

Reinforcement Learning China Summer School



RLChina 2021

Bayesian Brain

Prof. Jun Wang

UCL

August 17, 2021

Life and the universe



Life and the universe

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

What is “learning”?

- in biology, **learning** means
a change of behaviour as a result of experience
 - in classical conditioning¹, animals can learn to identify a useful pattern in the environment by associating one stimulus with another:
repeated given {*ring-a-bell*, *food*}, a dog will start to salivate (anticipate the upcoming of the food) when bell ringing again
- learned behaviours are **adaptive**, and thus are essential for animals to survive in the changing environment
 - e.g., they may learn not to eat certain foods if they have ever become ill after eating them
- more learned behaviours \implies more intelligent

¹Ivan Petrovitch Pavlov and William Gantt. “Lectures on conditioned reflexes: Twenty-five years of objective study of the higher nervous activity (behaviour) of animals.”. In: (1928).

Learning is not limited to biological systems

- *machine learning* is to answer the question of how **computer programmes** can improve their performance through **experience**
 - past experience exists in various forms, typically including
 - (1) human-labelled or unlabelled data sets,
 - (2) past interactions with the environments, or
 - (3) data collected from realistic simulators
- with experience, computers can be trained to identify the associations, spot the underlying patterns, and make accurate predictions and forecasts the futures
- as a field, machine learning studies fundamental theory and algorithms and provides a principled solution for the computational methods of learning

Three levels in any information processing system²

- **computational theory**

- what is the problem in a generic manner?
- what is the computing goal?
- what is the logic of strategy behind?

- **representation and algorithm**

- how can the identified computational problems be solved?
- what are the inputs/outputs and the algorithm for the transformation?

- **hardware implementation**

- how can the representation and algorithm be realised physically?
 - human: biological
 - AI: silicon using transistors

²David Marr. "Vision: A computational investigation into the human representation and processing of visual information". In: *Inc., New York, NY* 2.4.2 (1982).

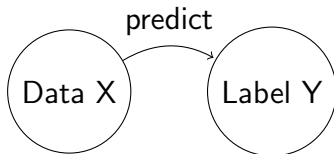
Human intelligent v.s. current AI³

*Dseitpe the fcat taht the letetres in tehese wrdos are jmbuled,
you are sitll albe to raed tehm, Bceasue the frsit and lsat ltertes
are in the rghit palce. your bairn can use tohse ceus to fgiure out
waht I'm syanig*

研究表明,汉字序顺并不定一影阅响读, 事证实明了当你看这完句话之后才发字现都乱是的

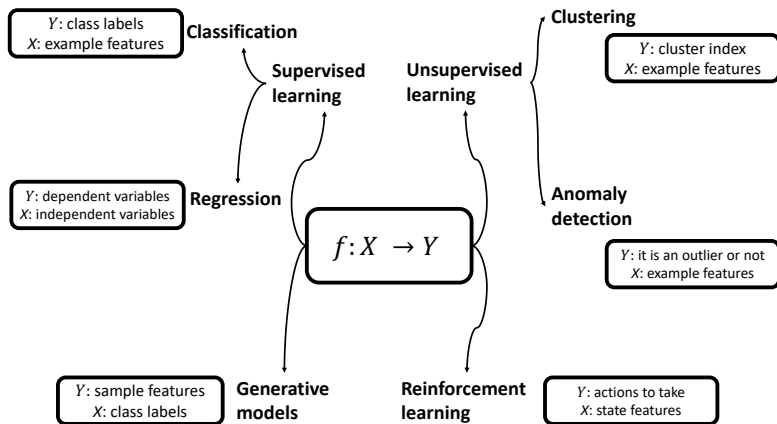
³Eliezer Sternberg. *NeuroLogic: The Brain's Hidden Rationale Behind Our Irrational Behavior*. Vintage, 2016.

Typical narrow machine learning



- mathematically, ml can loosely boil down to the question of finding a unknown function mapping $f : X \rightarrow Y$,
 - X is input feature space representing data points,
 - Y is label space representing the knowledge outputs
 - thus the mapping f represents a *knowledge discovery* process from a given data point $x \in X$ to the specific label $y \in Y$ associated with the data point: $y = f(x)$
- however, as f is not known a priori, the goal of the learning is to identify a hypothesis $h \in H$ from a predefined set H to approximate the unknown f , so that
 - the learned function h can *predict* the output variable $\hat{y}_{\text{new}} = h(x_{\text{new}})$ for a given new data point

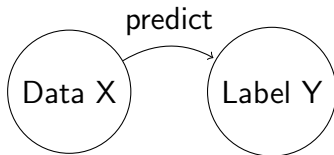
View **narrow** machine learning as function mapping



machine learning tasks are different in the experience available, the objective function, and the specific learning algorithms

Machine learning: reasoning under uncertainty

pattern recognition:



decision making (reinforcement learning):

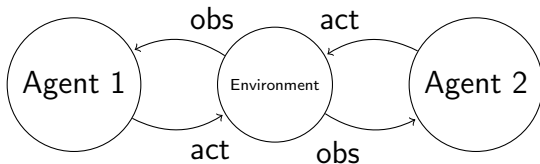
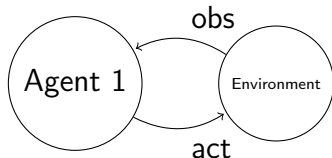


Table of Contents

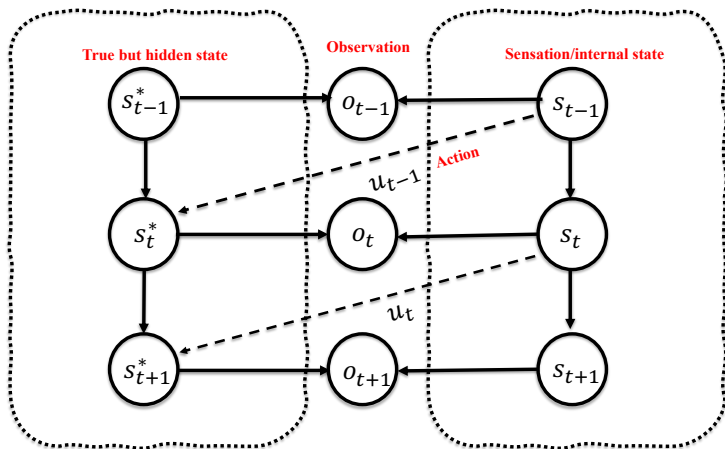
- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

A living system and its environment



Generative
process

Generative
model



Entropy

- *Entropy*: a measure of the average **surprise** (or uncertainty or disorder) of a random event sampled from a probability distribution or density
- definition: entropy of a random event X with possible outcomes x_1, \dots, x_n :

$$H(p) = - \sum_i p(X = x_i) \log p(X = x_i)$$

- low entropy means that, on average, the outcome is relatively predictable

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

Variational Bayes⁴/Free energy principle⁵

- an agent builds a **world model** by having an internal representation of the world s from observation o as:

$$p(s|o) = \frac{p(o|s)p(s)}{\int_s p(o|s)p(s) dx}$$

- typically the denominator is intractable and one can approximate the posterior using a tractable function family $q(s) \in \mathbb{Q}$ by minimising a dis-similarity measure, e.g., KL Divergence:

$$KL(q(s)||p(s|o)) = \mathbb{E}_{q(s)}[\log \frac{q(s)}{p(s|o)}]$$

⁴Michael I Jordan et al. “An introduction to variational methods for graphical models”. In: *Machine learning* 37.2 (1999), pp. 183–233.

⁵Karl Friston. “The free-energy principle: a unified brain theory?” In: *Nature reviews neuroscience* 11.2 (2010), pp. 127–138.

Variational Bayes⁶/Variational free energy⁷

- we then derive:

$$\begin{aligned}KL(q(s)||p(s|o)) &= \mathbb{E}_{q(s)}[\log \frac{q(s)}{p(s|o)}] \\&= \mathbb{E}_{q(s)}[\log q(s) - \log p(s, o) + \log p(o)] \\&= \mathbb{E}_{q(s)}[\log q(s) - \log p(s, o)] + \log p(o)\end{aligned}$$

- reorganising gives the measure of **surprise** (or negative model evidence or log marginal distribution):

$$\begin{aligned}-\log p(o) &= \mathbb{E}_{q(s)}[\log q(s) - \log p(s, o)] - D_{KL}[q(s)||p(s|o)] \\&\leq \mathbb{E}_{q(s)}[\log q(s) - \log p(s, o)] \equiv \mathbb{F}(q)\end{aligned}$$

where RHS is **variational free energy** (negative ELBO, Evidence Lower BOund)

⁶Michael I Jordan et al. "An introduction to variational methods for graphical models". In: *Machine learning* 37.2 (1999), pp. 183–233.

⁷Karl Friston. "The free-energy principle: a unified brain theory?" In: *Nature reviews neuroscience* 11.2 (2010), pp. 127–138.

Free energy principle⁹

- the VFE can be decomposed in three principal ways:

$$\begin{aligned}\mathbb{F}(q) &= \mathbb{E}_{q(s)}[\log q(s) - \log p(s, o)] \\ &= \underbrace{\mathbb{E}_{q(s)}[\log q(s)]}_{\text{Negative Entropy}} - \underbrace{\mathbb{E}_{q(s)}[\log p(s, o)]}_{\text{Energy}} \\ &= \underbrace{\mathbb{E}_{q(s)}[\log p(o|s)]}_{\text{Accuracy}} + \underbrace{D_{KL}[q(s)|p(s)]}_{\text{Complexity}} \\ &= \underbrace{D_{KL}[q(s)||p(s|o)]}_{\text{Posterior Divergence}} - \underbrace{\log p(o)}_{\text{Negative Log Model Evidence}}\end{aligned}$$

- free energy principle**: the goal of a living system is to minimise the free energy \mathbb{F} in order to avoid surprising observations (states) \rightarrow maintain **homeostasis**⁸ (thus remain alive)

⁸The process whereby an open or closed system regulates its internal environment to maintain its states within bounds.

⁹Karl Friston. "The free-energy principle: a unified brain theory?" In: *Nature reviews neuroscience* 11.2 (2010), pp. 127–138.

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop**
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

Extension to handle dynamics and control

- the real world is **dynamic** and thus we extend the model to have an observation $o = \{o_1, \dots, o_T\}$ and a hidden state at each point in time $s = \{s_1, \dots, s_T\}$
- also add **action policy** $\pi = \{u_1, \dots, u_T\}$ which interacts with the environment by altering the next hidden state
- the generative model: 1) $p(s_{t+1} | s_t, u_t), t > 1; p(s_1)$ 2) $p(o_t | s_t)$
- the **Expected Free Energy** (EFE), from time τ until the time horizon T :
$$\mathcal{G} = \mathbb{E}_{q(o_{\tau:T}, s_{\tau:T}, \pi)} [\ln q(s_{\tau:T}, \pi) - \ln \tilde{p}(o_{\tau:T}, s_{\tau:T})]$$
- where an agent's goals are encoded as (subjective) **desired distribution** over observations $\tilde{p}(o_{\tau:T})^{10}$; thus we have $\tilde{p}(o_\tau, s_\tau) \approx \tilde{p}(o_\tau) q(s_\tau | o_\tau)$

¹⁰in ML, an optimality variable is added to encode the desires

Expected free energy¹²

- a temporal **mean-field factorisation**¹¹ is assumed:
 - the variational function: $q(s_{T:T}, \pi) \approx q(\pi) \prod_{\tau} q(s_{\tau} | \pi)$
 - the gen. model: $\tilde{p}(o_{T:T}, s_{T:T}) \approx \prod_t \tilde{p}(o_{\tau}) q(s_{\tau} | o_{\tau})$
- as a result, they are independent between time steps:

$$\begin{aligned}\mathbb{G} &= \mathbb{E}_{q(o_{T:T}, s_{T:T}, \pi)} [\ln q(s_{T:T}, \pi) - \ln \tilde{p}(o_{T:T}, s_{T:T})] \\&= \mathbb{E}_{q(o_{T:T}, s_{T:T} | \pi) q(\pi)} [\ln q(s_{T:T} | \pi) + \ln q(\pi) - \ln \tilde{p}(o_{T:T}, s_{T:T})] \\&= \mathbb{E}_{q(\pi)} \left[\ln q(\pi) + \mathbb{E}_{q(o_{T:T}, s_{T:T} | \pi)} \left[\sum_{\tau} [\ln q(s_{\tau} | \pi) - \ln \tilde{p}(o_{\tau}, s_{\tau})] \right] \right] \\&= \mathbb{E}_{q(\pi)} \left[\ln q(\pi) - \left(- \sum_t \mathbb{G}_{\tau}(\pi) \right) \right] = D_{KL} \left[q(\pi) \parallel e^{-\sum_t \mathbb{G}_{\tau}(\pi)} \right] \\&\text{where } \mathbb{G}_{\tau}(\pi) = \mathbb{E}_{q(o_{\tau}, s_{\tau} | \pi)} [\ln q(s_{\tau} | \pi) - \ln \tilde{p}(o_{\tau}, s_{\tau})]\end{aligned}$$

¹¹David M Blei, Alp Kucukelbir, and Jon D McAuliffe. “Variational inference: A review for statisticians”. In: *Journal of the American statistical Association* 112.518 (2017), pp. 859–877.

¹²Beren Millidge, Alexander Tschantz, and Christopher L Buckley. “Whence the expected free energy?” In: *Neural Computation* 33.2 (2021), pp. 447–482.

Expected free energy

- thus optimising \mathbb{G} results in

$$q^*(\pi) = \text{SoftMax}\left(-\sum_t^T \mathbb{G}_\tau(\pi)\right)$$

- EFE at time τ , \mathbb{G}_τ , can be further decomposed as:

$$\begin{aligned}\mathbb{G}_\tau(\pi) &= \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau, s_\tau)] \\ &\approx \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau) - \ln q(s_\tau | o_\tau)] \\ &\approx \underbrace{-\mathbb{E}_{q(o_\tau | \pi)} [\ln \tilde{p}(o_\tau)]}_{\text{Extrinsic Value}} - \underbrace{\mathbb{E}_{q(o_\tau)} D_{KL}[q(s_\tau | o_\tau) \| q(s_\tau | \pi)]}_{\text{Epistemic Value}}\end{aligned}$$

- thus, optimal policies are obtained by minimising the sum of the expected free energies
- EFE is estimated using the generative model to roll out predicted futures, and compute the EFE of those futures

Expected free energy

- similar to variational free energy, EFE can be also decomposed as:

$$\begin{aligned}\mathbb{G}_\tau(\pi) &= \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau, s_\tau)] \\ &\approx \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau) - \ln q(s_\tau | o_\tau)] \\ &\approx \underbrace{-\mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln \tilde{p}(o_\tau)]}_{\text{Extrinsic Value}} - \underbrace{\mathbb{E}_{q(o_\tau)} D_{KL} [q(s_\tau | o_\tau) \| q(s_\tau | \pi)]}_{\text{Epistemic Value}}\end{aligned}$$

- note that an **epistemic uncertainty** refers to the deficiencies by a lack of knowledge or information; reducible with more data or a better model
- the second term measures the reduced entropy of s_τ when observed $o_\tau \rightarrow \max$ its value \rightarrow we intend to choose the policy such that $H[q(s|\pi)]$ is high and strong dependency between states and observation (aka $H[q(s|o)]$ is low)

Expected free energy

- one can also decompose it into the following:

$$\begin{aligned}\mathbb{G}_\tau(\pi) &= \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau, s_\tau)] \\ &= \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau | s_\tau) - \ln p(s_\tau)] \\ &= \underbrace{-\mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln \tilde{p}(o_\tau | s_\tau)]}_{\text{Accuracy}} + \underbrace{\mathbb{E}_{q(o_\tau | s_\tau)} [D_{KL}[q(s_\tau | \pi) \| p(s_\tau)]]}_{\text{Complexity}}\end{aligned}$$

- or this (typically used for computation):

$$\begin{aligned}\mathbb{G}_\tau(\pi) &= \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\ln q(s_\tau | \pi) - \ln \tilde{p}(o_\tau) - \ln q(s_\tau | o_\tau)] \\ &= \mathbb{E}_{q(o_\tau, s_\tau | \pi)} [\cancel{\ln q(s_\tau | \pi)} - \ln \tilde{p}(o_\tau) - \ln p(o_\tau | s_\tau) \\ &\quad - \cancel{\ln q(s_\tau | \pi)} + \ln q(o_\tau)] \\ &= \underbrace{D_{KL}[q(o_t | \pi) \| \tilde{p}(o_t)]}_{\text{Expected Cost}} + \underbrace{\mathbb{E}_{q(s_t | \pi)} [H[p(o_t | s_t)]]}_{\text{Expected Ambiguity}}\end{aligned}$$

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example**
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

A simple example¹³

- suppose there is an agent interacts with an environment
- **perception:**
 - the environment has two hidden states $s \in \{1, 2\}$, e.g., *there is food in your stomach (1) or not (2)*
 - while s is NOT measurable, there is an observation $o = \{1, 2\}$, e.g., you feeling fed (1) or hungry (2)
 - assume $p(o|s)$ is given as a 2x2 likelihood matrix (called "A" matrix) which maps states to observations, e.g., if you fed, you have food and vice versa
- **control:**
 - transition prob. $p(s_t|s_{t-1}, u)$ maps the previous state to the next one, depending on action $u = \{u_1, u_2\}$
 - parameterised by a separate 2x2 transition matrix ("B") for each action, e.g., *either go get food (u_1) or do nothing (u_2); if u_1 , will have food in the next state, regardless of whether we have it now, and vice versa*

¹³<https://medium.com/@solopchuk/tutorial-on-active-inference-30edcf50f5dc>

A simple example¹⁴

- preference (as a form of objective or reward)
 - we also have prior preferences $\tilde{p}(o)$, e.g., we like to be fed and not hungry, so we assign a higher probability to observation $o = 1$ fed. In other words, we express preferences over observations as probability $\tilde{p}(o)$

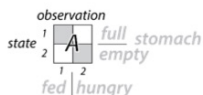
prior
preferences

$$p(o)$$



likelihood
(state-observation
mapping)

$$p(o|s)$$



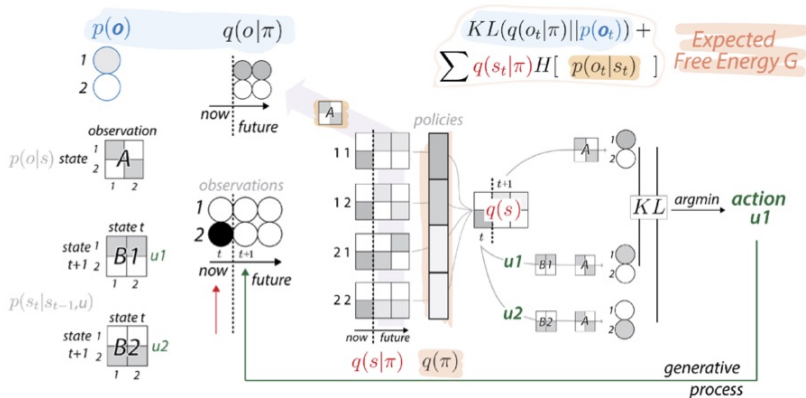
transition
(state - next state
mapping)

$$p(s_t | s_{t-1}, u)$$



¹⁴[https:](https://medium.com/@solopchuk/tutorial-on-active-inference-30edcf50f5dc)

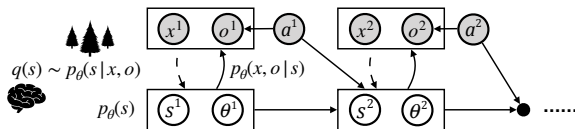
A simple example¹⁵



action selection: option 1) $q(s|\pi) \rightarrow q(o|\pi) \rightarrow \mathbb{G}(\pi) \rightarrow q(\pi)$
 option 2 $u_{t+1}^* = \arg \max_u D_{KL}[AS_{t+1}||A(B(u))S_t]$, Bayesian model average method

¹⁵<https://medium.com/@solopchuk/tutorial-on-active-inference-30edcf50f5dc>

Unifying perception and control (Bayesian brain)¹⁶



- the latent state s : true world configuration such as pixel assignment, the optimality o is a binary variable
- the perception model includes a bottom-up recognition model $q(s)$ and a top-down generative model $p(x, o, s)$ (decomposed into the likelihood $p(x, o|s)$ and the prior belief $p(s)$)
- the prior knowledge θ represents the physical law of the environment (the property of each object)
- control is performed by taking an action a to change the environment state.

¹⁶Minne Li et al. "Joint Perception and Control as Inference with an Object-based Implementation". In: (2020).

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference**
- 7 Multi-agent Variational Bayes
- 8 Conclusions

Bayesian decision principle

making decision under uncertainty can be interpreted as maximising expected utility in the face of uncertainty¹⁷

- the uncertainty is captured by the (hidden) *state of nature*: $z \in \mathbb{Z}$
- decisions (aka *actions*): $a \in \mathbb{A}$
- reward: $R(z, a)$
- the posterior distribution $p(z|o)$, where we can perform a statistical investigation to obtain information (denoted as $o = (o_1, o_2, \dots, o_n) \in O$) about the nature state θ

(Conditional Bayes Decision Principle)

$$a^*|o = \arg \max_a E_{\theta \sim p(z|o)}[R(o, a)], \quad (1)$$

where the principle is the only fundamentally correct analysis¹⁸.

Yet, the posterior can be difficult to obtain practically

¹⁷ John Von Neumann and Oskar Morgenstern. *Theory of Games and Economic Behavior*. Princeton University Press, 1953.

¹⁸ James O Berger. *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.

The duality between inference and control

- the filtering problem is shown to be the dual of the noise-free regulator (control) problem¹⁹

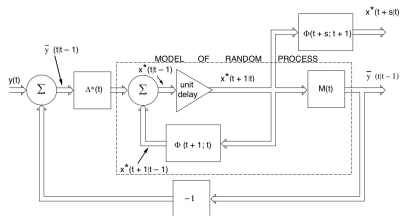


Figure: Optimal filter

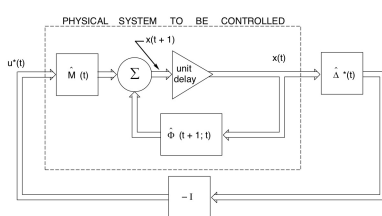
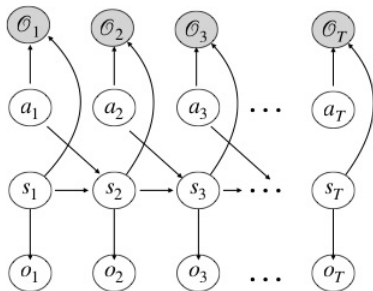


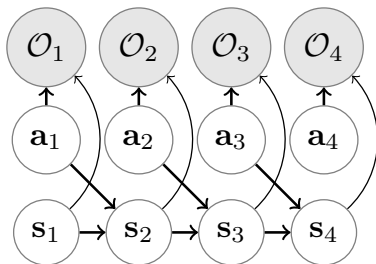
Figure: Optimal controller

¹⁹Rudolph Emil Kalman. "A new approach to linear filtering and prediction problems". In: (1960).

Control as inference:²⁰



(a) states s are hidden



(b) states s are observable

- the idea: add an optimality variable $\mathcal{O}_t = 1$ and assume a biased probability with reward r :

$$p(\mathcal{O}_t = 1 \mid s_t, a_t) = \exp(r(s_t, a_t))$$

²⁰Sergey Levine. “Reinforcement learning and control as probabilistic inference: Tutorial and review”. In: *arXiv preprint arXiv:1805.00909* (2018); Beren Millidge et al. “On the relationship between active inference and control as inference”. In: *International Workshop on Active Inference*. Springer. 2020, pp. 3–11.

Control as inference:²¹

- let us look at an MDP model (states are observable)
- the evidence: $\mathcal{O}_t = 1$ for all $t \in \{1, \dots, T\}$; thus ELBO is

$$\begin{aligned}\log p(\mathcal{O}_{1:T}) &= \log \iint p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\ &= \log \iint p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \frac{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\ &= \log E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[\frac{p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \right] \\ &\geq E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} [\log p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) - \log q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})]\end{aligned}$$

²¹Sergey Levine. "Reinforcement learning and control as probabilistic inference: Tutorial and review". In: *arXiv preprint arXiv:1805.00909* (2018).

Control as inference:²²

- we make use of the evidence $p(\mathcal{O}_{1:T}, s_{1:T}, a_{1:T}) = \left[p(s_1) \prod_{t=1}^T p(s_{t+1} | s_t, a_t) \right] \exp \left(\sum_{t=1}^T r(s_t, a_t) \right)$ and
- factorise the variational distribution as $q(\tau) \equiv q(s_{1:T}, a_{1:T}) = q(s_1) \prod_{t=1}^T q(s_{t+1} | s_t, a_t) q(a_t | s_t)$
- this leads to the final ELBO: $\log p(\mathcal{O}_{1:T}) \geq E_{q(s_{1:T}, a_{1:T})} \left[\sum_{t=1}^T r(s_t, a_t) - \log q(a_t | s_t) \right]$
- where we define $q(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t)$

²²Sergey Levine. “Reinforcement learning and control as probabilistic inference: Tutorial and review”. In: *arXiv preprint arXiv:1805.00909* (2018).

Control as inference:²³

- without using temporal mean-field factorisation as Active Inference, one can derive Value Iteration similar to a standard RL solution

- the ELBO can be decomposed recursively as:

$$\begin{aligned} & E_{q(s_t, a_t)} [r(s_t, a_t) - \log \pi(a_t | s_t)] + \\ & E_{q(s_t, a_t)} [E_{s_{t+1} \sim p(s_{t+1} | s_t, a_t)} [V(s_{t+1})]] = \\ & E_{q(s_t)} \left[-D_{\text{KL}} \left(\pi(a_t | s_t) \parallel \frac{1}{\exp(V(s_t))} \exp(Q(s_t, a_t)) \right) + V(s_t) \right] \end{aligned}$$

- where we define:

$$\begin{aligned} Q(s_t, a_t) &\equiv r(s_t, a_t) + E_{s_{t+1} \sim p(s_{t+1} | s_t, a_t)} [V(s_{t+1})] \\ V(s_t) &\equiv \log \int_{\mathcal{A}} \exp(Q(s_t, a_t)) da_t \end{aligned}$$

$$\rightarrow \pi(a_t | s_t) = \exp(Q(s_t, a_t) - V(s_t))$$

²³Sergey Levine. "Reinforcement learning and control as probabilistic inference: Tutorial and review". In: *arXiv preprint arXiv:1805.00909* (2018).

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes**
- 8 Conclusions

Agent modelling with maximum entropy objective²⁴

- each agent pursues the maximal cumulative reward

$$\max \eta^i(\pi_\theta) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^t R^i(s_t, a_t^i, a_t^{-i}) \right], \quad (2)$$

with actions (a_t^i, a_t^{-i}) sampled from policy $(\pi_{\theta^i}^i, \pi_{\theta^{-i}}^{-i})$

- a strategy profile $(\pi^{1*}, \dots, \pi^{n*})$ reaches optimum when:

$$\begin{aligned} & \mathbb{E}_{s \sim p_s, a_t^{i*} \sim \pi^{i*}, a_t^{-i*} \sim \pi^{-i*}} \left[\sum_{t=1}^{\infty} \gamma^t R^i(s_t, a_t^{i*}, a_t^{-i*}) \right] \\ & \geq \mathbb{E}_{s \sim p_s, a_t^i \sim \pi^i, a_t^{-i} \sim \pi^{-i}} \left[\sum_{t=1}^{\infty} \gamma^t R^i(s_t, a_t^i, a_t^{-i}) \right] \quad (3) \\ & \forall \pi \in \Pi, i \in (1 \dots n), \end{aligned}$$

where $\pi = (\pi^i, \pi^{-i})$ and agent i 's optimal policy is π^{i*}

²⁴Zheng Tian et al. "A regularized opponent model with maximum entropy objective". In: *IJCAI* (2019).

Agent modelling with maximum entropy objective

- a lower bound on the likelihood of optimality of agent i :

$$\log P(\mathcal{O}_{1:T}^i = 1 | \mathcal{O}_{1:T}^{-i} = 1) \geq \sum_t \mathbb{E}_{(s_t, a_t^i, a_t^{-i}) \sim q} [R^i(s_t, a_t^i, a_t^{-i}) + H(\pi(a_t^i | s_t, a_t^{-i})) - D_{\text{KL}}(\rho(a_t^{-i} | s_t) || P(a_t^{-i} | s_t))] \quad (4)$$

$$= \sum_t \mathbb{E}_{s_t} [\underbrace{\mathbb{E}_{a_t^i \sim \pi, a_t^{-i} \sim \rho} [R^i(s_t, a_t^i, a_t^{-i}) + H(\pi(a_t^i | s_t, a_t^{-i}))]}_{\text{MEO}}] - \underbrace{\mathbb{E}_{a_t^{-i} \sim \rho} [D_{\text{KL}}(\rho(a_t^{-i} | s_t) || P(a_t^{-i} | s_t))]}_{\text{Regulariser of } \rho}]. \quad (5)$$

- $\rho(a_t^{-i} | s_t, o_t^{-i} = 1)$ is agent i 's opponent model
- $\pi(a_t^i | s_t, a_t^{-i}, o_t^i = 1, o_t^{-i} = 1)$ is the agent i 's conditional policy at optimum ($o_t^i = o_t^{-i} = 1$) and
- $P(a_t^{-i} | s_t, o_t^{-i} = 1)$ is the prior of opponent model
- the prior $P(a_t^{-i} | s_t, o_t^{-i} = 1)$ is estimated empirically
- we drop ($o_t^i = 1, o_t^{-i} = 1$) in π, ρ and $P(a_t^{-i} | s_t)$

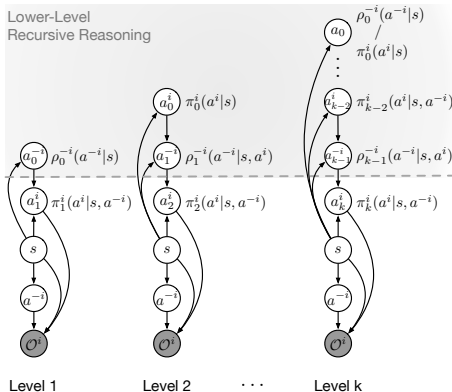
Recursive reasoning²⁵



- Example: in the “beauty contest” game, players are asked to pick numbers from 0 to 100, and the player whose number is closest to $2/3$ of the average wins a prize

²⁵Colin F Camerer, Teck-Hua Ho, and Juin-Kuan Chong. “A cognitive hierarchy model of games”. In: *The Quarterly Journal of Economics* 119.3 (2004), pp. 861–898.

Multi-agent generalized recursive reasoning²⁶



- unobserved opponent policies are approximated by ρ^{-i}
- agent i recursively reasons about opponents (grey area)
- in the recursion, agents with higher-level beliefs take the best response to the lower-level thinkers' actions.

²⁶Ying Wen et al. "Modelling Bounded Rationality in Multi-Agent Interactions by Generalized Recursive Reasoning". In: *IJCAI* (2020).

Table of Contents

- 1 Learning in biological and computerised systems
- 2 Perception-control loop: the problem
- 3 Active inference: perception
- 4 Active inference: perception-control loop
- 5 An illustrative example
- 6 Control as inference
- 7 Multi-agent Variational Bayes
- 8 Conclusions

Remarks

- AI: reasoning under uncertainty
- a learning system is an information system
- pattern recognition and decision making problems (include multiagent learning) can be unified and modelled by probabilistic inference such as Variational Bayes (e.g., active- inference or inference-as-control)
- as such information theory can be directly utilised
- however, we call for new research on information theory that deals with "intrinsic information exchange"

References I



James O Berger. *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.



David M Blei, Alp Kucukelbir, and Jon D McAuliffe. “Variational inference: A review for statisticians”. In: *Journal of the American statistical Association* 112.518 (2017), pp. 859–877.



Colin F Camerer, Teck-Hua Ho, and Juin-Kuan Chong. “A cognitive hierarchy model of games”. In: *The Quarterly Journal of Economics* 119.3 (2004), pp. 861–898.

References II



Karl Friston. “The free-energy principle: a unified brain theory?” In: *Nature reviews neuroscience* 11.2 (2010), pp. 127–138.



Michael I Jordan et al. “An introduction to variational methods for graphical models”. In: *Machine learning* 37.2 (1999), pp. 183–233.



Rudolph Emil Kalman. “A new approach to linear filtering and prediction problems”. In: (1960).



Sergey Levine. “Reinforcement learning and control as probabilistic inference: Tutorial and review”. In: *arXiv preprint arXiv:1805.00909* (2018).

References III



Minne Li et al. “Joint Perception and Control as Inference with an Object-based Implementation”. In: (2020).



David Marr. “Vision: A computational investigation into the human representation and processing of visual information”. In: *Inc., New York, NY* 2.4.2 (1982).



Beren Millidge et al. “On the relationship between active inference and control as inference”. In: *International Workshop on Active Inference*. Springer. 2020, pp. 3–11.

References IV



Beren Millidge, Alexander Tschantz, and Christopher L Buckley. “Whence the expected free energy?” In: *Neural Computation* 33.2 (2021), pp. 447–482.



Ivan Petrovitch Pavlov and William Gantt. “Lectures on conditioned reflexes: Twenty-five years of objective study of the higher nervous activity (behaviour) of animals.”. In: (1928).



Eliezer Sternberg. *NeuroLogic: The Brain's Hidden Rationale Behind Our Irrational Behavior*. Vintage, 2016.



Zheng Tian et al. “A regularized opponent model with maximum entropy objective”. In: *IJCAI* (2019).

References V



John Von Neumann and Oskar Morgenstern.
Theory of Games and Economic Behavior.
Princeton University Press, 1953.



Ying Wen et al. “Modelling Bounded Rationality
in Multi-Agent Interactions by Generalized
Recursive Reasoning”. In: *IJCAI* (2020).