticity in the environment. However, SPE can have high values even in the limit of $t \rightarrow \infty$ if the environment is stochastic (if $0 < p(s'|s, a) < 1$). Hence, disagreement and IG are robust with respect to the noisy-TV problem.