# **Reinforcement Learning for Language Model Training**

Polina Tsvilodub

**Reinforcement Learning** 



### LMs: Recap Transformers, Training & Inference

- $LM : X \mapsto \Delta(S)$
- transformers use self-attention to offer and retrieve relevant information
  - stacked transformer blocks and multi-head attention increase capacity
- LMs are trained to predict the next word using cross-entropy loss (via teacherforcing)
- decoding schemes are used for inference given a trained LM
  - different stochastic sampling regimes
- SOTA models exhibit 'in-context learning'
- advanced prompting techniques might improve LLMs' generalization performance



### Making LLMs useful Adaptation

- training a task-specific head on top of a model
  - e.g., span prediction layer on top of BERT with frozen BERT
  - on a dataset of ground truth input-output pairs for a particular task
- fine-tuning the model
  - further training part or entire model for a shorter time
  - on a dataset of ground truth input-output pairs for a particular task
- practical problem
  - training with standard supervision is impractical (data collection)
  - and inefficient (restricting "ground truth" to finite set of answers for open-ended tasks)
- direct demonstration of correct behaviour

RL is the solution: learn to achieve goal based on feedback from environment rather than





## **Reinforcement learning**

## Flavors of machine learning

- Unsupervised learning
  - e.g., clustering
- discover patterns in unlabeled data
  - 'given my inductive bias, what is the likely structure of the data?'
- Supervised learning also self-supervised learning aka behavioural cloning
- learn to output Y, given X, from labeled data
  - 'do as I show you'
- learning from demonstration

- Reinforcement learning
  - trial-and-error learning
- learning from interaction / experience
  - 'how do I optimally behave in order to maximize reward?'
  - or, 'how do i optimally achieve my goal?'
    - most natural way of learning?
    - tightly connected to the way organisms behave ("pleasure maximizers")

Christian (2020), Bishop (2009)



### **Reinforcement Learning: Overview** Introduction

- Reinforcement Learning: Computational formalisation of goal-directed learning and decision making under uncertainty
- Goal: Maximize rewards (by learning optimal behavior)
- Basic building blocks:
  - Agent
  - States
  - Actions
  - Transition function P
  - Reward
  - Policy



### Associative RL



Sutton & Barto (2018, p. 48, Fig 3.1), image source





### **Reinforcement Learning: Overview** Introduction

- Reinforcement Learning: Computational formalisation of goal-directed learning and decision making under uncertainty
- Goal: Maximize rewards (by learning optimal behavior)
- Basic building blocks:
  - Agent
  - States:  $S_t \in S$  for t = 0, 1, 2, 3, ...
  - Actions:  $A_t \in A(s)$
  - Transition function:  $P(s' \mid s, a)$
  - Reward:  $R_{t+1} \in R$
  - Policy:  $\pi(S_t) = P(A_t | S_t)$



### Associative RL



Sutton & Barto (2018, p. 48, Fig 3.1), image source





Finite MDPs: 
$$(S, A, T, R)$$
  
1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$   
2.  $A_t \in A(s)$   
3.  $R_{t+1} \in R$   
4.  $T(s' | s, a) = \sum_{r' \in R} P(s', r' | s, a)$   
5. Horizon H, discount factor  $0 \le \gamma \le 1$ 

Expected rewards for state s and action a: r(a)

Markov property:  $S_{t-1}$ ,  $A_{t-1}$  include all information about past agent-environment interactions that are relevant for  $S_t, R_t$ 



$$(s, a) = \mathbb{E}(R_t | S_{t-1} = s, A_{t-1} = a) = \sum_{r \in R} r \sum_{s' \in S} P(s', r | s, a)$$



Finite MDPs: 
$$(S, A, T, R)$$
 Goa

 1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$ 
 $G_t =$ 

 2.  $A_t \in A(s)$ 
 Form

 3.  $R_{t+1} \in R$ 
 epis

 4.  $T(s' | s, a) = \sum_{r' \in R} P(s', r' | s, a)$ 
 $G_t =$ 



al: maximize returns until goal achieved  $= R_{t+1} + R_{t_2} + \ldots + R_T$ 

mally: maximize expected discounted rewards over sode





Finite MDPs: 
$$(S, A, T, R)$$
  
1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$   
2.  $A_t \in A(s)$   
3.  $R_{t+1} \in R$   
4.  $T(s'|s, a) = \sum_{r' \in R} P(s', r'|s, a)$   
Goal: maximize discounted returns  
 $G_t = R_{t+1} + \gamma R_{t_2} + \gamma^2 R_{t+3} + ... + \gamma^{T-t-1} R_T = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$   
 $= R_{t+1} + \gamma G_{t+1}$ 



- and/or actions are: Optimal state-value function:  $V_{\pi}^*(s) = \max \mathbb{E}[G_t | S_t =$  $= \max_{a} \sum_{r}^{n} P(s', r \mid s, a)$

We can identify optimal way to behave if we know what good particular states

$$= s] = \max_{\pi} \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s]$$
$$)[r + \gamma G_{t+1} | S_t = s] \text{ for all } s$$

• Optimization problem: (computationally) find optimal policy  $\pi^*(S_t) = P(A_t | S_t)$ 





Finite MDPs: 
$$(S, A, T, R)$$
  
1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$   
2.  $A_t \in A(s)$   
3.  $R_{t+1} \in R$   
4.  $T(s' | s, a) = \sum_{r' \in R} P(s', r' | s, a)$   
deterministic optimal policy deterministic optimal policy



$$V^{*}(4,3)$$
  
 $V^{*}(3,3)$   
 $V^{*}(2,3)$   
 $V^{*}(1,1)$   
 $V^{*}(4,2)$ 

 $\gamma = 1$  $\gamma = 0.9$ 

 $V^{*}(4,3) = 1$ = 1  $V^{*}(3,3) = 0.9$ = 1 $V^{*}(2,3) = 0.81$ = 1 $V^{*}(1,1) = 0.9^5 = 0.59$ = 1 = -1 $V^{*}(4,2) = -1$ 



Finite MDPs: 
$$(S, A, T, R)$$
  
1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$   
2.  $A_t \in A(s)$   
3.  $R_{t+1} \in R$   
4.  $T(s' \mid s, a) = \sum_{r' \in R} P(s', r' \mid s, a)$   
Goal: maximize discounted returns  
 $G_t = R_{t+1} + \gamma R_{t_2} + \gamma^2 R_{t+3} + ... + \gamma^{T-t-1} R_T = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$   
 $= R_{t+1} + \gamma G_{t+1}$ 



and/or actions are:

Optimal action-value function:  $Q_{\pi}^{*}(s,a) = \max_{\pi} \mathbb{E}[G_{t} | S_{t} = s, A_{t} = a] = \max_{\pi} \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_{t} = s, A_{t} = a]$  $= \sum_{r} P(s',r | s,a)[r + \gamma \max_{a'} Q^{*}(s',a') | S_{t} = s, A_{t} = a] \text{ for all } s, a$ S', r

We can identify optimal way to behave if we know what good particular states





Finite MDPs: 
$$(S, A, T, R)$$
Goal: maximize discounted returns1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$  $G_t = R_{t+1} + \gamma R_{t_2} + \gamma^2 R_{t+3} + ... + \gamma^{T-t-1} R_T = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$ 2.  $A_t \in A(s)$  $= R_{t+1} + \gamma G_{t+1}$ 

**CALC** MDPs: 
$$(S, A, T, R)$$
  
1.  $S_t \in S$  for  $t = 0, 1, 2, 3, ...$   
2.  $A_t \in A(s)$   
3.  $R_{t+1} \in R$   
Goal: maximize discounted returns  
 $G_t = R_{t+1} + \gamma R_{t_2} + \gamma^2 R_{t+3} + ... + \gamma^{T-t-1} R_T = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$   
 $= R_{t+1} + \gamma G_{t+1}$ 

3. 
$$R_{t+1} \in R$$
  
4.  $T(q' \mid q, q) - \sum D(q', q' \mid q, q)$ 

4. 
$$T(s' | s, a) = \sum_{r' \in R} P(s', r' | s, a)$$



- Can be estimated from experience!

• Optimization problem: (computationally) find optimal policy  $\pi^*(S_t) = P(A_t | S_t)$ • Optimal policy  $\pi^* : \pi \ge \pi' \Leftrightarrow v^*_{\pi^*}(s) \ge v_{\pi'}(s)$  for all *s* 

• e.g., value iteration methods for episode tasks





### **Multi-Armed Bandits** Finding the best restaurant

### Finite MDPs: (S, A, T, R)1. $S_t \in S$ where $S = \{s\}$ $Q_{\pi}^{*}$ *t* is iteration over repeated choice now 2. $A(s) = \{a_0, a_1, \dots\}$ $Q_{\pi}^{*}$ 3. $R_t \in R$ 4. $T(s'|s,a) = \sum P(s',r'|s,a)$ $r' \in R$ $A_t = a_i$



Goal: maximize reward over repeated action selection Optimal action-value function:

Optimization problem: (computationally) find optimal policy  $\pi^*(S_t) = P(A_t \mid S_t)$ 

### Can be estimated from experience over repeated decision-making!

• e.g., with action-value methods like the sample-average method central problem: exploration vs. exploitation



# RL Algorithms Approximating Optimal Policy







### **RL Algorithms** Approximating Optimal Policy







### **Policy-Gradient Methods** Introduction

- so far: deriving optimal policy from estimated value function
  - coming up with value functions might be difficult
  - state-value function doesn't prescribe actions
  - action-value functions require argmax
- idea: optimize policy directly, such that expected reward is maximized
  - think: optimize model with respect to objective function
- goal: find optimal  $\theta$ 
  - $\max_{\theta} \mathbb{E}_{\pi_{\theta}}[G_t]$
- $\blacktriangleright$  recall LM optimization: tweak  $\theta$  so as to minimize loss
  - Gradient descent:  $\theta_{new} = \theta_{old} \alpha \nabla L_{\theta}$
  - Now: gradient ascent:  $\theta_{new} = \theta_{old} + \alpha \nabla L_{\theta}$





### **Policy-Gradient Methods** Policy-gradient theorem

- goal: find optimal  $\theta$ 
  - Now: gradient ascent:  $\theta_{new} = \theta_{old} + \alpha \nabla L_{\theta}$
- we write  $\tau$  for a sequence of states, actions, rewards and  $R(\tau)$  for (discounted) return  $L(\theta) = \sum P(\tau, \theta) R(\tau)$
- sample-based policy gradient estimation  $\nabla L(\theta) = \nabla \sum P(\tau, \theta) R(\tau) = \sum \nabla_{\theta} P(\tau, \theta) R(\tau)$  $= \sum_{\sigma} \frac{P(\tau, \theta)}{P(\tau, \theta)} \nabla_{\theta} P(\tau, \theta) R(\tau)$  $= \sum P(\tau,\theta) \frac{\nabla_{\theta} P(\tau,\theta)}{P(\tau,\theta)} R(\tau) = \sum P(\tau,\theta) \nabla_{\theta} \log P(\tau,\theta) R(\tau)$  $\mathcal{T}$  $\approx \frac{1}{m} \sum_{i=1}^{m} \nabla_{\theta} \log P(\tau^{i}, \theta) R(\tau^{i})$

 $V\log(f(x)) = Vf(x)/f(x)$ 



### **Policy-Gradient Methods** Language models as policies

Policy gradient estimation:  $\nabla L(\theta) = \sum P(\tau, \theta) \nabla_{\theta}$ 

- policy P: language model
- trajectories  $\tau$ : generations from language model
- ▶  $\log P(\tau^i, \theta)$ : log probability of a generation  $\tau^i$  under the language model
- $R(\tau^i)$ : reward for generation  $\tau^i$

$$\log P(\tau, \theta) R(\tau) \approx \frac{1}{m} \sum_{i=1}^{m} \nabla_{\theta} \log P(\tau^{i}, \theta) R(\tau^{i})$$

Sutton & Barto (2018)



### **Summary** Reinforcement learning

- the central framework for formalizing RL problems are Markov Decision Processes (MDPs)
- task of RL is to solve MDP such that the expected return is maximized
  - and to find the optimal policy
- classical solution methods for MDPs include estimation of optimal state- and action-value functions
- policy gradient methods directly optimize the policy such that the expected return is maximized
  - can be applied to LMs!



## Announcements

### Solutions to exercises will be on Moodle on November 15th!

### Next class (November 15th) online only!