Project Ideas
HW 4    Friday

```
y |
  |_____
      x
```

```
xs = [ ]
ys = [ ]
x = 0
for i in 1:10000
      reset!.(env)                          ┐
      while !terminated(env)                │  learning
            s = observe(env)                │  episode
            a = explore(s)    ←             │
            r = act!(env, a)                │
            x += 1                          │
            update Q  ←                     ┘
      if i % 100 == 0
            push!(xs, x)
            push!(ys, eval(env, Q))
plot(xs, ys)

function eval(env, Q)
      rsum = 0
      for i in 1:1000                       ┐
            reset!(env)                     │
            while !terminated(env)          │  evaluation
                  s = observe(env)          │  episode
                  a = argmax(Q(s,a))        │
                  rsum += act!(env,a)       │
      return rsum/1000                      ┘
```

## Last Time
### Neural Networks

$$f_\theta(x) = W_1 \sigma\left(W_2 \sigma\left(W_3 x + b_3\right) + b_2\right) + b_1$$

$$\theta = (W_1, b_1, W_2, b_2, W_3, b_3)$$

Stochastic Gradient Descent with backprop

## This Time
### RL with Neural Networks

# DQN    DPG

Approximate $Q(s,a)$ with $Q_\theta(s,a)$

### Review: Q-learning update

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

### Deep Q learning

x,y

```
loop
    s = observe(env)
    a = explore(s)
    r = act!(env, s)
    s' = observe(env)
```
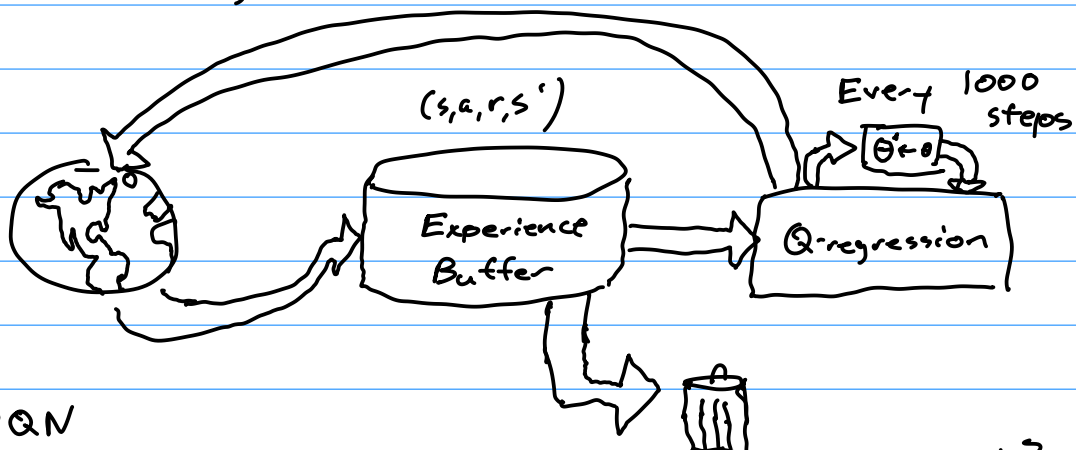$$y = r + \gamma \max_{a'} Q(s', a')$$
$$\theta \leftarrow \theta + \alpha \nabla_\theta Q(s,a)(y - Q(s,a))$$

This won't work ... at all.

$(s,a,r,s')$

1. Samples highly correlated } use experience buffer
2. Size-1 batches
3. Moving Target } periodically freeze $\theta$

Every 1000 steps

$(s,a,r,s')$

$\theta' \leftarrow \theta$

Experience Buffer

Q-regression

"Classic" DQN

$$\ell(s,a,r,s') = \left(r + \gamma \max_{a'} Q_{\theta'}(s',a') - Q_\theta(s,a)\right)^2$$

## Rainbow

$$r \leftarrow \gamma Q_1(s', \underset{a'}{\arg\max} \, Q_2(s';a'))$$

- Double Q-Learning
- Prioritized replay
  └ priority proportional to last TD error
- Dueling networks
  - Value network + advantage network
  · $Q(s,a) = V(s) + A(s,a)$

$(s, a, r, s')$

- Multi-step learning
  $$\left(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n \underset{a'}{\max} Q_\theta \cdot (s_{t+n}, a') - Q(s,a)\right)^2$$
- Distributional RL
  predict a distribution of returns
- Noisy nets
  - add noise to neural network

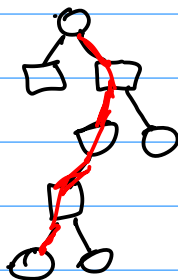## Breakout Room

What is the difference between MCTS, SARSA-$\lambda$?

$Q(s,a)$   $N(s,a)$

| MCTS | SARSA-$\lambda$ |
|---|---|
| Rollouts | Q-values |
| UCB1 | $\varepsilon$-greedy |
| └ N used for exploration | N values used for eligibility |
| Tree | No tree |
| Model-Based | Model-Free |
| Online - decision at current state | Offline - policy |

log trick
causality
baselines

DPG

$$\nabla_\theta U(\theta) = E_\tau \left[ \sum_{k=1}^{d} \nabla_\theta \log \pi_\theta(a^{(k)}|s^{(k)}) \, \gamma^{(k-1)} \left( r_{to\text{-}go}^{(k)} - r_{base}(s^{(k)}) \right) \right]$$

$$U(\theta') \approx U(\theta) + \nabla U(\theta)^T (\theta' - \theta)$$

$$g(\theta,\theta') = \tfrac{1}{2}(\theta'-\theta)^T I (\theta'-\theta) = \tfrac{1}{2}\|\theta'-\theta\|_2^2$$

maximize $U(\theta) + \nabla U(\theta)^T (\theta'-\theta)$
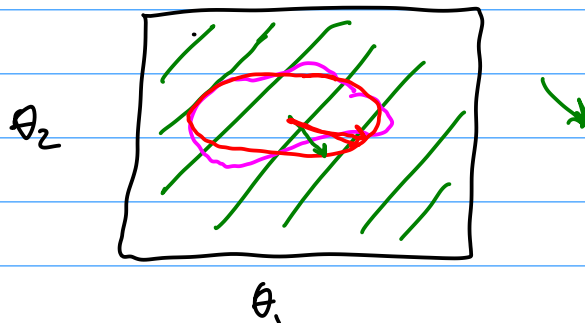$\theta'$

subject to $g(\theta,\theta') \le e$

analytic solution

$$\theta' = \theta + u\sqrt{\frac{2e}{u^T u}} = \theta + \sqrt{2e}\,\frac{u}{\|u\|}$$

$$u = \nabla U(\theta)$$

## Natural Gradient

$$g(\theta,\theta') = D_{KL}(p(\cdot|\theta) \| p(\cdot|\theta')) \le e$$



$\theta_2$

$\theta_1$

$$g(\theta,\theta') = \tfrac{1}{2}(\theta'-\theta)^T F_\theta (\theta'-\theta) \le \varepsilon$$

↰ Taylor approximation

$$F_\theta = E_\tau \left[ \nabla \log p(\tau|\theta) \nabla \log p(\tau|\theta)^T \right]$$

maximize $\nabla U(\theta)^T (\theta'-\theta)$
$\theta'$

subject to $\tfrac{1}{2}(\theta'-\theta)^T F_\theta (\theta'-\theta) = e$

analytical

$$\theta' = \theta + u\sqrt{\frac{2e}{\nabla U(\theta)^T u}}$$

TRPO
PPO

$$u = F_\theta^{-1} \nabla U(\theta)$$