

# Praktikum: Selbstlernende Systeme

Acrot-Critic

---

08.01.2019

Universität Augsburg  
Institut für Informatik  
Lehrstuhl für Organic Computing

1. Actor-Critic
2. A2C und A3C
3. Implementierung von A2C
4. Quellen

# Actor-Critic

---

## Value based

- Q-learning, Deep Q-Learning
- Value function: mappt jedes Zustands-Aktions-Paar auf einen Wert
- beste Aktion → größter Wert
- funktioniert gut mit einem finiten Set von Aktionen

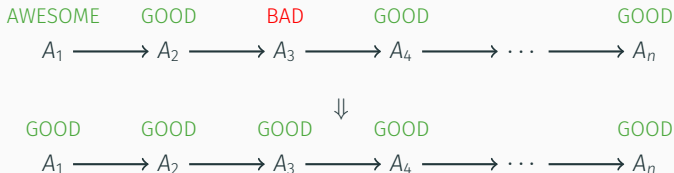
## Policy based

- REINFORCE
- direktes parametrisieren und optimieren der Policy
- keine Value function
- nützlich wenn der Aktionsraum kontinuierlich oder stochastisch ist
- Problem: gute Score function (gesamter EpisodenReward)

beide Methoden haben Nachteile  
⇒ hybride Methode: **Actor-Critic**  
zwei neuronale Netze:

- **Actor:** kontrolliert wie sich der Agent verhält (policy-based)
- **Critic:** misst wie gut die gewählte Aktion ist (value-based)

- Berechnung des Rewards erfolgt am Ende der Episode
- man könnte Annehmen: großes  $R(t) \rightarrow$  alle gewählten Aktionen waren gut



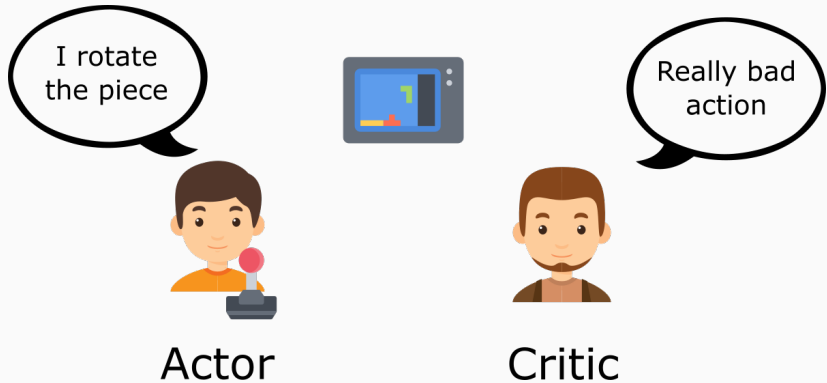
- Alle Aktionen werden durchschnittlich als *GOOD* gewertet
- gesamter Reward entscheidend
- $\rightarrow$  viel Samples nötig  $\rightarrow$  langsames Lernen

- Update in jedem Zeitschritt (TD Learning)

$$\Delta\theta = \alpha * \nabla_{\theta} * (\log \pi(S_t, A_t, \theta)) * R(t)$$

$$\Delta\theta = \alpha * \nabla_{\theta} * (\log \pi(S_t, A_t, \theta)) * Q(S_t, A_t)$$

- gesamter Reward  $R(t)$  kann nicht benutzt werden
- Critic Modell  $\rightarrow$  approximiert die Value function
- Value function ersetzt den Reward





- zu Beginn kein Wissen wie man spielt → zufällig Aktionen ausprobieren
- Critic überwacht und gibt Feedback → Policy updaten
- Critic lernt besseres Feedback zu geben → Critic updaten
- Actor:

$$\pi(s, a, \theta)$$

- Critic:

$$\hat{Q}(s, a, w)$$

- $\Rightarrow$  2 Gewichts Sets  $\theta$  und  $w$

$\theta$  und  $w$  müssen separat upgedatet werden:

## Policy Update:

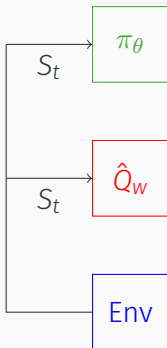
$$\Delta\theta = \alpha \nabla_{\theta} (\log \pi_{\theta}(a|s)) \underbrace{\hat{Q}_w(s, a)}_{\substack{\text{Q learning} \\ \text{function} \\ \text{approximation}}}$$

## Value Update:

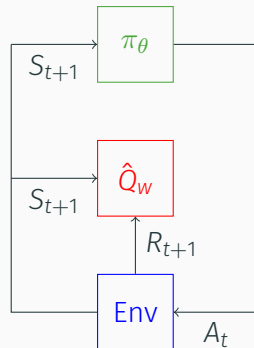
$$\Delta w = \beta \underbrace{(R(s_t, a_t) + \gamma \hat{Q}_w(s_{t+1}, a_{t+1}) - \hat{Q}_w(s_t, a_t))}_{\text{TD error}} \underbrace{\nabla_w \hat{Q}_w(s_t, a_t)}_{\text{Gradient of value function}}$$

$\alpha$  = Policy Learningrate

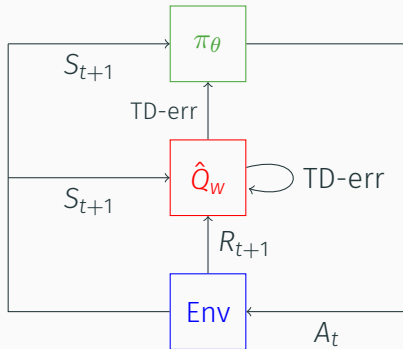
$\beta$  = Value Learningrate



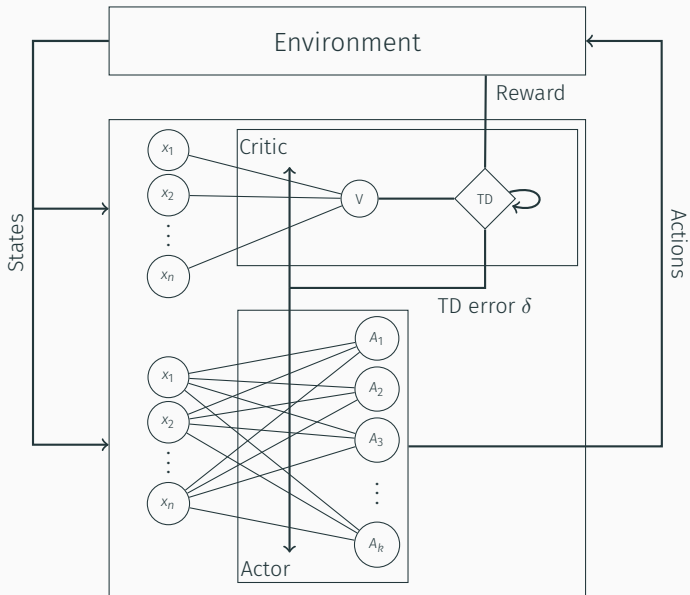
- in jedem Zeitschritt  $t$
- gebe Zustand  $S_t$  vom Environment zu Actor und Critic



- Actor gibt eine Aktion  $A_t$  aus
- Actor und Critic bekommen neuen Zustand  $S_{t+1}$
- Critic bekommt Reward  $R_{t+1}$



- Critic berechnet den Wert den diese Aktion erbringt
- Actor und Critic verbessern ihre Hypothese



## A2C und A3C

---

REINFORCE-Update:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) v_t$$

- Return  $v_t = Q_{\pi_{\theta}}(s_t, a_t)$
- $v_t$  = zukünftiger diskontierter Reward beginnend bei  $t$
- jeder sofortige Reward = Zufallsvariable (Env und  $\pi$  stochastisch)
- $\Rightarrow v_t$  hängt von allen zukünftigen Zuständen ab  $\rightarrow$  hohe Varianz (Schwankung)

- Wenn man eine **Baseline**  $B(s)$  von  $v_t$  abzieht:
  - Kann man die Varianz verringern
  - ohne dabei das Ergebnis zu verändern
  - Eine gute Baseline ist die Value function  $B(s) = V_{\pi\theta}(s)$
- $V_{\pi\theta}(s)$  ist der erwartete Reward beginnend in Zustands  $s$  zum Zeitpunkt  $t$
- $v_t - V_{\pi\theta}(s)$  entfernt die Abhängigkeit von allen Zukünftigen Zuständen von  $v_t$
- Das Ergebnis hängt nur von Zustand  $s$  ab  $\rightarrow$  kleinere Varianz



$$A(s, a) = \underbrace{Q(s, a)}_{Q\text{-Value}} - \underbrace{V(s)}_{\text{avg Value}}$$

- Verbesserung im Vergleich zum durchschnittlichen Wert des Zustands wenn man Aktion  $a$  wählt
- extra Reward wenn man  $a$  wählt
- $A(s, a) > 0 \rightarrow$  Gradient in diese Richtung verschieben
- $A(s, a) < 0 \rightarrow$  Gradient andere Richtung verschieben

**Problem:** zwei Value functions -  $Q(s, a)$  und  $V(s)$

**Lösung:** TD error

$$A(s, a) = Q(s, a) - V(s)$$

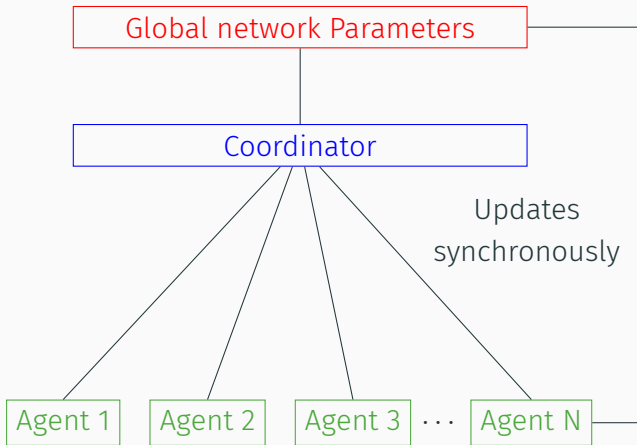
$$= \underbrace{r + \gamma V(s') - V(s)}_{\text{TD Error}}$$

- **A3C - Asynchronous Advantage Actor Critic:**
  - kein Experience Replay
  - unterschiedliche Agenten laufen parallel
  - mehrere Instanzen des Environments
  - Jeder Worker updatet das globale Netz asynchron
- **A2C - Advantage Actor Critic:**
  - genauso wie A3C
  - Außer: synchrones Update des globalen Netzwerks
  - Warten bis alle Worker fertig sind mit dem Training
  - Gradient-Durchschnitt verwenden

- Jeder Agent kommuniziert unabhängig mit den globalen Variablen
- es kann passieren, dass die Agenten mit unterschiedlichen Versionen der Policy arbeiten
- $\Rightarrow$  nicht ganz optimales aggregiertes Update



A2C verwendet einen Koordinator der abwartet bis alle Agenten eine Episode beendet haben



⇒ Training zusammenhängender und schneller

# Implementierung von A2C

---

- normalerweise: 1-step return -  $Q(s, a), V(s)$

$$V(s_0) \leftarrow r_0 + \gamma V(s_1)$$

- mann kann auch mehr steps verwenden

$$V(s_0) \leftarrow r_0 + \gamma r_1 + \gamma^2 V(s_2)$$

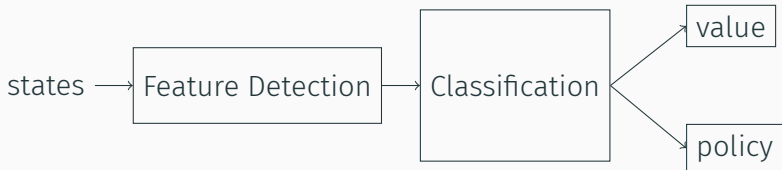
- Allgemein: N-step return

$$V(s_0) \leftarrow r_0 + \gamma r_1 + \cdots + \gamma^n V(s_n)$$

- Änderungen werden schneller propagiert
  - unerwarteter Reward tritt auf
  - 1-step return: Value function ändert sich langsam (1 Schritt rückwärts/Iteration)
  - n-step return: Veränderung wird pro Iteration n Schritte zurückpropagiert
- höhere Varianz → Value hängt von einer Kette von Aktionen ab, welche in viele unterschiedliche Zustände führen können
- Verzögert



- Feature Detection und Klassifikation → Verständnis vom Problem
- Bei zwei Netzen muss das Problem unabhängig voneinander 2x gelernt werden ⇒ **doppelter Trainingsaufwand**
- Lösung: ein Netz mit Verzweigung am Ende (2 Outputs)



$$L = L_{\pi} + c_{val}L_{val} + c_H H$$

$L_{\pi}$ : Policy Loss

$c_{val}$ : Konstante (0.5)

$L_{val}$ : Value Loss

$c_H$ : Konstante (0.01)

$H$ : Policy Entropy

$$L_{\pi} = -J(\pi)$$

$$J(\pi) = E[A(s, a) \nabla \log \pi(a|s)]$$

Durchschnitt aller Samples eines Batches:

$$L_{\pi} = -\frac{1}{n} \sum_{i=1}^n A(s_i, a_i) * \log \pi(a_i, s_i)$$

Achtung:  $A(s, a)$  ist hier eine **Konstante**: `tf.stop_gradient()`

Analog zu Q-Learning

n-step scenario:

$$V(s_0) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^{n-1} r_{n-1} + \gamma^n V(s_n)$$

Fehler/Advantage:

$$err = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^{n-1} r_{n-1} + \gamma^n V(s_n) - V(s_0)$$

Mean squarred error:

$$L_{val} = \frac{1}{n} \sum_{i=1}^n err_i^2$$

- Entropy zur Loss function hinzuzufügen verbessert die Exploration und verhindert vorzeitiges konvergieren zu einer suboptimalen Policy

$$H(\pi(s)) = - \sum_{k=1}^n \pi(s)_k * \log \pi(s)_k$$

$\pi(s)_k$ : Wahrscheinlichkeit für die  $k$ -te Aktion in zustand  $s$

- $\min H = [1, 0, 0, 0]$
- $\max H = [0.25, 0.25, 0.25, 0.25]$
- $\Rightarrow$  durch das maximieren ( $-H$ ) der Entropy bewegen wir die Policy weg vom deterministischen Zustand
- Durchschnitt über alle Samples

$$L_{reg} = - \frac{1}{n} \sum_{i=1}^n H(\pi(s))$$

- managed die einzelnen Instanzen des Environments
- jedes Environment wird von einem Worker gesteuert (eigener Thread)
- Kommunikation:

```
1 from multiprocessing import Process, Pipe
```

- Aktionen der Worker:
  1. reset
  2. step
  3. close

# Quellen

---

- <https://medium.freecodecamp.org/an-introduction-to-policy-gradients-with-cartpole-and-doom-495b5ef2207f>