CS 2770: Computer Vision Self-Supervised and Embodied Learning

Prof. Adriana Kovashka University of Pittsburgh April 9, 2019

## Motivation

- What's the data we've learned from thus far?
- Labeled static datasets
  - Expensive to obtain
  - Doesn't match how humans learn
- Alternatives
  - Unsupervised learning (no labels) next time
  - Self-supervised learning ("fake"/emergent labels)
  - Embodied/active learning (agents in environments)

## Self-supervised learning

## Unsupervised Visual Representation Learning by Context Prediction

#### Carl Doersch, Alexei Efros and Abhinav Gupta ICCV 2015

## ImageNet + Deep Learning







- Image Retrieval
- Detection (RCNN)
- Segmentation (FCN)
- Depth Estimation

## ImageNet + Deep Learning



## **Context as Supervision**

[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often anded, without resontment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal raile but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet sweepers of Lilliput, but she knew that most adult visitors would







## Semantics from a non-semantic task







## Architecture



Doersch et al., "Unsupervised Visual Representation Learning by Context Prediction", ICCV 2015

## What is learned?



## **Pre-Training for R-CNN**



Pre-train on relative-position task, w/o labels

Doersch et al., "Unsupervised Visual Representation Learning by Context Prediction", ICCV 2015

[Girshick et al. 2014]



Doersch et al., "Unsupervised Visual Representation Learning by Context Prediction", ICCV 2015

#### Shuffle and Learn: Unsupervised Learning using Temporal Order Verification

Ishan Misra, C. Lawrence Zitnick, and Martial Hebert ECCV 2016



Fig. 1: (a) A video imposes a natural temporal structure for visual data. In many cases, one can easily verify whether frames are in the correct temporal order (shuffled or not). Such a simple sequential verification task captures important spatiotemporal signals in videos. We use this task for unsupervised pre-training of a Convolutional Neural Network (CNN). (b) Some examples of the automatically extracted positive and negative tuples used to formulate a classification task for a CNN.



Fig. 2: (a) We sample tuples of frames from high motion windows in a video. We form positive and negative tuples based on whether the three input frames are in the correct temporal order. (b) Our triplet Siamese network architecture has three parallel network stacks with shared weights upto the fc7 layer. Each stack takes a frame as input, and produces a representation at the fc7 layer. The concatenated fc7 representations are used to predict whether the input tuple is in the correct temporal order.

Table 2: Mean classification accuracies over the 3 splits of UCF101 and HMDB51 datasets. We compare different initializations and finetune them for action recognition.

Dataset	Initialization	Mean Accuracy
UCF101	Random (Ours) Tuple verification	38.6 <b>50.2</b>
HMDB51	Random UCF Supervised (Ours) Tuple verification	13.3 15.2 <b>18.1</b>

# Learning image representations tied to ego-motion

#### Dinesh Jayaraman and Kristen Grauman ICCV 2015

#### The kitten carousel experiment [Held & Hein, 1963]



## Problem with today's visual learning

**Status quo**: Learn from "disembodied" bag of labeled snapshots.

# **Our goal:** Learn in the context of acting and moving in the world.





## Our idea: Ego-motion ↔ vision

**Goal:** Teach computer vision system the connection: "how I move" ↔ "how my visual surroundings change"



#### **Ego-motion motor signals**

**Unlabeled video** 

## **Ego-motion** ↔ **vision**: view prediction



#### After moving:



## **Ego-motion** ↔ **vision** for recognition

Learning this connection requires:

- Depth, 3D geometry
- Semantics
- Context

Also key to recognition!

Can be learned without manual labels!

**Our approach:** unsupervised feature learning using egocentric video + motor signals

## Approach idea: Ego-motion equivariance

Invariant features: unresponsive to some classes of transformations

 $\mathbf{z}(g\mathbf{x}) \approx \mathbf{z}(\mathbf{x})$ 

Equivariant features : predictably responsive to some classes of transformations, through simple mappings (e.g., linear) "equivariance map"

 $\mathbf{z}(g\mathbf{x}) \approx \mathbf{M}_{g}\mathbf{z}(\mathbf{x})$ 

#### Invariance <u>discards</u> information; equivariance <u>organizes</u> it.

## Approach idea: Ego-motion equivariance

## **Training data**

# Unlabeled video + motor signals



#### Equivariant embedding organized by ego-motions

Pairs of frames related by similar ego-motion should be related by same feature transformation

## **Approach overview**

**Our approach:** unsupervised feature learning using egocentric video + motor signals

- 1. Extract training frame pairs from video
- 2. Learn ego-motion-equivariant image features
- 3. Train on target recognition task in parallel

## **Training frame pair mining**

#### **Discovery of ego-motion clusters**



#### forward distance



## **Ego-motion equivariant feature learning**



## Summary



## **Results: Recognition**

#### Learn from unlabeled car video (KITTI)















Geiger et al, IJRR '13

# Exploit features for static scene classification (SUN, 397 classes)



Xiao et al, CVPR '10

## **Results: Recognition**

Do ego-motion equivariant features improve recognition?



6 labeled training examples per class

### The Curious Robot: Learning Visual Representations via Physical Interactions

Lerrel Pinto, Dhiraj Gandhi, Yuanfeng Han, Yong-Lae Park, and Abhinav Gupta ECCV 2016

## **Embodied representations**



**Physical Interaction Data** 



Conv Layer1 Filters



**Conv3 Neuron Activations** 

#### **Learned Visual Representation**



**Conv5 Neuron Activations** 

Pinto et al., "The Curious Robot: Learning Visual Representations via Physical Interactions", ECCV 2016

## Grasping



Fig. 2. Examples of successful (left) and unsuccesful grasps (right). We use a patch based representation: given an input patch we predict 18-dim vector which represents whether the center location of the patch is graspable at  $0^{\circ}$ ,  $10^{\circ}$ , ...  $170^{\circ}$ .
# Pushing

Objects and push action pairs



Fig. 4. Examples of initial state and final state images taken for the push action. The arrows demonstrate the direction and magnitude of the push action.

Pinto et al., "The Curious Robot: Learning Visual Representations via Physical Interactions", ECCV 2016

# Poking

Objects and poke tactile response pairs



Fig. 6. Examples of the data collected by the poking action. On the left we show the object poked, and on the right we show force profiles as observed by the tactile sensor.

#### Pose/viewpoint invariance



Fig. 7. Examples of objects in different poses provided to the embedding network.

Pinto et al., "The Curious Robot: Learning Visual Representations via Physical Interactions", ECCV 2016

## **Representations from interactions**



Fig. 8. Our shared convolutional architecture for four different tasks.

## Classification/retrieval performance



Fig. 10. The first column corresponds to query image and rest show the retrieval. Note how the network learns that cups and bowls are similar (row 5).

## Classification/retrieval performance

Table 1. Classification accuracy on ImageNet Household, UW RGBD and Caltech-256

	Household	UW RGBD	Caltech-256
Root network with random init.	0.250	0.468	0.242
Root network trained on robot tasks ( <b>ours</b> )	0.354	0.693	0.317
AlexNet trained on ImageNet	0.625	0.820	0.656

#### Table 2. Image Retrieval with Recall@k metric

	Instance level			Category level				
	k=1	k=5	k=10	k=20	k=1	k=5	k=10	k=20
Random Network	0.062	0.219	0.331	0.475	0.150	0.466	0.652	0.800
Our Network	0.720	0.831	0.875	0.909	0.833	0.918	0.946	0.966
AlexNet	0.686	0.857	0.903	0.941	0.854	0.953	0.969	0.982

#### Object-Graphs for Context-Aware Category Discovery

#### Yong Jae Lee and Kristen Grauman CVPR 2010

## Goal



 Discover *new* object categories, based on their relation to categories for which we have trained models

Lee and Grauman, "Object-Graphs for Context-Aware Category Discovery", CVPR 2010

# Existing approaches

Previous work treats unsupervised visual discovery as an appearance-grouping problem.



Can you identify the recurring pattern?

## Our idea

How can seeing previously learned objects in novel images help to discover *new* categories?



Can you identify the recurring pattern?

## Our idea

Discover visual categories within unlabeled images by modeling interactions between the unfamiliar regions and familiar objects.



Can you identify the recurring pattern?

#### Context-aware visual discovery









Lee and Grauman, "Object-Graphs for Context-Aware Category Discovery", CVPR 2010



#### Clusters from region-region affinities



Lee and Grauman, "Object-Graphs for Context-Aware Category Discovery", CVPR 2010

#### **Object Discovery Accuracy**







MSRC-v2



**PASCAL 2008** 



MSRC-v0



Corel

## **Examples of Discovered Categories**



## Discussion

- Many types of supervision have been tried
- What other types of supervision "for free" can we use?
- How do we know if a certain supervision type would work?
- Can we make this type of learning perform on par with supervised learning?

#### **Embodied learning**

#### **Reinforcement Learning**



#### **Cart-Pole Problem**



**Objective**: Balance a pole on top of a movable cart

State: angle, angular speed, position, horizontal velocityAction: horizontal force applied on the cartReward: 1 at each time step if the pole is upright

### Atari Games



**Objective**: Complete the game with the highest score

State: Raw pixel inputs of the game state Action: Game controls e.g. Left, Right, Up, Down Reward: Score increase/decrease at each time step

#### Go



Objective: Win the game!

State: Position of all piecesAction: Where to put the next piece downReward: 1 if win at the end of the game, 0 otherwise

#### How can we mathematically formalize the RL problem?



## **Markov Decision Process**

- Mathematical formulation of the RL problem
- Markov property: Current state completely characterises the state of the world

Defined by:  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$ 

- ${\mathcal S}\,$  : set of possible states
- ${\cal A}\,$  : set of possible actions
- ${\boldsymbol{\mathcal{R}}}$  : distribution of reward given (state, action) pair
- $\mathbb{P}$ : transition probability i.e. distribution over next state given (state, action) pair
- $\gamma$ : discount factor

## **Markov Decision Process**

- At time step t=0, environment samples initial state  $s_0 \sim p(s_0)$
- Then, for t=0 until done:
  - Agent selects action a<sub>t</sub>
  - Environment samples reward  $r_t \sim R(. | s_t, a_t)$
  - Environment samples next state  $s_{t+1} \sim P(. | s_t, a_t)$
  - Agent receives reward r<sub>t</sub> and next state s<sub>t+1</sub>
- A policy u is a function from S to A that specifies what action to take in each state
- **Objective**: find policy u\* that maximizes cumulative discounted reward:







Set a negative "reward" for each transition (e.g. r = -1)

**Objective:** reach one of terminal states (greyed out) in least number of actions

#### A simple MDP: Grid World



**Random Policy** 

**Optimal Policy** 

Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung

#### The optimal policy u\*

We want to find optimal policy u\* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?

#### The optimal policy u\*

We want to find optimal policy u\* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)? Maximize the **expected sum of rewards!** 

Formally: 
$$\pi^* = \arg \max_{\pi} \mathbb{E}\left[\sum_{t \ge 0} \gamma^t r_t | \pi\right]$$
 with  $s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$ 

#### Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths)  $s_0$ ,  $a_0$ ,  $r_0$ ,  $s_1$ ,  $a_1$ ,  $r_1$ , ...

#### How good is a state?

The value function at state s, is the expected cumulative reward from following the policy from state s:  $V^{T}(\cdot) = \mathbb{E}\left[\sum_{i=1}^{n} t_{i} + \frac{1}{2} \right]$ 

$$V^{\pi}(s) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi
ight]$$

#### How good is a state-action pair?

The **Q-value function** at state s and action a, is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^{\pi}(s,a) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi
ight]$$

## **Bellman equation**

The optimal Q-value function Q\* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s,a) = \max_{\pi} \mathbb{E}\left[\sum_{t \ge 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi
ight]$$

Q\* satisfies the following Bellman equation:

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

**Intuition:** if the optimal state-action values for the next time-step Q\*(s',a') are known, then the optimal strategy is to take the action that maximizes the expected value of  $r + \gamma Q^*(s',a')$ 

The optimal policy u\* corresponds to taking the best action in any state as specified by Q\*

# Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

 $Q(s,a;\theta) \approx Q^*(s,a)$  function parameters (weights)

If the function approximator is a deep neural network => deep q-learning!

### **Q-network Architecture**



**Current state s<sub>t</sub>: 84x84x4 stack of last 4 frames** (after RGB->grayscale conversion, downsampling, and cropping)

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 Initialise 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  $(\phi_i, a_i, r_i, \phi_{i+1})$  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

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory  $\mathcal{D}$  to capacity N Initialize replay memory, Q-network Initialize action-value function Q with random weights for episode = 1, M do Initialise 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

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory  $\mathcal{D}$  to capacity N Initialize action-value function Q with random weights ——— Play M episodes (full games) for episode = 1, M do Initialise 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

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 Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ Initialize state for t = 1, T do (starting game With probability  $\epsilon$  select a random action  $a_t$ screen pixels) at the otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ beginning of each Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ episode 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
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 Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ for t = 1, T do For each timestep t With probability  $\epsilon$  select a random action  $a_t$ of the game 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

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 Initialise 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  $a_t$ With small probability, otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ select a random Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ action (explore), Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ otherwise select Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ greedy action from Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ current policy 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

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 Initialise 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  $x_{t+1}$ Take the action  $(a_i)$ , 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  $\mathcal{D}$ and observe the Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ reward r, and next Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ state s<sub>t+1</sub> Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3 end for end for

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 Initialise 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 in Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ replay memory 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

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 Initialise 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  $\mathcal{D}$ Experience Replay: Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Sample a random minibatch of transitions Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3 from replay memory end for and perform a gradient end for descent step

### **Policy Gradients**

What is a problem with Q-learning? The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

But the policy can be much simpler: just close your hand Can we learn a policy directly, e.g. finding the best policy from a collection of policies?

#### **Policy Gradients**

Formally, let's define a class of parameterized policies:  $\Pi = \{\pi_{\theta}, \theta \in \mathbb{R}^m\}$ 

For each policy, define its value:

$$J( heta) = \mathbb{E}\left[\sum_{t \ge 0} \gamma^t r_t | \pi_{ heta}
ight]$$

We want to find the optimal policy  $\theta^* = \arg \max_{\theta} J(\theta)$ 

How can we do this?

Gradient ascent on policy parameters!

Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung

# REINFORCE Algorithm (orig. Williams 1992)

## Gradient estimator: $\nabla_{\theta} J(\theta) \approx \sum_{t \ge 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

#### Interpretation:

- If r(r) is high, push up the probabilities of the actions seen
- If r(r) is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. But in expectation, it averages out!

#### **Policy Gradients**



### **Policy Gradients**

- Loss:  $\sum i Ai \log p(y_i | x_i)$
- x<sub>i</sub> = state
- y<sub>i</sub> = sampled action
- A<sub>i</sub> = "advantage" e.g. +1/-1 for win/lose in simplest version, or discounted, or improvement over "baseline"

### Policy Gradients vs Q-Learning

- Estimating exact value of state-action pairs vs choosing what actions to take (value not important)
- Policy gradients can handle continuous action spaces (Gaussian policy)
- Step-by-step (did I correctly estimate the reward at this time) vs delayed feedback (run policy and wait until game terminates)
- Policy gradients suffers from high variance and instability; might want to make gradients smaller (e.g. relative to a baseline)

### Actor-Critic Algorithm

We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q-function).

- The actor decides which action to take, and the critic tells the actor how good its action was and how it should adjust
- Also alleviates the task of the critic as it only has to learn the values of (state, action) pairs generated by the policy
- Can also incorporate Q-learning tricks e.g. experience replay
- **Remark:** we can define by the **advantage function** how much an action was better than expected  $A^{\pi}(s, a) = Q^{\pi}(s, a) V^{\pi}(s)$

#### Actor-Critic Algorithm



Policy Update: 
$$\Delta heta = lpha \ 
abla_ heta (log \ \pi_ heta(s,a)) \hat{q}_w(s,a)$$

q learning function approximation (estimate action value)

Value update: 
$$\Delta w = \beta \left( \frac{R(s, a) + \gamma \hat{q}_w(s_{t+1}, a_{t+1}) - \hat{q}_w(s_t, a_t)}{\mathsf{TD} \text{ error}} \right) \nabla_w \hat{q}_w(s_t, a_t)$$
Gradient of our value function
Gradient of our value function

**Thomas Simonini** 

#### Example Q-network



Figure 3: Architecture of our proposed Q-network. It receives histories of top-ranked images, positive and negative proxy images, and taken actions. It predicts the best action given a specific state. Inputs are denoted with dotted lines. Please see text for further explanation.

### Example Q-network (actions)



Figure 1: We learn how to intelligently combine different forms of user feedback for interactive image search, and find the user's desired content in fewer iterations. The *image search* section depicts our search agent that predicts an appropriate action at a certain iteration. For example, our agent selects free-form attribute feedback for iteration 1, and sketching for iteration 2. The *actions* section presents the three possible interactions (actions) of our agent.

#### **REINFORCE** in action: Recurrent Attention Model (RAM)



### RL for object detection



Figure 1. A sequence of actions taken by the proposed algorithm to localize a cow. The algorithm attends regions and decides how to transform the bounding box to progressively localize the object.

Caicedo and Lazebnik, "Active Object Localization with Deep Reinforcement Learning", ICCV 2015

### **RL** for object detection



Figure 2. Illustration of the actions in the proposed MDP, giving 4 degrees of freedom to the agent for transforming boxes.

$$R_a(s,s') = sign\left(IoU(b',g) - IoU(b,g)\right) \qquad R_{\omega}(s,s') = \begin{cases} +\eta & \text{if } IoU(b,g) \ge \tau \\ -\eta & \text{otherwise} \end{cases}$$



Caicedo and Lazebnik, "Active Object Localization with Deep Reinforcement Learning", ICCV 2015

### **RL** for navigation



Fig. 1. The goal of our deep reinforcement learning model is to navigate towards a visual target with a minimum number of steps. Our model takes the current observation and the image of the target as input and generates an action in the 3D environment as the output. Our model learns to navigate to different targets in a scene without re-training.

#### **RL** for navigation



Figure 1: Our goal is to use scene priors to improve navigation in unseen scenes and towards novel objects. (a) There is no mug in the field of view of the agent, but the likely location for finding a mug is the cabinet near the coffee machine. (b) The agent has not seen a mango before, but it infers that the most likely location for finding a mango is the fridge since similar objects such as apple appear there as well. The most likely locations are shown with the orange box.



Figure 2: **Overview of the architecture.** Our model to incorporate semantic knowledge into semantic navigation. Specifically, we learn a policy network that decides an action based on the visual features of the current state, the semantic target category feature and the features extracted from the knowledge graph. We extract features from the parts of the knowledge graph that are activated.

### **RL for question-answering**



Figure 1: Embodied Question Answering – EmbodiedQA– tasks agents with navigating rich 3D environments in order to answer questions. These agents must jointly learn language understanding, visual reasoning, and goal-driven navigation to succeed.

#### **RL** for question-answering



Figure 4: Our PACMAN navigator decomposes navigation into a planner and a controller. The planner selects actions and the controller executes these actions a variable number of times. This enables the planner to operate on shorter timescales, strengthening gradient flows.

#### What's next?