

Distributed Deep Q-Learning

Hao Yi Ong, Kevin Chavez, and Augustus Hong

CME 323, Stanford University

June 3, 2015

Outline

Introduction

Mathematical formulation

Serial algorithm

Distributed algorithm

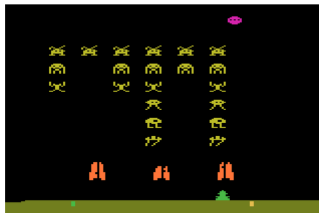
Numerical experiments

Conclusion

Motivation

- ▶ long-standing challenge of reinforcement learning (RL)
 - control with high-dimensional sensory inputs (e.g., vision, speech)
 - shift away from reliance on hand-crafted features
- ▶ utilize breakthroughs in deep learning for RL [M⁺13, M⁺15]
 - extract high-level features from raw sensory data
 - learn better representations than handcrafted features with neural network architectures used in supervised and unsupervised learning
- ▶ create fast learning algorithm
 - train efficiently with stochastic gradient descent (SGD)
 - distribute training process to accelerate learning [DCM⁺12]

Success with Atari games



Goals

distributed deep RL algorithm

- ▶ robust neural network agent
 - must succeed in challenging test problems
- ▶ control policies with high-dimensional sensory input
 - obtain better internal representations than handcrafted features
- ▶ fast training algorithm
 - efficiently produce, use, and process training data

Outline

Introduction

Mathematical formulation

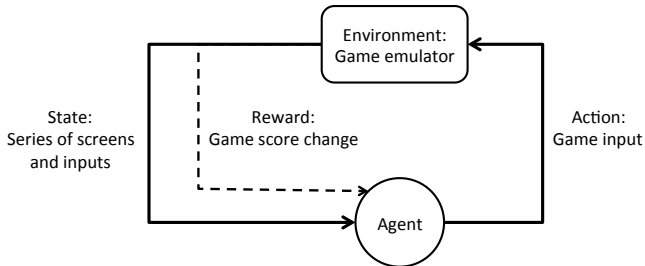
Serial algorithm

Distributed algorithm

Numerical experiments

Conclusion

Playing games



objective: learned policy maximizes future rewards

$$R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'},$$

- ▶ discount factor γ
- ▶ reward change at time t' $r_{t'}$

State-action value function

- ▶ basic idea behind RL is to estimate

$$Q^*(s, a) = \max_{\pi} \mathbf{E}[R_t \mid s_t = s, a_t = a, \pi],$$

where π maps states to actions (or distributions over actions)

- ▶ optimal value function obeys Bellman equation

$$Q^*(s, a) = \mathbf{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right],$$

where \mathcal{E} is the MDP environment

Q-network

- ▶ trained by minimizing a sequence of loss functions

$$L^{(i)}(\theta^{(i)}) = \mathbf{E}_{s,a \sim \rho(\cdot)} \left[\left(y^{(i)} - Q(s, a; \theta^{(i)}) \right)^2 \right],$$

with

- iteration number i , i th network parameters $\theta^{(i)}$
 - target $y^{(i)} = \mathbf{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta^{(i-1)}) \mid s, a \right]$
 - “behavior distribution” (exploration policy) $\rho(s, a)$
- ▶ architecture varies according to application

Outline

Introduction

Mathematical formulation

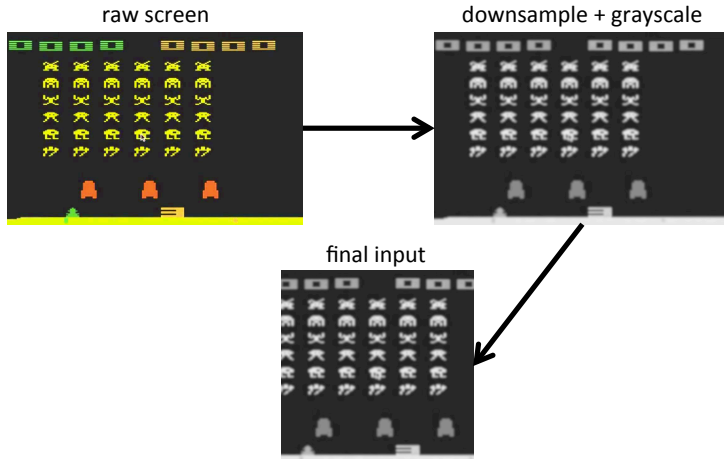
Serial algorithm

Distributed algorithm

Numerical experiments

Conclusion

Preprocessing



Q-learning

- ▶ optimize Q-network loss function via

$$Q(s, a) := Q(s, a) + \alpha \left(r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)$$

- ▶ trains optimal policy using “behavior policy” (off-policy)
 - learns policy $\pi^*(s) = \operatorname{argmax}_a Q(s, a; \theta)$
 - uses an ϵ -greedy strategy (behavior policy) for state-space exploration

Experience replay

a kind of short-term memory

- ▶ store agent's experiences at each time step

$$e_t = (s_t, a_t, r_t, s_{t+1})$$

- ▶ experiences form a replay memory dataset

$$\mathcal{D} = \{e_1, \dots, e_N\},$$

where N is the fixed memory capacity

- ▶ execute Q-learning updates with samples of experience

$$e \sim \mathcal{D}$$

Serial deep Q-learning

given replay memory \mathcal{D} with capacity N

initialize Q-networks Q, \hat{Q} with same random weights θ

repeat until timeout

initialize frame sequence $s_1 = \{x_1\}$ and preprocessed state $\phi_1 = \phi(s_1)$
for $t = 1, \dots, T$

1. select action $a_t = \begin{cases} \max_a Q(\phi(s_t), a; \theta) & \text{w.p. } 1 - \epsilon \\ \text{random action} & \text{otherwise} \end{cases}$
2. execute action a_t and observe reward r_t and frame x_{t+1}
3. append $s_{t+1} = (s_t, a_t, x_{t+1})$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
4. store experience $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}
5. uniformly sample minibatch $(\phi_j, a_j, r_j, \phi_{j+1}) \sim \mathcal{D}$
6. set $y_j = \begin{cases} r_j & \text{if } \phi_{j+1} \text{ terminal} \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta) & \text{otherwise} \end{cases}$
7. perform gradient descent step for Q on minibatch
8. every C steps reset $\hat{Q} = Q$

Outline

Introduction

Mathematical formulation

Serial algorithm

Distributed algorithm

Numerical experiments

Conclusion

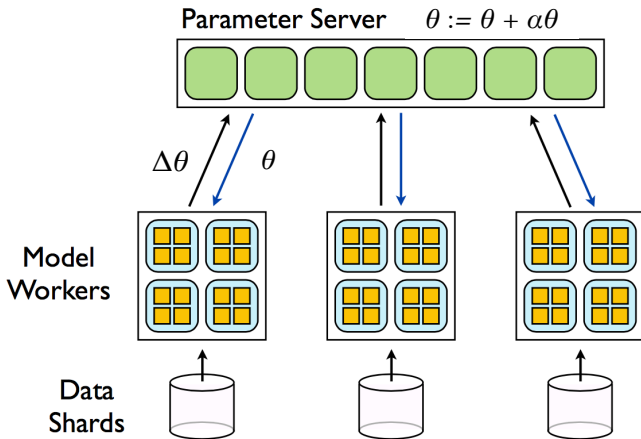
Model parallelism

for each Q-network

- ▶ partition model across CPUs/GPUs
 - up to availability of CPU/GPU resources
 - uses Caffe deep learning framework
- ▶ **[[Kevin/Hao Yi: How does caffe use CPU/GPU resources? How does complexity scale for implementation? Answers question of how our algorithm scale for model.]]**

Data parallelism

downpour SGD: generic asynchronous distributed SGD



Implementation

- ▶ data shards are generated locally on each model worker in real-time
 - data is stored independently for each worker
 - since game emulation is simple, generating data is fast
 - simple fault tolerance approach: regenerate data if worker dies
- ▶ algorithm scales very well with data
 - since data lives locally on workers, no data is sent

Implementation

- ▶ bottleneck is parameter update time on parameter server
 - e.g., if parameter server gradient update takes 2 ms, then we can only do up to 500 updates per second (using buffers, etc.)
- ▶ trade-off between parallel updates and model staleness
 - because worker is likely using a stale model, the updates are “noisy” and not of the same quality as in serial implementation

Implementation

communication pattern

- ▶ one-to-all and all-to-one, but asynchronous for every minibatch
- ▶ like multiple asynchronous all-reduces

Outline

Introduction

Mathematical formulation

Serial algorithm

Distributed algorithm

Numerical experiments

Conclusion

Evaluation



Snake

► parameters

- snake length grows with number of apples eaten
- one apple at any time, regenerated once eaten
- $n \times n$ array, with walled-off world
- want to maximize score, equal to snake length

► complexity

- four possible states for each cell: {empty, head, body, apple}
- state space cardinality is $O(n^8)$
- four possible actions: {north, south, east, west}

Results

Outline

Introduction

Mathematical formulation

Serial algorithm

Distributed algorithm

Numerical experiments

Conclusion

Summary

References

- ▶ Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, et al.
Large scale distributed deep networks.
In Advances in Neural Information Processing Systems, pages 1223–1231, 2012.
- ▶ V. Mnih et al.
Playing Atari with deep reinforcement learning.
arXiv preprint arXiv:1312.5602, 2013.
- ▶ V. Mnih et al.
Human-level control through deep reinforcement learning.
Nature, 518(7540):529–533, 2015.

Theoretical complications

deep learning algorithms require

- ▶ huge training datasets
- ▶ independence between samples
- ▶ fixed underlying data distribution

Deep Q-learning

avoids theoretical complications

- ▶ greater data efficiency
 - each experience potentially used in many weight updates
- ▶ reduce correlations between samples
 - randomizing samples breaks correlations from consecutive samples
- ▶ experience replay averages behavior distribution over states
 - smooths out learning
 - avoids oscillations or divergence in gradient descent