


**TU/e** EINDHOVEN  
UNIVERSITY OF  
TECHNOLOGY

## Module 11: Sequence modeling and reinforcement learning

5LSM0: Convolutional neural networks for computer vision

Fons van der Sommen

Electrical Engineering / VCA research group



## Administrative

### Paper sessions

- Next Thursday (March 28<sup>th</sup>) hours 1-4
- Please indicate a preferred time-slot on Canvas (also check for more details)
- Upload presentations through Canvas Assignments by Wednesday March 27<sup>th</sup> 23:59
  - *Otherwise penalty on grade for presentation*

### Final assignment

- Detailed announcement next week...



### Last time(s)

**Auto-encoders**

**Variational Auto-encoders**

$$\mathcal{L}(x^{(i)}, \theta, \phi) = \mathbb{E}_z[\log p_\theta(x^{(i)}|z)] - D_{KL}(q_\phi(z|x^{(i)})||p_\theta(z))$$

3 SLSMO Module 11: Sequence modeling and reinforcement learning

### Last time(s)

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log \left( 1 - D_{\theta_d}(G_{\theta_g}(z)) \right) \right]$$

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log \left( 1 - D_{\theta_d}(G_{\theta_g}(z)) \right) \right]$$

Gradient ascent for both  $\theta_d$  and  $\theta_g$

$$\max_{\theta_g} \left[ \mathbb{E}_{z \sim p(z)} \log \left( D_{\theta_d}(G_{\theta_g}(z)) \right) \right]$$

**Generative Adversarial Networks (GANs)**

4 SLSMO Module 11: Sequence modeling and reinforcement learning

## Last time(s)

### Self-supervised learning

image → Coherence/Incoherence → learning

image → de-color → learning → colored

image → context → learning → Missing tile

Unlabeled data may reveal structure

Beware of pit-falls!

5 5LSM0 Module 11: Sequence modeling and reinforcement learning

**TU/e**

## This time

### Sequential modeling

**How can we model sequential relations using deep neural networks?**

- Recurrent Neural Networks (RNNs)

**How can we optimize gradient flow over the sequential dimension? (e.g. time)**

- Long Short-Term Memory (LSTM)

### Reinforcement learning

**Can we learn to play Atari games with deep learning?:D**

- (And do other cool stuff too...)
- Deep Q-learning

6 5LSM0 Module 11: Sequence modeling and reinforcement learning

**TU/e**

## Sequential modeling

5LSM0 Module 11: Sequence modeling and reinforcement learning

### So far...

#### CNNs to do

- Image classification; Object detection/segmentation; Density estimation for data generation

#### Mostly fixed size input → fixed size output

- If we didn't we required several passes through network or auxiliary methods

#### How would we do

- Image captioning [one image to variable number of words]
- Sentiment classification [variable number of words to one sentiment]
- Translation [variable number of words to variable number of words]

#### Recurrent Neural Networks (RNNs)



8

5LSM0 Module 11: Sequence modeling and reinforcement learning

## Recurrent neural networks

**one to one**

Example:

**one to many**

Image captioning

**many to one**

Text sentiment classification

**many to many**

Machine translation

**many to many**

Video frame classification

9 SLSMO Module 11: Sequence modeling and reinforcement learning

## Recurrent neural networks

**Recurrence relation**

- For a certain time  $t$ :

Function with parameter  $W$     **Input**

$$h_t = f_W(h_{t-1}, x_t)$$

New state                      Old state

**Straightforward option:**

- Two linear relations:
  - old state  $\rightarrow$  new state
  - Input  $\rightarrow$  new state
- Linear output relation:
  - New state  $\rightarrow$  output

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

**Note:**  
parameters  $W$   
fixed over time

10 SLSMO Module 11: Sequence modeling and reinforcement learning

## Recurrent neural networks

**RNN computational graphs**

- Unrolling the hidden states over time

Final loss  
 $L = \sum L_i$

Note: all  $f_W$  blocks use the same weight matrix  $W$

11 SLSM0 Module 11: Sequence modeling and reinforcement learning

## Recurrent neural networks

**RNN computational graphs**

Many to many

Many to one

One to many

Use single input to initialize the hidden state of the model

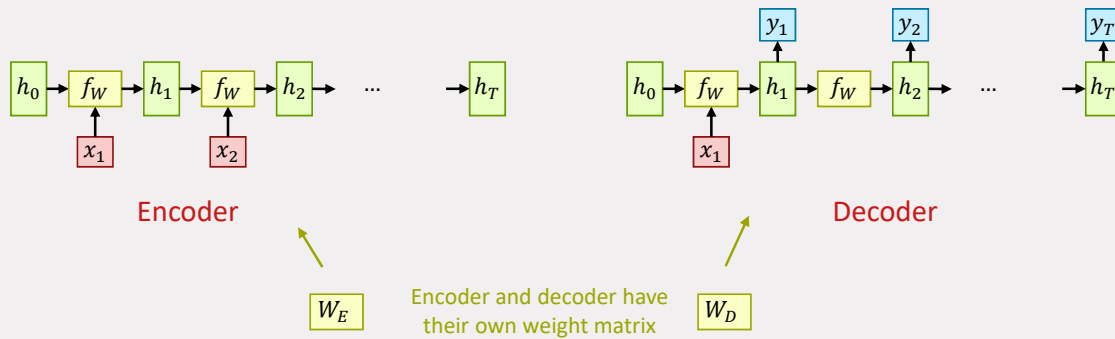
Final hidden state summarized all of the context of the entire sequence

12 SLSM0 Module 11: Sequence modeling and reinforcement learning

# Recurrent neural networks

## RNN computational graphs: Sequence to sequence

- Many-to-one followed by one-to-many



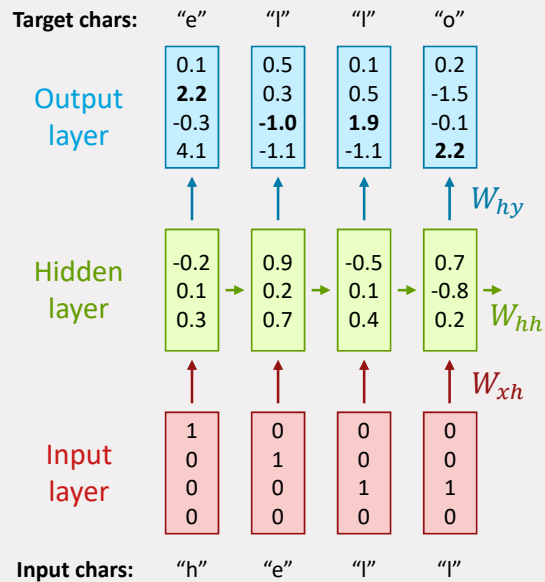
# Recurrent neural networks

## RNN computational graphs: Sequence to sequence

- Many-to-one followed by one-to-many

### Example: character-level language model

- Predict next letter
- Vocabulary: [h, e, l, o]
- Training sequence: "hello"



## Recurrent neural networks

**RNN computational graphs:**  
Sequence to sequence

- Many-to-one followed by one-to-many

**Example: character-level language model**

- Predict next letter
- Vocabulary: [h, e, l, o]
- Training sequence: "hello"

Q: Why not feed back max sample?

**At test time we can sample!**

- Start with a letter, and sample from output to use for next input in time

Sample:	"e"	"l"	"l"	"o"
Softmax	0.03 0.13 0.00 0.84	0.25 0.20 0.05 0.50	0.11 0.17 0.68 0.03	0.11 0.02 0.08 0.79
Output layer	0.1 2.2 -0.3 4.1	0.5 0.3 -1.0 -1.1	0.1 0.5 1.9 -1.1	0.2 -1.5 -0.1 2.2
Hidden layer	-0.2 0.1 0.3	0.9 0.2 0.7	-0.5 0.1 0.4	0.7 -0.8 0.2
Input layer	1 0 0 0	0 1 0 0	0 0 1 0	0 0 1 0
Input chars:	"h"	"e"	"l"	"l"

\*Example from cs231n lecture 10 slides 34-36



## Recurrent neural networks

**RNN training: backpropagation through time**

**Problem: backprop could take ages for long sequences**

\*Example from cs231n lecture 10 slides 34-36





## Recurrent neural networks

RNN training: **truncated** backpropagation through time

17 5LSM0 Module 11: Sequence modeling and reinforcement learning **TU/e**

## Recurrent neural networks

RNN training: **truncated** backpropagation through time

18 5LSM0 Module 11: Sequence modeling and reinforcement learning **TU/e**

## Recurrent neural networks

**RNN training: truncated backpropagation through time**

The diagram illustrates the training of an RNN using truncated backpropagation through time. It shows a sequence of 16 time steps. Each time step consists of an input (red box), a hidden state (green box), and an output (blue box). The forward pass (indicated by grey arrows) processes the sequence from left to right. A pink box highlights a subset of 4 time steps (steps 9-12) where a loss is calculated. The backward pass (indicated by grey arrows) flows from right to left, but is truncated at the start of the pink box, meaning gradients are not propagated back to time steps before step 9.

19 SLSM0 Module 11: Sequence modeling and reinforcement learning **TU/e**

## Recurrent neural networks

**RNN training: truncated backpropagation through time**

Q: What other training approach looks similar?

This diagram is identical to the one on slide 19, showing the forward and backward passes of an RNN with a truncated backward pass over a subset of time steps.

20 SLSM0 Module 11: Sequence modeling and reinforcement learning **TU/e**

# Cool stuff you can do with RNNs

**PANDARUS:**  
 Alas, I think he shall be come approached and the day  
 When little strain would be attain'd into being never fed,  
 And who is but a chain and subjects of his death,  
 I should not sleep.

**Second Senator:**  
 They are away this miseries, produced upon my soul,  
 Breaking and strongly should be buried, when I perish  
 The earth and thoughts of many states.

**DUKE VINCENTIO:**  
 Well, your wit is in the care of side and that.

**Second Lord:**  
 They would be ruled after this chamber, and  
 my fair nues begun out of the fact, to be conveyed,  
 Whose noble souls I'll have the heart of the wars.

**Clown:**  
 Come, sir, I will make did behold your worship.

**VIOLA:**  
 I'll drink it.

**VIOLA:**  
 Why, Salisbury must find his flesh and thought  
 That which I am not aps, not a man and in fire,  
 To show the reining of the raven and the wars  
 To grace my hand reproach within, and not a fair are hand,  
 That Caesar and my goodly father's world;  
 When I was heaven of presence and our fleets,  
 We spare with hours, but cut thy council I am great,  
 Murdered and by thy master's ready there  
 My power to give thee but so much as hell:  
 Some service in the noble bondman here,  
 Would show him to her wine.

**KING LEAR:**  
 O, if you were a feeble sight, the courtesy of your law,  
 Your sight and several breath, will wear the gods  
 With his heads, and my hands are wonder'd at the deeds,  
 So drop upon your lordship's head, and your opinion  
 Shall be against your honour.

Write in Shakespeare style



21 SLSMO Module 11: Sequence modeling and reinforcement learning

\*Example from cs231n lecture 10 slides 34-36



# Cool stuff you can do with RNNs

For  $\mathbb{A}^1_{\mathbb{Z}} = \mathbb{A}^1_{\mathbb{Z}} \times_{\mathbb{Z}} \mathbb{A}^1_{\mathbb{Z}}$  where  $\mathcal{L}_{\mathbb{Z}} = 0$ , hence we can find a closed subset  $\mathcal{H}$  in  $\mathcal{H}$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get  
 $S = \text{Spec}(R) = U \times_X U \times_X U$   
 and the comparably in the fibre product covering we have to prove the lemma generated by  $\coprod \mathcal{Z} \times_{\mathbb{Z}} U \rightarrow V$ . Consider the maps  $M$  along the set of points  $\text{Sch}_{\text{fpf}}$  and  $U \rightarrow U$  is the fibre category of  $S$  in  $U$  in Section ?? and the fact that any  $U$  affine, see Morphisms, Lemma ?? . Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $\text{Sh}(G)$  such that  $\text{Spec}(R) \rightarrow S$  is smooth or an  
 $U = \bigcup U_i \times_S U_i$   
 which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,S}$  is a scheme where  $x, x', s'' \in S'$  such that  $\mathcal{O}_{X,S'} \rightarrow \mathcal{O}_{X',S'}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_d(\mathcal{L}'/S')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}_{\mathbb{Z}}$  is a covering of  $\mathcal{X}'$ , and  $\mathcal{T}_i$  is an object of  $\mathcal{F}_{X/S}$  for  $i > 0$  and  $\mathcal{F}_p$  exists and let  $\mathcal{F}_i$  be a presheaf of  $\mathcal{O}_X$ -modules on  $C$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that  
 $\bar{M}^* = \mathcal{I}^* \otimes_{\text{Spec}(\mathbb{Z})} \mathcal{O}_{S,S} - i_X^{-1} \mathcal{F}$   
 is a unique morphism of algebraic stacks. Note that  
 $\text{Arrows} = (\text{Sch}/S)_{\text{fpf}}, (\text{Sch}/S)_{\text{fpf}}$   
 and  
 $V = \Gamma(S, \mathcal{O}) \rightarrow (U, \text{Spec}(A))$   
 is an open subset of  $X$ . Thus  $U$  is affine. This is a continuous map of  $X$  is the inverse, the groupoid scheme  $S$ .

*Proof.* See discussion of sheaves of sets.  $\square$

The result for prove any open covering follows from the less of Example ?? . It may replace  $S$  by  $N_{\text{open}, \text{étale}}$ , which gives an open subspace of  $X$  and  $T$  equal to  $S_{\text{par}}$ , see Descent, Lemma ?? . Namely, by Lemma ?? we see that  $R$  is geometrically regular over  $S$ .

**Lemma 0.1.** Assume (5) and (3) by the construction in the description.  
 Suppose  $X = \text{lim}|X|$  (by the formal open covering  $X$  and a single map  $\text{Proj}_{\mathbb{Z}}(A) = \text{Spec}(B)$  over  $U$  compatible with the complex  
 $\text{Set}(A) = \Gamma(X, \mathcal{O}_{X, \mathbb{Z}})$ .)  
 When in this case of to show that  $\mathcal{Q} \rightarrow \mathcal{C}_{\mathbb{Z}/X}$  is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If  $T$  is surjective we may assume that  $T$  is connected with residue fields of  $S$ . Moreover there exists a closed subspace  $Z \subset X$  of  $X$  where  $U$  in  $X'$  is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem  
 (1)  $f$  is locally of finite type. Since  $S = \text{Spec}(R)$  and  $Y = \text{Spec}(R)$ .)

*Proof.* This is form all sheaves of sheaves on  $X$ . But given a scheme  $U$  and a surjective étale morphism  $U \rightarrow X$ . Let  $U \cap U = \coprod_{i=1, \dots, n} U_i$  be the scheme  $X$  over  $S$  at the schemes  $X_i \rightarrow X$  and  $U = \text{lim}, X_i$ .  $\square$

The following lemma surjective retrocomposes of this implies that  $\mathcal{F}_{\mathbb{Z}_n} = \mathcal{F}_{\mathbb{Z}_n} = \mathcal{F}_{X, \dots, 0}$ .

**Lemma 0.2.** Let  $X$  be a locally Noetherian scheme over  $S$ .  $E = \mathcal{F}_{X/S}$ . Set  $\mathcal{I} = \mathcal{I}_i \subset \mathcal{I}_n$ . Since  $\mathcal{I}^n \subset \mathcal{I}^n$  are nonzero over  $i_0 \leq p$  is a subset of  $\mathcal{J}_{n,0} \circ \bar{A}_2$  works.

**Lemma 0.3.** In Situation ?? . Hence we may assume  $q' = 0$ .

*Proof.* We will use the property we see that  $p$  is the next functor (??). On the other hand, by Lemma ?? we see that  
 $D(\mathcal{O}_X) = \mathcal{O}_X(D)$   
 where  $K$  is an  $F$ -algebra where  $\delta_{n+1}$  is a scheme over  $S$ .  $\square$

Produce mathy LaTeX source code



22 SLSMO Module 11: Sequence modeling and reinforcement learning

\*Example from cs231n lecture 10 slides 34-36



## Cool stuff you can do with RNNs

```
static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << 1))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000ffffffff) & 0x000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = &of_changes[PAGE_SIZE];
    rek_controls(offset, idx, &offset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}
```

Produce code



23 5LSM0 Module 11: Sequence modeling and reinforcement learning

\*Example from cs231n  
lecture 10 slides 34-36

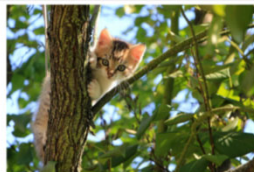
TU/e

## Cool stuff you can do with RNNs

Image  
captioning



A cat sitting on a  
suitcase on the floor



A cat is sitting on a tree  
branch



A dog is running in the  
grass with a frisbee



A white teddy bear sitting in  
the grass



Two people walking on  
the beach with surfboards



A tennis player in action  
on the court



Two giraffes standing in a  
grassy field



A man riding a dirt bike on  
a dirt track



24 5LSM0 Module 11: Sequence modeling and reinforcement learning

\*Example from cs231n  
lecture 10 slides 34-36

TU/e

## RNN gradient flow

**What happens during RNN training?**

During backwards pass need to multiply by  $W$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$= \tanh\left((W_{hh} \quad W_{xh}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

$$= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

We can stack  $h_{t-1}$  and  $x_t$  together

25 5LSM0 Module 11: Sequence modeling and reinforcement learning

## RNN gradient flow

**What happens during RNN training?**

Backprop involves successive multiplication by  $W^T$

Largest singular value  $> 1$   
**Exploding gradient**

➔


Hack: gradient clipping

Largest singular value  $< 1$   
**Vanishing gradient**

➔

Change architecture: LSTMs

26 5LSM0 Module 11: Sequence modeling and reinforcement learning

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation (1997) 

## RNN gradient flow

### Long Short-Term Memory (LSTM)

Vanilla RNN


$$h_t = \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

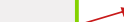
LSTM


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma(\cdot) \\ \sigma(\cdot) \\ \sigma(\cdot) \\ \tanh(\cdot) \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$


$$h_t = o \odot \tanh(c_t)$$

Cell state 

Hidden state 



27 SLSMO Module 11: Sequence modeling and reinforcement learning





1997



28 SLSMO Module 11: Sequence modeling and reinforcement learning

SNOWTOP GAPS



## Long Short Term Memory (LSTM)

Introduction of a "cell state"

The diagram shows a yellow box labeled 'W' on the left. To its right are two stacked boxes: a red one labeled 'x' and a green one labeled 'h'. A grey arrow points to the right, leading to another yellow box labeled 'W'. To its right are four stacked boxes representing activation functions: 'sigmoid' (white), 'sigmoid' (white), 'sigmoid' (white), and 'tanh' (white). To the right of these are four colored boxes representing gates: 'i' (blue), 'f' (pink), 'o' (yellow), and 'g' (purple).

i	Input gate
f	Forget gate
o	Output gate
g	... Gate

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$
  

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma(\cdot) \\ \sigma(\cdot) \\ \sigma(\cdot) \\ \tanh(\cdot) \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

29 SLSM0 Module 11: Sequence modeling and reinforcement learning

## Long Short Term Memory (LSTM)

$c_t = f \odot c_{t-1} + i \odot g$

$h_t = o \odot \tanh(c_t)$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma(\cdot) \\ \sigma(\cdot) \\ \sigma(\cdot) \\ \tanh(\cdot) \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

The diagram shows a yellow rounded rectangle representing the LSTM cell. On the left, a blue arrow labeled  $h_{t-1}$  and a red arrow labeled  $c_{t-1}$  enter. On the right, a blue arrow labeled  $h_t$  and a red arrow labeled  $c_t$  exit. Inside the cell, a blue arrow labeled  $x_t$  enters from the bottom. It goes to a 'stack' box, which then branches to four sigmoid ( $\sigma(\cdot)$ ) and one tanh ( $\tanh(\cdot)$ ) activation functions. These produce gates  $f$ ,  $i$ ,  $g$ , and  $o$ . A weight matrix  $W$  is multiplied (indicated by a circle with an 'x') with the stacked input. The outputs of the sigmoid functions are multiplied (indicated by circles with dots) with the previous cell state  $c_{t-1}$  and the current cell state  $c_t$ . The outputs of the tanh function and the  $i$  gate are multiplied together. The results are summed (indicated by a circle with a '+'). The output of the  $o$  gate is multiplied (indicated by a circle with a dot) with the sum to produce the new cell state  $c_t$ . The output of the  $o$  gate is multiplied (indicated by a circle with a dot) with the new cell state  $c_t$  to produce the new hidden state  $h_t$ .

During backprop from  $c_t$  to  $c_{t-1}$ :

- only elementwise multiplication by  $f$
- no matrix multiplication by  $W$

30 SLSM0 Module 11: Sequence modeling and reinforcement learning



## Long Short Term Memory (LSTM)

← Uninterrupted gradient flow

Q: What CNN architecture looks similar?

31 SLSM0 Module 11: Sequence modeling and reinforcement learning

Vanilla RNN

LSTM

32 SLSM0 Module 11: Sequence modeling and reinforcement learning



## Temporal modeling summary

### Recurrent Neural Networks (RNNs)

- Add recurrent relation in the network to model sequential relations
- Unroll the network over time and apply standard backpropagation to train
- In practice: vanilla RNNs don't work that well due to vanishing/exploding gradients

### Long Short-Term Memory (LSTM)

- Introduce an additional cell state, governing the input and output of the hidden state
  - *Input-, output-, forget-, "gating"- gates*
- Uninterrupted gradient flow



33 SLSMO Module 11: Sequence modeling and reinforcement learning

**TU/e**

**TU/e** EINDHOVEN  
UNIVERSITY OF  
TECHNOLOGY

## Reinforcement learning

Slides from cs231n lecture 14 [1-62]

([http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture14.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture14.pdf))

SLSMO Module 11: Sequence modeling and reinforcement learning

## Machine learning so far...

### Supervised learning

- We have data and we have labels

### Semi-/self supervised learning

- We have data and we kinda create our own labels

### Unsupervised learning

- We have data, but no labels

### Reinforcement learning

- We have no data..
  - *(But we do have an environment that can provide us with data! And even with rewards!)*



35 5LSM0 Module 11: Sequence modeling and reinforcement learning

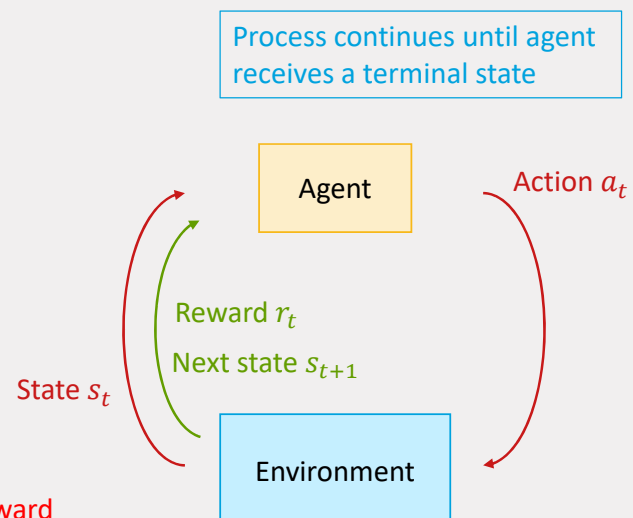
TU/e

## Reinforcement learning

### Reinforcement learning (RL) setup

- Environment
  - Gives agent a state  $s_t$
- Agent
  - Takes action  $a_t$
  - Receives reward  $r_t$  and next state  $s_{t+1}$  from environment

→ Goal: take actions that maximize reward

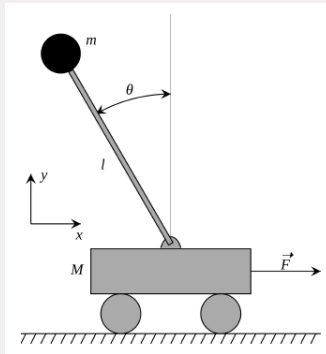


36 5LSM0 Module 11: Sequence modeling and reinforcement learning

TU/e

## Reinforcement learning

### Example: cart-pole problem



#### Objective:

- Balance a pole on top of a movable cart

#### State:

- Angle, angular speed, position, horizontal velocity

#### Action:

- Horizontal force applied on the cart

#### Reward:

- 1 at each time step if the pole is upright

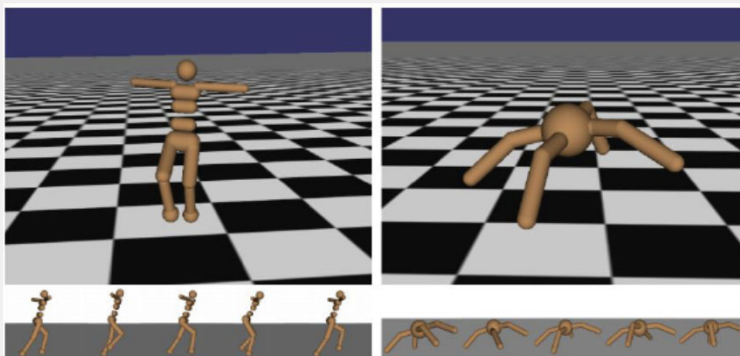


37 5LSM0 Module 11: Sequence modeling and reinforcement learning

TU/e

## Reinforcement learning

### Example: Robot locomotion



#### Objective:

- Make the robot move forward

#### State:

- Angle and position of the joints

#### Action:

- Torques applied on joints

#### Reward:

- 1 at each time step upright + forward movement



38 5LSM0 Module 11: Sequence modeling and reinforcement learning

TU/e

## Reinforcement learning

### Example: 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:

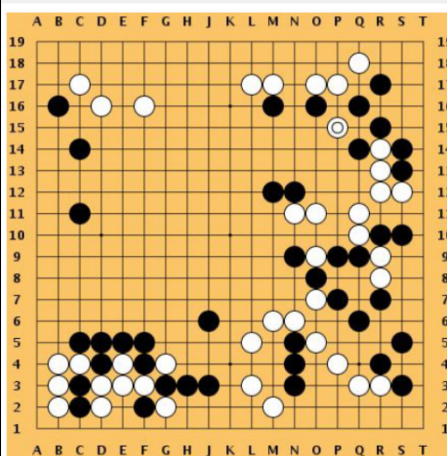
- Score increase/decrease at each time step



39 5LSM0 Module 11: Sequence modeling and reinforcement learning

TU/e

## Reinforcement learning



### Example: Go

#### Objective:

- Win the game!

#### State:

- Position of all the pieces

#### Action:

- Where to put the next piece down

#### Reward:

- 1 if win the game, 0 otherwise



40 5LSM0 Module 11: Sequence modeling and reinforcement learning

TU/e

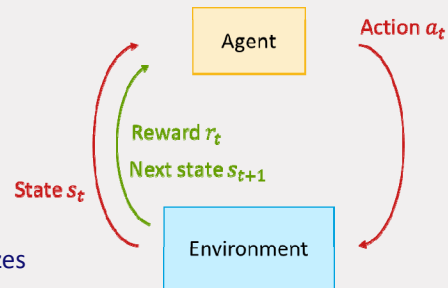
## Reinforcement learning

### How to mathematically formalize this?

#### Markov Decision Process (MDP)

- Markov Property: Current state completely characterizes the state of the world.
- Defined by  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

- $\mathcal{S}$  Set of possible states
- $\mathcal{A}$  Set of possible actions
- $\mathcal{R}$  Distribution over reward, given state-action pair
- $\mathbb{P}$  Transition probability: distribution over next state given state-action pair
- $\gamma$  Discount factor: how much do we value early vs. late rewards



41 SLSMO Module 11: Sequence modeling and reinforcement learning

TU/e

## Reinforcement learning

### How does the MDP work?

- At time step  $t = 0$ , environment samples initial state  $s_0$  from initial state distribution  $p(s_0)$
- Then for  $t = 0$  until done:
  - *Agent selects action  $a_t$*
  - *Environment samples reward  $r_t \sim \mathcal{R}(\cdot | s_t, a_t)$*
  - *Environment samples next state*
  - *Agent receives reward  $r_t$  and next state  $s_{t+1}$*

- A policy  $\pi$  is a function from  $S$  to  $A$  that specifies what action to take in each state

- **Objective**: find policy  $\pi^*$  that maximizes cumulative discounted reward:  $\sum_{t \geq 0} \gamma^t r_t$







42 SLSMO Module 11: Sequence modeling and reinforcement learning

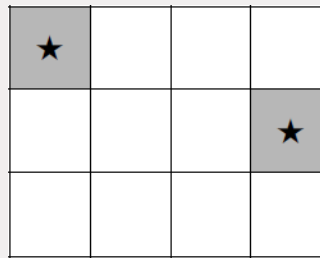
TU/e

## Reinforcement learning

### A simple MDP: Grid World

Actions = {

1. Right 
  2. Left 
  3. Up 
  4. Down 
- }



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

**Objective:** reach one of the terminal states in the least number of actions

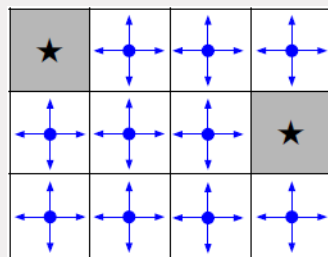


43 5LSM0 Module 11: Sequence modeling and reinforcement learning

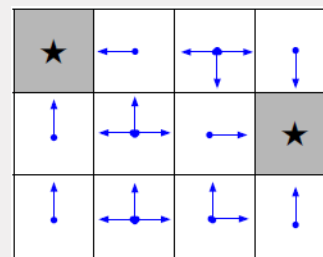
TU/e

## Reinforcement learning

### A simple MDP: Grid World



Random policy



Optimal policy



44 5LSM0 Module 11: Sequence modeling and reinforcement learning

TU/e

## Reinforcement learning

### Want to find the optimal policy $\pi^*$

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

→ Maximize the **expected sum of the rewards!**

- Formally:

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid \pi \right] \quad \text{with} \quad s_0 \sim p(s_0), a_t \sim \pi(\cdot \mid s_t), s_{t+1} \sim \mathbb{P}(\cdot \mid s_t, a_t)$$



## Reinforcement learning

### Definitions: value function & Q-value function

#### How good is a state?

- The **value function** at state  $s$ , is the expected cumulative reward from following the policy from that state  $s$ :

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

#### 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:

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



## Reinforcement learning

### Bellman equation

- The optimal  $Q$ -value function  $Q^*$  is the maximum expected cumulative reward achievable from a given state-action pair:

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

- $Q^*$  satisfies the **Bellman equation**:

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

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

The optimal policy  $\pi^*$  corresponds to taking the best action in any state specified by  $Q^*$



## Reinforcement learning

### Value iteration algorithm: use Bellman equation as an iterative update

$$Q_{i+1}^*(s, a) = \mathbb{E} \left[ r + \gamma \max_{a'} Q_i^*(s', a') \mid s, a \right]$$

- $Q_i$  will converge to  $Q^*$  as  $i \rightarrow \infty$

### Problem: **this is not scalable:**

- We must compute  $Q(s, a)$  for every state-action pair.
- If state is e.g. current game state pixels, this is computationally infeasible to compute for the entire state space!

Deep Q-learning

**Solution: use a function approximator to estimate  $Q(s, a)$ , e.g. a neural network!**





## Reinforcement learning

We want to find a Q-function that satisfies the Bellman equation

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

**Forward pass**

- Loss function

$$L_i(\theta_i) = \mathbb{E}_{s, a \sim \mathcal{R}} \left[ (y_i - Q(s, a; \theta_i))^2 \right]$$

where

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

**Backward pass**

Gradient update (with respect to Q-function parameters  $\theta$ )

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

Iteratively try to make the Q-value close to the target value  $y_i$  it should have

If Q-function corresponds to optimal  $Q^*$  (and optimal policy  $\pi^*$ )



## Reinforcement learning

Case study: 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

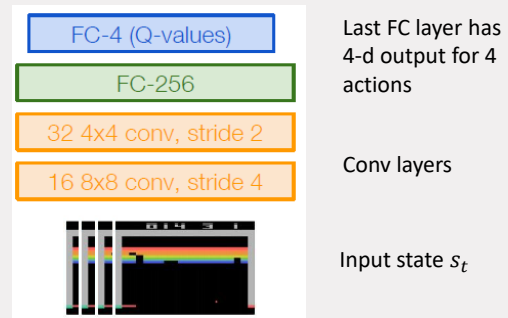


## Reinforcement learning

### Case Study: Playing Atari Games

- $Q(s, a; \theta)$  neural network with weights  $\theta$
- Current state  $s_t$ :  $84 \times 84 \times 4$  stack of last 4 frames
  - (after RGB->grayscale conversion, downsampling, and cropping)

A single feedforward pass to compute Q-values for all actions from the current state → efficient!



## Reinforcement learning

### Learning from batches of consecutive samples is problematic:

- Samples are correlated → inefficient learning
- Current Q-network parameters determines next training samples
  - (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand side) → can lead to bad feedback loops

### Address these problems using **experience replay**

- Continually update a replay memory table of transitions  $(s_t, a_t, r_t, s_{t+1})$  as game (experience) episodes are played
- Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples



[Mnih et al. NIPS Workshop 2013; Nature 2015]

## Reinforcement learning

---

**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  $\mathcal{D}$ 
    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
  
```

Play M episodes (full games)


For each time-step of the game

With small probability, select a random action, otherwise select greedy action from current policy


Take the action  $a_t$ , and observe the reward  $r_t$  and next state  $s_t$

Store transition in replay memory

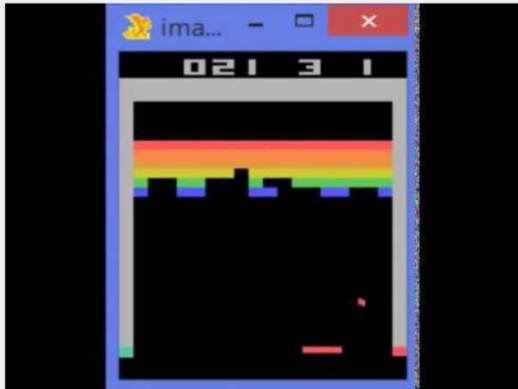
Experience Replay: Sample a random minibatch of transitions from replay memory and perform a gradient descent step




53 SLSMO Module 11: Sequence modeling and reinforcement learning




## Reinforcement learning



<https://www.youtube.com/watch?v=V1eYniJ0Rnk>



54 SLSMO Module 11: Sequence modeling and reinforcement learning



# Reinforcement learning

## Summary

### Agent in an environment

- The agent takes actions and receives rewards from the environment
- The environment also controls the state of the agent

### Q-function evaluates the value of a certain state-action pair

- Infeasible to compute...
- We can use neural networks to approximate the Q-function!

