# ✳ Introduction.

기존: 수동으로 설계했던 feature, linear한 value function이 policy 를 이용한 RL이 주로 사용됐다.
    그래서 이런 feature representation의 quality에 성능이 매우 의존적이었다.

현재의 RL challenge; High dimensional input인 vision과 speech과같은 agent가 받게되는 control해야 함.
  → 최근 딥러닝의 발전으로 CV 나 speech recognition 이 큰 발전이 일어났고.
    raw sensory data로부터 high-level feature를 추출할 수 있게 되었다.
    그래서 이런 강화학습에도 적용해보려고 한다.

  → 그러나 딥러닝 강화학습에 적용하는데 많은 어려움이 있었는데,
    i) 기존 DL은 대량의 hand-labeled data가 필요한 반면,
       RL은 sparse, noisy, delay 된 scalar reward에서 학습한다.
       특히 action과 reward 간의 delay가 큰 경우가 있다.
    ii) DL 신경망은 data sample간에 서로 independent하다고 가정하는데,
       RL은 state간 연관되어 있다.
    iii) RL 의 행동은 reward 최대화하는 방향 data distribution이 변화한다. (데이터-약점)

 ⇒ CNN에 위와 유사하게. Q-learning 변기법으로 학습하고, SGD로 가중치 update한다.
   ii), iii) 같은 위험성을 극복하기 위해 replay 메모리 사용.

⇒ Testbed로 Atari 게임으로 사용했다.
   Atari 게임은 high-dimensional visual input을 받고, 인간에게도 어려운 다양한 task를 풀게한다.
   목표는 single NN agent로 만드는것.
   이 agent는 사람이 사용하는 동일한 input 받고, visual feature 와와 보상 설계없이 시작한다.
   ex, video input, reward, terminal signal, set of possible actions만 가진다. hyperparameter 설계까지 동일하게 유지한다...

# ✳ Background

$$a \in \{1, ... k\}$$
Each time step. agent는 action set 중에서 action을 고른다.
• Active emulator는 조작된 internal states game score(reward) 를 변경한다.
• ε (environment) 는 stochastic하다. Internal state는 agent로부터 오는데 Agent는 그저 screen의 vector과 raw pixel를 받는다.
• MDP 상황을 가정.
   + action과 observation의 sequence-정보가 가진. $S_t = x_1, a_1, x_2, ...$
   + 각각 별의 time-step 에게 terminal state에서 끝난다고 가정

• 기타 Assumption
   – Discounting rate 포함.
   – Optimal action value $Q^*(s, a)$ 고려.
     → state s'에서 모든 가능한 action 중에 optimal한 action a'을 다시 고르면, 현재 $Q^*(s', a')$ 를 고른다.

$$Q^*(s, a) = \mathbb{E}_{s' \sim \varepsilon}\left[r + \gamma \max_{a'} Q^*(s', a') \,\middle|\, s, a\right]$$

   – 기존의 value iteration 에서는, $Q_{i+1}(s, a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s', a') \,|\, s, a\right]$ 반복적으로 $i \to \infty$ 하면 $Q_i \to Q^*$ 라고 수렴한다고 알려짐.
     그러나 action value function이 각 sequence에 대해 분리되어 추정이 어렵다.
     따라서 $Q(s, a; \theta) \approx Q^*(s, a)$ 인 approximator를 대신에서 만든 사용한다.
     function approximator는 weight가 $\theta$인 Neural network function을 Q-network라고 사용한다.
     Q-network는 loss function $L_i(\theta_i)$ 를 반복적으로 변화를 최소화한다.

$$L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)}\left[\left(y_i - Q(s, a; \theta_i)\right)^2\right]$$
$$y_i = \mathbb{E}\left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) \,|\, s, a\right]$$
$\rho(s, a)$ : behavior distribution

   – Supervised learning과 달리, target과 network weight에 의해 변화된다.

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \varepsilon}\left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)\right) \nabla_{\theta_i} Q(s, a; \theta_i)\right]$$

# ☀Related Work

- 원래 가장 성공적인 RL의 예시는 TD-gammon이었다. (value function approximator로 one-hidden layer를 가진 multi-layer perceptron 사용)
  그러나 예전 chess 연구에서 밝혀졌듯이 여기서 우리의 randomness는 state exploration과 value function을 smooth하게 한 TD-gammon의 특성 때문에 얻어지던 성공이 아니냐는 말들 돌았다.
- 더구나, model-free RL (ex. Q-learning)를 비선형 approximator or off-policy와 결합하면 많은 Q-network가 수렴하지 않고 발산한다는 게 밝혀졌다.

- 그래서, 최근에는 DL의 발전에 따라 DL과 RL을 combine 하려는 시도가 이뤄지고 있다.
  먼저 divergence 문제를 gradient temporal-difference 라며 다뤄봤다.
  이 방법론 fixed policy를 nonlinear function approximator로 평가 혹은 control policy를 linear function approximator의 조합으로 수렴하기게 한다. ✓
  그러나 이런 기법은 아직으로 nonlinear control로도 확장되지 않았다. (향후 할 것이다.)

- 우리 연구와 가장 비슷한 기법이 얼마전 연구된 neural fitted Q-learning (NFQ)이다.
  NFQ는
    - RPROP 알고리즘으로 less function을 최적화 한다.
    - 그러나 여기서 사용하는 batch update는 data set의 크기에 따라 계산 비용이 커진다.
    - 그래서 batch update 대신 Stochastic gradient update를 쓴다.
    - 또한 NFQ는 autoencoders 가능 pixels의 low-dimensional로 coding 다음 NFQ에 넣는다.
    - 반면 이 연구는 visual input (raw pixel)를 바로 input으로 사용하고, 최종에서 몇가지 강화과정을 한 사용한다.

- 설명 아니다

<br/>

# ☀DQN

- CV와 음성 인식이 발전하면서 raw input은 바로 받아서 처리할 수 있게 됐다.
  DNN은 SGD에 의한 lightweight update을 이용한다.
  → 우리의 목표도 RGB image으로 직접 받아서 SGD update에 사용하여 훈련하며, DNN과 RL과 결합할 모델을 만들 것이다.

- Tesauro의 TD-gammon이 온 성공과 뱃추다
  여기서 value function을 추정하는 parameter를 on-policy으로 update했다.
  20년 전이고 이 결과가 이전보다 뛰어났었으므로 따라서, hardware의 진보가 받은 20년 후인 지금, 현대 DNN과 RL을 결합한다면 성공한 버전을 볼 수 있을 것이다.

- 그래서 기존 연구와 다르게, 우린 experience replay라는 걸 사용한다.
  experience$_t = (s_t, a_t, r_t, s_{t+1})$
  dataset $= e_1, e_2, \dots e_N$
- 알고리즘의 내부 loop에서 DQN sample randomly 갖고와서 Q-learning update or minibatch update를 수행한다.
  이 다음, agent는 $\epsilon$-greedy에 의해 action을 한다.
  또한 input로 정해진 history를 사용하면 길이가 계속 바뀌니까 $\phi$ function 이용해 길이를 고정시킨다.
      (계속 늘어나는 여부) (바로 직전의 frame만 쓰는 식)



Algorithm 1 Deep Q-learning with Experience Replay
Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
for episode $= 1, M$ do
  Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
  for $t = 1, T$ do
    With probability $\epsilon$ select a random action $a_t$
    otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
    Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
    Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
    Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
    Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
    Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
    Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
  end for
end for

- 장점 1. 모든 step의 experience가 나중에 잠재적으로 쓰일수 있으니 data efficiency 좋음.

- 장점 2. 연속된 samples 학습하는 건 samples 간 연관성 때문에 효율이 나쁨.
           그러니 이렇게 experience를 저장해 두고 randomly 꺼내 쓰는 게 효율이 나음.

- 장점 3. off-policy 효능.
           Data distribution의 매우 큰 편차를 평균화시켜 없앤다.
           on-policy일 경우, 그건 지금의 optimal한 data를 강조하므로 편향되어 따라서 local optimum에 빠지거나 diverge할 수 있다.

- 한계점.
  현재 버전에서는 최근 N개의 experience만 저장하고, update 시 randomness가 거의 없다.
  → 그러니 균형적 업데이트를 자주말로 말하는.
  또한 중요한 uniform sampling로 replay memory에서 특정 순으로 부여해 부분에 역점이 있지가 된다.

# ☀ Preprocessing & Model Architecture.

## ✱ Experiments

- 7개의 서로 다른 Atari games 데모으로 실험을 진행했다. (사전 정보 없이)
  그래도 항상 algorithm과 hyperparameter는 사용하여 모델의 간편함을 보였다.
  역시 하나, reward 크기가 scaling 없다. 성공은 0, 신규 +1, 아니면 0. → error 추이로, learning rate 안 바꿔도 되고, agent의 performance로는 민감성을 준.

- size=32의 미니배치를 가진 RMSProp 알고리즘 사용.
  ∈는 1에서 선형이 0.1로 선형적으로 감소 후 0.1 유지.

- Frame-skipping technique.
  매 frame 보고 action 취하기 아니라 action을 고르는 거 사이에 많이 길게기 때문에 빠른거라서.
  따라서 k번의 frame마다 action 취해 취하고, skip된 frame은 직전의 action을 취함.
  Space Invader라는 게임은 k=3이고 나머진 k=4.

── Training & Stability



- 지도학습과 달리, 강화학습은 학습 도중에 progress를 평가하기 어렵다.
  각각 average total reward로 하면, 약간의 weight change에도 policy가 state의 방문에 영향이 여러 noise될 수.
  그래서 action-value function Q로 두면, 특정 state에서의 항상 policy를 따랐을 때 reward 얻어나 변동 적 예측 가능.
  → 훨씬 stable하고 변화도 거의 없었어.
  ⇒ 비록 이론적인 수렴은 보장은 없지만, RL과 SGD를 사용하여 large N.N을 안정적으로 학습할 수 있다.

── Visualizing the Value Function.



A 적이 나타내나 predicted value가 증가하는 모습

B 적을 명중시 직전, predicted value 급증함

C 적이 순간 난 후면 rewards 면 기대가 크로 사라지므로 기대가 크로 내림

── Main Evaluation.

| | B. Rider | Breakout | Enduro | Pong | Q*bert | Seaquest | S. Invaders |
|---|---|---|---|---|---|---|---|
| Random | 354 | 1.2 | 0 | −20.4 | 157 | 110 | 179 |
| Sarsa [3] | 996 | 5.2 | 129 | −19 | 614 | 665 | 271 |
| Contingency [4] | 1743 | 6 | 159 | −17 | 960 | 723 | 268 |
| DQN | 4092 | 168 | 470 | 20 | 1952 | 1705 | 581 |
| Human | 7456 | 31 | 368 | −3 | 18900 | 28010 | 3690 |
| HNeat Best [8] | 3616 | 52 | 106 | 19 | 1800 | 920 | 1720 |
| HNeat Pixel [8] | 1332 | 4 | 91 | −16 | 1325 | 800 | 1145 |
| DQN Best | 5184 | 225 | 661 | 21 | 4500 | 1740 | 1075 |

- Sarsa나 Contingency는 게임에 대한 사전 지식을 알 수 있으므로 DQN이 성능이 더 낫다.