

Reinforcement learning for large state spaces ("RL in practice")

- so far: by learning state-value function $Q^*(s, a)$, we can derive an optimal policy:
$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

- problem: Q^* was represented as a table of size $|S| \times |A|$

this may be impractical: - too big to be stored explicitly
 - impossible to get enough training data to fill each entry (\Rightarrow analogous to multi-dimensional histogram)

\Rightarrow learn Q^* (or rather: \hat{Q} approx.) by regression:

- generalize from events that we have seen in τS to events we have not seen, provided there is sufficient similarity
 - define features:
 - $X \equiv s$ (or $X \equiv o$ in partially observable problems) and have regression alg. figure out what is important in s
- example: deep Q-learning of Atari games (famous paper by Deep Mind): s is content of the screen, $X \equiv s$, deep neural network automatically learns more informative features from the raw pixel values

- $X = \psi(s)$ or $X = \psi(o)$ with some hand-crafted or learned feature extraction function $\psi(\cdot)$

- define response:
$$Y = \mathbb{E} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t-1} r_{t'} \right]$$
 expected return

examples:

- MC value estimation:

$$Y_{\tau,t} = \sum_{t'=t+1}^{\infty} \gamma^{t'-t-1} \underbrace{R_{\tau,t'}}_{\text{actually observed rewards in timestep } t'}$$

actually observed rewards in timestep t'
of episode τ (τ = game index, t = move index)

- TD (temporal difference) Q-learning

$$Y_{\tau,t} = R_{\tau,t} + \gamma \max_a Q(s_{\tau,t+1}, a)$$

↑ current guess of Q-function

- n-step TD Q-learning

$$Y_{\tau,t} = \sum_{t'=t+1}^{t+n} \gamma^{t'-t-1} R_{\tau,t'} + \gamma^n \max_a Q(s_{\tau,t+n}, a)$$

- SARSA:

$$Y_{\tau,t} = R_{\tau,t+1} + \gamma Q(s_{\tau,t+1}, a_{\tau,t+1})$$

• optimization algorithm: alternating optimization
repeat until convergence:

(1) compute target reference $Y_{\tau,t}$ using current guess for Q (2) optimize Q using squared loss

→ new guess $Q' = \arg \max_{\tilde{Q}} \sum_{TS} \mu(s_{\tau,t}) (Y_{\tau,t} - \tilde{Q}(s_{\tau,t}, a_{\tau,t}))^2$

↑ weights

usually solved by batch gradient descent
random subset of TS → improve our current by updating parameters

- common update strategies:
 - "online": apply a gradient step after every move
 - "episodic": gradient step(s) after every game
 - "offline": TS is an "experience buffer" of size K
 (when new experiences arrive (= moves are executed), new data are added to TS and ~~if~~ old possibly removed (if K would otherwise be exceeded))
- elements in TS: $(s_{\tau,t}, a_{\tau,t}, s_{\tau,t+1}, r_{\tau,t+1})$

in each update round: ("experience replay")

- sample a batch of size k' from TS
- compute $x_{\tau,t}, y_{\tau,t}$ for selected elements
- perform gradient descent with squared loss

variant: prioritized experience replay:

- sample batch of size $k'' > k'$
- compute $x_{\tau,t}, y_{\tau,t}$ and residual $y_{\tau,t} - Q(x_{\tau,t}; \theta_{\tau,t})$
- only keep k' instances with biggest squared residual $\hat{=} p(s_{\tau,t})$ big for difficult decisions
- gradient descent with squared loss

alternating optimization: keep two models: old guess $Q \Rightarrow$ compute $y_{\tau,t}$
 working guess $Q' \Rightarrow$ apply updates to
 every once in a while, update $Q \leftarrow Q'$ (frequency of this is a hyper-parameter)

Models for Q-function regression:

- linear value approximation

$$Q(X(s), a) = X(s) \cdot \beta_a \quad \text{for each action: one parameter vector } \beta_a$$

this only works if we define good features! $X = X(s)$

(generally a hard engineering problem)

training updates: training batch $B = (\beta_{a_1}, \dots, \beta_{a_N})$

training rate (hyperparameter) \uparrow sub batch for each action

\Rightarrow gradient update:

$$\beta_a \leftarrow \beta_a + \frac{\alpha}{N} \sum_{T: t \in B_a} X(s_{T,t})^T (y_{T,t} - X(s_{T,t}) \cdot \beta_a)$$

\uparrow size of sub batch B_a \uparrow fixed feature extraction functions

- feature definition: usual preprocessing methods like PCA and clustering and their kernel-version can be applied before manual feature design (or sometimes instead)
- (monte carlo) tree search: roll-out each game (from current state) to the end with random actions \Rightarrow game out-come gives information on the current state
- data augmentation: exploit the symmetries of a game: in Bomberman, rotating by 90° or mirroring are still valid game states
 \Rightarrow each training move gives information about 8 game states
 (alternative: define features to be invariant under rotation and mirroring)

- Regression Forest: two possibilities
 - 1) use as input $(X, a) \Rightarrow$ predict $Q(X, a)$
 - 2) use as input X and predict a vector of values, for one entry for each action
- problem: updating an existing forest is difficult (on-line learning of RF open prob.)
 - \Rightarrow learn a new RF from scratch in every iteration & batch size must be big enough to get good results &
 - \Rightarrow in ensemble of trees: use different batches for every tree, ~~we~~ can even ^{retain} keep some trees from the current forest and ~~update~~ only the rest
- Gradient Boosting: alternative method to create an ensemble:
 - RF: each tree is independent of all other trees
 - GB: a new member of the ensemble is trained to correct the mistakes of the ensemble so far $\Rightarrow h_{i+1} = f(TS, \{h_i, \sum_{j=1}^i\})$
 new member conditioned on old members, not independent
 - ~~modified~~ modification to alternating optimization alg.:
 instead of doing gradient updates on the existing model, train a new ensemble member on the current residuals, keeping existing ensemble fixed
 - current residuals $\beta_{T,t} = Y_{T,t} - Q(S_{T,t}, a_{T,t})$
 \leftarrow current ensemble response
 - $$h_{i+1} = \arg \min_h \sum_{\text{batch}} w_{T,t} (\beta_{T,t} - h(S_{T,t}, a_{T,t}))^2$$
 - \Rightarrow new ensemble $Q'(S, a) = Q(S, a) + \sum_{i=1}^n w_i h_i(S, a)$
 \leftarrow weight

advantage: h_i can be very simple non-linear expressors, e.g.

- decision stumps ($\hat{=}$ decision trees with very small depth, e.g. 1)
- rounding of a ~~linear~~ linear model ($\text{sign}(x \cdot p)$)

weight schedule: $w_i = \frac{w_0}{1 + p \cdot i}$: new ensemble members get lower weight

• Deep Q-learning (DQN $\hat{=}$ deep Q network)

- doesn't define features manually, but sets $x \equiv s$ (or $x \equiv 0$) and have the first layers of the network learn good features

- $Q(s, a)$ is represented as $Q_\theta(s, a)$ with θ : the parameters of the neural network

value estimates computed with stored guess Q_θ

- gradient updates: $\theta \leftarrow \theta - \alpha \nabla_\theta \left[\frac{1}{|B|} \sum_{\tau_i \in B} (Y_{\tau_i} - Q_\theta(s_{\tau_i}, a_{\tau_i}))^2 \right]$

use autograd in pytorch or tensorflow

current parameters look as before

(can use a stored older guess to compute $Y_{\tau_i} = r_{\tau_i+1} + \gamma Q_\theta(s_{\tau_i+1}, a_{\tau_i+1})$)

SARSA

$Y_{\tau_i} = r_{\tau_i+1} + \gamma \max_a Q_\theta(s_{\tau_i+1}, a)$

Q-learning

every once in a while, update $\theta' \leftarrow \theta$

Training strategies

- self play: by playing with your current agent against itself or an older version of itself, create arbitrary amount of training data
 - ~~crucial~~ crucial: ensure that ~~an~~ action selection is partially random to avoid bad convergence to local optima \Rightarrow do not use greedy policy

$$\pi(s) = \arg \max_a Q(s, a)$$

common choice: - Softmax policy $\pi(s) \sim \text{Softmax}_a Q(s, a)$

$$\text{Softmax}(u) = \frac{e^{u_i}}{\sum_{i'} e^{u_{i'}}} \quad \text{or with "temperature"} \quad \frac{e^{u_i/\beta}}{\sum_{i'} e^{u_{i'}/\beta}}$$

$$\hat{=} \text{discrete}_a(\exp(Q(s, a)))$$

($\beta = 0 \Rightarrow$ greedy policy, $\beta \Rightarrow \infty$ uniform random actions)

- IAUU exploration: improved ϵ -greedy strategy

$$a \sim \begin{cases} \arg \max_a Q(s, a) & \text{with prob } 1 - \epsilon \\ \text{softmax}(-\text{count}(s, a)/\beta) & \text{with prob } \epsilon \end{cases}$$

prefer actions we haven't tried much in the past

- reward shaping: often, rewards are sparse ($\hat{=}$ rare, e.g. just win/loss at the end of each episode)
 - \Rightarrow hard to learn, slow convergence

- idea: add auxiliary rewards during training to speed-up convergence

but: bad auxiliary rewards lead to bad policies

Common failure modes: - infinite loop of positive local rewards

A = start	B
C = monster	D = treasure

$r_{A \rightarrow B} = r_{B \rightarrow A} = 1$
 \Rightarrow agent can make infinite return by going in an infinite loop $A \rightarrow B$

- reward gaming: intended effect of the auxiliary reward is put to ridiculous extremes

ex: robot should learn to stretch out its arm \Rightarrow auxiliary reward for placing an object far away from the "body"

unintended consequence: robot learns to throw objects instead of placing them softly

- theory: potential-based reward shaping is necessary and sufficient for the optimal policy to remain invariant

shaped rewards:
$$r'_{\tau, t+1} = r_{\tau, t+1} + F(s_{\tau, t}, s_{\tau, t+1})$$

$$F(s_{\tau, t}, s_{\tau, t+1}) = \gamma \phi(s_{\tau, t+1}) - \phi(s_{\tau, t})$$

↑
 potential: only depends on a single state and independent of actions

- under a greedy policy $\phi(s) = V^*(s)$ (optimal value function) is best

- but: in an exploration/exploitation setting during training, this may lead to many losses (eaten by monster) : to avoid this, the auxiliary rewards should be biased towards "safe" behavior \Rightarrow good $\phi(s)$ are an "art".