

value function: $v_\pi: S \rightarrow \mathbb{R}, v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$
 Bellman equation:
 $v_\pi(s) = \mathbb{E}_\pi[G_t + \gamma v_\pi(S_{t+1}) | S_t = s]$
 $= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$
 $= \sum_a \pi(a|s) \sum_{s'} p(s', r | s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']]$
 $= \sum_a \pi(a|s) \sum_{s'} p(s', r | s, a) [r + \gamma v_\pi(s')]$, for all $s \in S$
 Action-value function: $q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$
 Greedy policy: $\pi(s) = \arg \max_a q_\pi(s, a) = \arg \max_a [r(s, a) + \gamma v_\pi(p(s, a))]$
 Discounted function: $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

TD(0): $V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$ choose A_t acc to policy
 SARSA: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$
 Q-learning: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$
 Expected reward $J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)}[r(\tau)] = \int_{\tau} r(\tau) p(\tau; \theta) d\tau$
 REINFORCE trick: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau = \int_{\tau} r(\tau) p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} d\tau$
 $= \int_{\tau} r(\tau) p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) d\tau$
 we can estimate $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) d\tau$ (RL)
 $= \mathbb{E}_{\tau \sim p(\tau; \theta)}[r(\tau) \nabla_{\theta} \log p(\tau; \theta)]$ with Monte Carlo sampling
 independence from transition probabilities can be seen as below

$p(\tau; \theta) = \prod_{t=0}^{\infty} p(S_{t+1} | S_t, a_t) \pi(a_t | S_t)$
 $\log p(\tau; \theta) = \sum_{t=0}^{\infty} \log p(S_{t+1} | S_t, a_t) + \log \pi(a_t | S_t)$
 $\nabla_{\theta} \log p(\tau; \theta) = \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi(a_t | S_t)$
 $\nabla_{\theta} J(\theta) \approx \sum_{t=0}^{\infty} A^{\pi}(S_t, a_t) \nabla_{\theta} \log \pi(a_t | S_t)$, $A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$

- bootstrapping: not waiting for the final outcome
 both SARSA and Q-learning use e-greedy methods
 actor-critic combines policy gradients and Q-learning
 actor: policy, critic: q-function \rightarrow tells how good the action
 \rightarrow decides which action to take was and now actor should adjust

TRPO (trust region policy opt.): Del of old and new policies, hard constraint, idea is to control the step length (don't change too much)
 PPO: (proximal policy opt.): TRPO with 1st order optimization and soft constraint, we can do better using clipped objective rather than the DPPG (deep deterministic policy gradients): actor-critic off-policy
 DRL, combines deep Q-learning and DPG, random noise when selecting a

AR models: ppo is tractable, easy to train, easy (but slower) to sample, no natural latent variable repr (but doable i.e. VAE, STCN)
 - major downside is sequential (slow) generation
 - work for both cont. and disc. data (GANS hard to learn disc.)
 - training more stable than GANS
 - because we train via NLL, we have a direct measure for comparison. Thus, also makes it straightforward to apply on domains such as compression and prob. planning/exploration

Dynamic programming
 + exact methods
 + guaranteed to converge (finite iter.)
 - need to know transition matrix
 + easy to implement
 - need to iterate over whole state space
 - requires memo + prop. to size of state space
 • value iter. typically more efficient
 • policy eval expensive as it may be also iterative over state space
 • bootstraps, do not sample
 • uses value of neighboring states to update current state (given π)

Monte Carlo sampling (RL)
 + unbiased estimate
 + no need to know system dynamics (experience-based)
 - high variance
 - exploration/exploitation dilemma
 - need to know state
 - slow for long episodes (MC must wait until the final outcome to update)
 + easy to focus on small subset of states
 • first-visit and every-visit MC both converge to vals if visits \rightarrow infinity
 • does not bootstrap, samples

Temporal Differences
 + less variance than MC due to bootstrap
 + more sample efficient
 + no need to know prob. matrix
 - biased due to bootstrapping
 - exploration/exploitation dilemma
 - can behave poorly in stochastic envs
 • combination of DP and MC ideas
 • bootstraps and samples
 - exploitation may lead to local minima

Policy optimization
 - directly optimizes desired $q_{\pi}(s, a)$
 - more compatible with modern ML (NNs)
 - more versatile and flexible
 - more competitive with aux. objectives

Dynamic programming
 - indirect, exploits problem structure and self-consistency
 - more sample efficient (than they were)
 - more compatible with off-policy and exploration

Issues in DRL research
 - Exploitation: agents can get stuck in local minima (exploit behaviour that is irreversible)
 - Reward: modifying the shaped reward may lead to undesirable behaviour (needs careful design)
 - Domain knowledge: even small imper. differences may have a big impact
 - Sample efficiency: require enormous amount of samples, or we update and improve policies we collect better data and throw away the rest

softmax $Q_t = \frac{e^{Q_t}}{\sum_i e^{Q_i}} \rightarrow Q_t = \frac{e^{Q_t}}{\sum_i e^{Q_i}}$

$\frac{\partial Q_t}{\partial Q_i} = \frac{\partial}{\partial Q_i} \left(\frac{e^{Q_t}}{\sum_j e^{Q_j}} \right) = \frac{-e^{Q_i} \cdot e^{Q_t}}{(\sum_j e^{Q_j})^2} = -Q_i \cdot Q_t$ (for $i \neq t$)
 $\frac{\partial Q_t}{\partial Q_i} = \frac{e^{Q_i} \sum_j e^{Q_j} - e^{Q_i} e^{Q_i}}{(\sum_j e^{Q_j})^2} = \frac{e^{Q_i} (\sum_j e^{Q_j} - e^{Q_i})}{(\sum_j e^{Q_j})^2}$
 $= \frac{e^{Q_i}}{\sum_j e^{Q_j}} \left(\frac{\sum_j e^{Q_j}}{\sum_j e^{Q_j}} - \frac{e^{Q_i}}{\sum_j e^{Q_j}} \right) = Q_i (1 - Q_i)$ (for $i = t$)

softmax grad

$\delta_i^{(l)}$: gradient respect to the i th unit in the l layer
 $\delta^{(l)}$: gradient that flows from layer $l+1$ to l
 $\delta_i^{(l)} = \frac{\partial C}{\partial z_i^{(l)}} = \sum_{j=1}^N \frac{\partial C}{\partial z_j^{(l+1)}} \frac{\partial z_j^{(l+1)}}{\partial z_i^{(l)}} = \sum_{j=1}^N \delta_j^{(l+1)} \frac{\partial z_j^{(l+1)}}{\partial z_i^{(l)}}$
 $\delta^{(l)} = \frac{\partial C}{\partial z^{(l)}} = \sum_{j=1}^N \delta_j^{(l+1)} \frac{\partial z_j^{(l+1)}}{\partial z^{(l)}} = \frac{\partial C}{\partial z^{(l+1)}} \frac{\partial z^{(l+1)}}{\partial z^{(l)}}$
 $z_j^{(l+1)} = \sigma(\theta_j^{(l+1)} \cdot z^{(l)}) \Leftrightarrow z^{(l+1)} = \sigma(Q^{(l+1)} \cdot z^{(l)})$
 $\langle 1 \rangle \quad \langle 1, m \rangle \quad \langle m, 1 \rangle \quad \langle n, 1 \rangle \quad \langle n, m \rangle \quad \langle m, n \rangle$
 $\frac{\partial C}{\partial z_i^{(l)}} = \frac{\partial C}{\partial z_i^{(l)}} \frac{\partial z_i^{(l+1)}}{\partial z_i^{(l)}} \quad (i\text{th unit in the } l \text{ layer})$
 $\frac{\partial C}{\partial z^{(l)}} = \sum_{j=1}^N \frac{\partial C}{\partial z_j^{(l+1)}} \frac{\partial z_j^{(l+1)}}{\partial z^{(l)}} = \sum_{j=1}^N \delta_j^{(l+1)} \frac{\partial z_j^{(l+1)}}{\partial z^{(l)}}$

$z^{(l)} = W^{(l)} z^{(l-1)} + b^{(l)} = \sum_m \sum_n W_{m,n}^{(l)} z_{m,n}^{(l-1)} + b^{(l)}$ (fwd pass)
 $\delta_{i,j}^{(l-1)} = \frac{\partial C}{\partial z_{i,j}^{(l-1)}} = \sum_i \sum_j \frac{\partial C}{\partial z_{i,j}^{(l)}} \frac{\partial z_{i,j}^{(l)}}{\partial z_{i,j}^{(l-1)}} = \sum_i \sum_j \delta_{i,j}^{(l)} \frac{\partial}{\partial z_{i,j}^{(l-1)}} \sum_m \sum_n W_{m,n}^{(l)} z_{m,n}^{(l-1)} + b^{(l)}$
 $= \sum_i \sum_j \delta_{i,j}^{(l)} \frac{\partial}{\partial z_{i,j}^{(l-1)}} \sum_m \sum_n W_{m,n}^{(l)} z_{m,n}^{(l-1)} + b^{(l)}$
 $= \sum_i \sum_j \delta_{i,j}^{(l)} W_{i,j}^{(l)} (change of vars: i = i' - m \rightarrow m = i' - i)$
 $= \sum_i \sum_j \delta_{i,j}^{(l)} W_{i,j}^{(l)} = \delta^{(l)} * W_{i,j}^{(l)} = \delta^{(l)} * rot_{180}(W_{i,j}^{(l)})$
 $\frac{\partial C}{\partial W_{m,n}^{(l)}} = \sum_i \sum_j \frac{\partial C}{\partial z_{i,j}^{(l)}} \frac{\partial z_{i,j}^{(l)}}{\partial W_{m,n}^{(l)}} = \sum_i \sum_j \delta_{i,j}^{(l)} \frac{\partial}{\partial W_{m,n}^{(l)}} \left(\sum_m \sum_n W_{m,n}^{(l)} z_{m,n}^{(l-1)} + b^{(l)} \right)$
 $= \sum_i \sum_j \delta_{i,j}^{(l)} z_{i-m, j-n}^{(l-1)} = \delta^{(l)} * z_{i-m, j-n}^{(l-1)} = \delta^{(l)} * rot_{180}(z_{i-m, j-n}^{(l-1)})$

chain backprop

Prx2prx: $L(G, D) = L_{GAN}(G, D) + \lambda L_L(G)$
 $- L_{GAN}(G, D) = \mathbb{E}_{x, y} [\log D(x, y)] + \mathbb{E}_{x, z} [1 - \log D(x, G(x, z))]$
 $- L_L(G) = \mathbb{E}_{x, y, z} [\|y - G(x, z)\|_2]$
 CycleGAN: $L(G, F, D_X, D_Y)$
 $= L_{GAN}(G, D_X, x, y) + L_{GAN}(F, D_Y, y, x) + \lambda L_{cyc}(G, F)$
 $L_{cyc} = \mathbb{E}_{x \sim p_{data}(x)} [\|F(G(x)) - x\|_2] + \mathbb{E}_{y \sim p_{data}(y)} [\|G(F(y)) - y\|_2]$



