### Model-Free Control

(Reinforcement Learning)

### and Deep Learning

### MARC G. BELLEMARE

Google Brain (Montréal)



Starters				
Onion Bhaji Thinky sliced onion in gram flour and herb, deep fried	£2.45			
Samosa (Meat or Vegetable) Tilangular shoped savoury filed with spicy meat or vegetable	£2.45			
Prawn on Puree Prawn cooked in herts & spices, served on a thin fried bread	£3.25			
King Prawn on Puree Large Frawn cooked in herbs & spices, served on a thin filed bread	£4.25			
King Prawn Butterfly Buttered king prowns with herbs and spices and deep fied	£4.25			
Tandoori Chicken (Qir) Spring chicken matinated in herbs & spices, cooked over charcoal tandoor	£3.75			
Chicken or Lamb Tikka Diced chicken or lamb marinated in yoghutt and spices and then backecued in the clay oven	£3.25			
Sheek Kebab Minced lamb pungently spiced, skewered and barbecued in the tandoor oven	£3.25			
Shami Kebab Minced lamb in herbs and spices then fried in ghee	£3.50			
Chicken Chat or Aloo Chat (sour) Small plicy pieces of chicken or potatoes spiced with hot and sour source	£2.95			
Prawn Cocktail Succutent prawns in our own mayorinate	£2.95			
Dall Soup or Mulligatwany Soup (Lentils)	£2.25			
Stuffed Tomatoes Tomatoes stuffed with vegetables or minced meat	£3.25			
Chicken or Vegetable Pakora Chicken titko or vegetable deep fried in gram four coating	£3.25			
Plain or Massala Papadum	60p			
Mixed Kebab	£4.95			
Kebab Roll	£3.75			

















THE ARCADE LEARNING ENVIRONMENT (BELLEMARE ET AL., 2013)



- 33,600 (discrete) dimensions
- Up to 108,000 decisions/episode (30 minutes)
- 60+ games: heterogenous dynamical systems



#### DEEP LEARNING: AN AI SUCCESS STORY



$$\mathcal{M} := \langle \mathcal{X}, \mathcal{A}, R, P, \gamma \rangle$$
$$Q^*(x, a) = r(x, a) + \gamma \mathop{\mathbf{E}}_{P} \max_{a' \in \mathcal{A}} Q^*(x', a')$$
$$\hat{Q} = \Pi \mathcal{T}^{\pi} \hat{Q}$$
$$\|\hat{Q} - Q^{\pi}\|_D \le \frac{1}{1 - \gamma} \|\Pi Q^{\pi} - Q^{\pi}\|_D$$



Practice

Theory

### WHERE HAS MODEL-FREE CONTROL BEEN SO SUCCESSFUL?

- Complex dynamical systems
  - Black-box simulators
  - High-dimensional state spaces
  - Long time horizons
  - Opponent / adversarial element









- Are simulations reasonably cheap? model-free
- Is the notion of "state" complex? model-free
- Is there partial observability? maybe model-free
- Can the state space be enumerated? value iteration
- Is there an explicit model available? model-based

### OUTLINE OF TALK



### WHAT'S REINFORCEMENT LEARNING, ANYWAY?



"ALL GOALS AND PURPOSES ... CAN BE THOUGHT OF AS THE MAXIMIZATION OF SOME VALUE FUNCTION"

- SUTTON & BARTO (2017, IN PRESS)



- At each step t, the agent
  - Observes a state
  - Takes an action
  - Receives a reward

### THREE LEARNING PROBLEMS

# **Stochastic** approximation **Policy evaluation Optimal** control **Function** approximation

• Formalized as a Markov Decision Process:

$$\mathcal{M} := \langle \mathcal{X}, \mathcal{A}, R, P, \gamma \rangle$$

- R, P reward, transition functions
- γ discount factor
- A trajectory is a sequence of interactions with the environment

 $x_1, a_1, r_1, x_2, a_2, \ldots$ 

- Policy  $\pi$ : a probability distribution over actions:  $a_t \sim \pi(\cdot \mid x_t)$  If deterministic:  $a_t = \pi(x_t)$
- Transition function:  $x_{t+1} \sim P(\cdot | x_t, a_t)$
- Value function  $Q^{\pi}(x, a)$ : total discounted reward

$$Q^{\pi}(x,a) = \mathbf{E}_{P,\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(x_t, a_t) \,|\, x_0, a_0 = x, a \right]$$

• As a vector in space of value functions:  $Q^{\pi} \in \mathcal{Q}$ 

"Maximize value function": find

$$Q^*(x,a) := \max_{\pi} Q^{\pi}(x,a)$$

• Bellman's equation:

$$Q^{\pi}(x,a) = \mathbf{E}_{P,\pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} r(x_{t},a_{t}) \,|\, x_{0}, a_{0} = x, a \right]$$
$$= r(x,a) + \gamma \mathbf{E}_{P,\pi} Q^{\pi}(x',a')$$

• Optimality equation:

$$Q^*(x,a) = r(x,a) + \gamma \operatorname{E}_{P} \max_{a' \in \mathcal{A}} Q^*(x',a')$$



### Bellman Operator

$$\mathcal{T}^{\pi}Q(x,a) := r(x,a) + \gamma \mathop{\mathbf{E}}_{\substack{x' \sim P \\ a' \sim \pi}} Q(x',a')$$

- The Bellman operator is a  $\gamma$ -contraction:  $\|\mathcal{T}^{\pi}Q - Q^{\pi}\|_{\infty} \leq \gamma \|Q - Q^{\pi}\|_{\infty} |_{Q_k}$
- Fixed point:

$$Q^{\pi} = \mathcal{T}^{\pi} Q^{\pi}$$



### BELLMAN OPTIMALITY OPERATOR

$$\mathcal{T}Q(x,a) := r(x,a) + \gamma \mathop{\mathbf{E}}_{x' \sim P} \max_{a' \in \mathcal{A}} Q(x',a')$$

- Also a  $\gamma$ -contraction (beware! different proof):  $\|\mathcal{T}Q - Q^*\|_{\infty} \leq \gamma \|Q - Q^*\|_{\infty} \qquad Q_k$
- Fixed point is optimal v.f.:

$$Q^* = \mathcal{T}Q^* \ge Q^\pi$$



1.Value iteration:

 $Q_{k+1}(x,a) \leftarrow \mathcal{T}Q_k(x,a) = r(x,a) + \gamma \mathbf{E}_P \max_{a' \in \mathcal{A}} Q_k(x',a')$ 

### 2.Policy iteration:

a.  $\pi_k = \underset{\pi}{\arg \max} \mathcal{T}^{\pi}Q_k(x, a)$  b.  $Q_{k+1}(x, a) \leftarrow Q^{\pi_k}(x, a)$ 3. Optimistic policy iteration:

$$Q_{k+1}(x,a) \leftarrow (\mathcal{T}^{\pi_k})^m Q_k(x,a) = \underbrace{\mathcal{T}^{\pi_k} \cdots \mathcal{T}^{\pi_k}}_{m \text{ times}} Q_k(x,a)$$

#### POLICY ITERATION



### MODEL-FREE REINFORCEMENT LEARNING

- Typically no access to P, R
- Two options:
  - Learn a model (not in this talk)
  - Model-free: learn  $Q^{\pi}$  or  $Q^*$  directly from samples



- For all x, a, sample  $x' \sim P(\cdot \,|\, x, a), a' \sim \pi(\cdot \,|\, x')$
- The SARSA algorithm:

$$Q_{t+1}(x,a) \leftarrow (1-\alpha_t)Q_t(x,a) + \alpha_t \hat{\mathcal{T}}_t^{\pi} Q_t(x,a)$$
  
=  $(1-\alpha_t)Q_t(x,a) + \alpha_t \left(r(x,a) + \gamma Q_t(x',a')\right)$   
=  $Q_t(x,a) + \alpha_t \left(r(x,a) + \gamma Q_t(x',a') - Q_t(x,a)\right)$   
TD-error  $\delta$ 

•  $\alpha_t \in [0,1)$  is a step-size (sequence)

• The Q-Learning algorithm: max. at each iteration

$$Q_{t+1}(x,a) \leftarrow (1-\alpha_t)Q_t(x,a) + \alpha_t \big( r(x,a) + \gamma \max_{a' \in \mathcal{A}} Q_t(x',a') - Q_t(x,a) \big)$$

- Both converge under Robbins-Monro conditions
- Not trivial! Interleaved learning problems



### Asynchronous Updates

• The asynchronous case: learn from trajectories

 $x_1, a_1, r_1, x_2, a_2, \dots \sim \pi, P$ 

Apply update at each step:

 $Q(x_t, a_t) \leftarrow Q_t(x, a) + \alpha_t (r_t + \gamma Q(x_{t+1}, a_{t+1}) - Q(x_t, a_t))$ 

- This is the setting we usually deal with
- Convergence even more delicate



OPEN QUESTIONS/ AREAS OF ACTIVE RESEARCH

- Rates of convergence [1]
- Variance reduction [2]
- Convergence guarantees for multi-step methods [3, 4]
- Off-policy learning: control from fixed behaviour [3, 4]

[1] Konda and Tsitsiklis (2004)

- [2] Azar et al., Speedy Q-Learning (2011)
- [3] Harutyunyan, Bellemare, Stepleton, Munos (2016)
- [4] Munos, Stepleton, Harutyunyan, Bellemare (2016)

## Stochastic approximation

### Policy evaluation

### Optimal control

## Function approximation

## Stochastic approximation

### Policy evaluation

### **Optimal** control

# Function approximation

### (VALUE) FUNCTION APPROXIMATION

• Parametrize value function:

 $Q^{\pi}(x,a) \approx Q(x,a,\theta)$ 

- Learning now involves a projection step  $\Pi$ :  $\Pi \mathcal{T}^{\pi}Q(x, a, \theta_k) : \theta_{k+1} \leftarrow \arg\min_{\theta} \left\| \mathcal{T}^{\pi}Q_k(x, a, \theta_k) - Q(x, a, \theta) \right\|_D$
- This leads to additional, compounding error
- Can cause divergence



### Some classic Results [1]

- Linear approximation:  $Q^{\pi}(x,a) \approx \theta^{\top} \phi(x,a)$
- SARSA converges to  $\hat{Q}$  satisfying

$$\hat{Q} = \Pi \mathcal{T}^{\pi} \hat{Q} \quad \|\hat{Q} - Q^{\pi}\|_{D} \leq \frac{1}{1 - \gamma} \|\Pi Q^{\pi} - Q^{\pi}\|_{D}$$

• Q-Learning may diverge!

[1] Tsitsiklis and Van Roy (1997)

OPEN QUESTIONS/ AREAS OF ACTIVE RESEARCH

- Convergent, linear-time optimal control [1]
- Exploration under function approximation [2]
- Convergence of multi-step extensions [3]

[1] Maei et al. (2009)
[2] Bellemare, Srinivasan, Ostrovski, Schaul, Saxton, Munos (2016)
[3] Touati et al. (2017)



















### Deep Learning



#### Deep Learning

 $\nabla_{\theta} \mathcal{L}(\theta)$ 



Graphic by Volodymyr Mnih



Mnih et al., 2015

- Value function as a Q-network  $Q(x,a,\theta)$
- Objective function: mean squared error

$$\mathcal{L}(\theta) := \mathbf{E} \left[ \left( \underbrace{r + \gamma \max_{a' \in \mathcal{A}} Q(x', a', \theta)}_{\text{target}} - Q(x, a, \theta) \right)^2 \right]$$

• Q-Learning gradient:

$$\nabla_{\theta} \mathcal{L}(\theta) = \mathbf{E} \left[ \left( r + \gamma \max_{a' \in \mathcal{A}} Q(x', a', \theta) - Q(x, a, \theta) \right) \nabla_{\theta} Q(x, a, \theta) \right]$$

- Naive Q-Learning oscillates or diverges
- 1. Data is sequential
  - + Successive samples are non-iid
- 2. Policy changes rapidly with Q-values
  - + May oscillate; extreme data distributions
- 3. Scale of rewards and Q-values is unknown
  - + Naive gradients can be large; unstable backpropagation

- 1. Use experience replay
  - + Break correlations, learn from past policies
- 2. Target network to keep target values fixed
  - Avoid oscillations
- 3. Clip rewards
  - + Provide robust gradients

Equivalent to planning with empirical model

- Build dataset from agent's experience
  - Take action according to ε-greedy policy
  - Store (x, a, r, x', a') in replay memory D
  - Sample transitions from D, perform asynchronous update:

 $\mathcal{L}(\theta) = \mathbf{E}_{\substack{x,a,r,x',a'\sim\mathcal{D}}} \left[ \left( r + \gamma \max_{a'\in\mathcal{A}} Q(x',a',\theta) - Q(x,a,\theta) \right)^2 \right]$ 

Effectively avoids correlations within trajectories



- To avoid oscillations, fix parameters of target in loss function
  - Compute targets w.r.t. old parameters

$$r + \gamma \max_{a' \in \mathcal{A}} Q(x', a', \theta^-)$$

• As before, minimize squared loss:

$$\mathcal{L}(\theta) = \mathbf{E}_{\mathcal{D}} \left[ \left( r + \gamma \max_{a' \in \mathcal{A}} Q(x', a', \theta^{-}) - Q(x, a, \theta) \right)^2 \right]$$

• Periodically update target network:

$$\theta^- \leftarrow \theta$$



Based on material by David Silver

- Clip rewards in range [-1, +1]
- Ensures gradients are well-conditioned
- Also prevents value overestimation
- No longer can tell small, large rewards apart

	Q-learning	Q-learning	Q-learning	Q-learning
			+ Replay	+ Replay
		+ Target Q		+ Target Q
Breakout	3	10	241	317
Enduro	29	142	831	1006
River Raid	1453	2868	4103	7447
Seaquest	276	1003	823	2894
Space Invaders	302	373	826	1089



### Some Recent Research

### ACTIVE RESEARCH: OFF-POLICY METHODS

 Reusing data (e.g. from experience replay) can diverge with approximation:

$$Q(x, a, \theta) \stackrel{\nabla_{\theta}}{\leftarrow} r(x, a) + \gamma \mathop{\mathbf{E}}_{x' \sim P_a} \max_{a' \in \mathcal{A}} Q(x', a', \theta)$$

- Can use importance sampling ratio:  $\frac{\pi(a \mid s)}{\mu(a \mid s)}$
- But variance is high
- Also safety issues: how to guarantee performance?

Precup, Sutton, and Singh (2000) Thomas and Brunskill (2016)

### ACTIVE RESEARCH: MULTI-STEP METHODS

• Greater accuracy [1] from multi-step returns:

$$\mathcal{T}^{\lambda}Q(x,a) := \sum_{k=0}^{\infty} \lambda^{k} \left[ \underbrace{\sum_{t=0}^{k} \gamma^{t} r(x_{t},a_{t}) + \gamma^{k+1} Q(x_{k+1},a_{k+1})}_{\text{n-step return}} \right]$$

 $= Q(x,a) + \sum_{t=0}^{\infty} (\lambda \gamma)^t \delta(x_t, a_t)$ 

- Retrace( $\lambda$ ) [2] both off-policy and multi-step

$$\mathcal{R}Q(x,a) := Q(x,a) + \sum_{t=0}^{\infty} (\lambda\gamma)^t \left(\prod_{s=0}^{t-1} c_s\right) \delta(x_t, a_t) \quad c_s := \min\left\{1, \frac{\pi(a_s \mid x_s)}{\mu(a_s \mid x_s)}\right\}$$

- Convergence surprisingly nontrivial, even without value approximation
  - Tsitsiklis and Van Roy (1997)
     Munos, Stepleton, Harutyunyan, Bellemare (2016)

### ACTIVE RESEARCH: GAP-INCREASING OPERATORS

- Action gap:  $\max_{a' \in \mathcal{A}} Q^*(x, a') Q^*(x, a)$
- New operators that increase action gap, e.g.

 $\tilde{\mathcal{T}}Q(x,a) := \mathcal{T}Q(x,a) - \beta \left[ \max_{a' \in \mathcal{A}} Q(x,a') - Q(x,a) \right], \quad \beta \in [0,1)$ 

- Not necessarily contraction operators
- Suboptimal Q-values may not converge
- Yet: guaranteed convergence:

$$\lim_{k \to \infty} \max_{a \in \mathcal{A}} (\tilde{\mathcal{T}})^k Q(x, a) = \max_{a \in \mathcal{A}} Q^*(x, a)$$

Bellemare, Ostrovski, Guez, Thomas, Munos (2016)



# Stochastic approximation

### Policy evaluation

### **Optimal control**



### Function approximation

### Model-Free Control with Deep Learning MARC G. BELLEMARE



G. Ostrovski



Arthur Guez



D. Saxton



Rémi Munos



T. Stepleton



S. Srinivasan



T. Schaul



J. Veness



Y. Naddaf



M. Bowling



Philip Thomas



A. Harutyunyan