



INTERNSHIP PROGRESS

Modeling curriculum learning



TABLE OF CONTENTS

3	WEEKLY OBJECTIVES
4	LAST WEEK'S OBJECTIVES
5	CURRENT STATE
10	QUESTIONS
13	THIS WEEK'S OBJECTIVES

WEEKLY OBJECTIVES

	FEBRUARY	MARCH	APRIL	MAY
W1	/	Level 2: accuracy RL Level 2: learning progress RL	Unify the Level 1 models Account for catastrophic interference	Finishing touches Written report
W2	/	Create hypotheses and model ideas for Level 2 & 3	Level 3: integrate both accuracy and learning progress	Written report
W3	Level 1: implement tasks and neural network	Prepare the presentation	Level 3: explore more options (chaining effect? Between-task learning?)	(exams)
W4	Finish 3 models for Level 1 Level 2: accuracy RL learning	LAB PRESENTATION Integrate all suggestions	Level 3: adjust and compare model performances	(exams)



Past week



Upcoming week

LAST WEEK'S OBJECTIVES



LEARN TO TRACK MODEL PARAMETERS

Such as accuracy, loss, performance



②

START LEVEL 2: REINFORCEMENT LEARNING

Research teacher level networks, and implement them using accuracy

④

UPLOAD CODE TO GITHUB

Improve understanding of the current model, search for potential improvements.



③

SUMMARISE FINDINGS

Write a short summary of findings, questions and plan for next week.



⑤

① IMPLEMENT ALL 3 TASKS

Easy - Hard - Impossible
Give them the same structure




```
n_input, n_output = train_x.shape[1], 1
```

```
model = Sequential()  
model.add(Dense(  
    units=4, activation='tanh', input_shape=(n_input,)))  
model.add(Dense(  
    n_output, activation='sigmoid'))  
model.build()
```

```
model.compile(  
    optimizer = tf.keras.optimizers.Adam(learning_rate=0.1),  
    loss = tf.keras.losses.BinaryCrossentropy(),  
    metrics = ['accuracy'])
```

```
##Y_DATA SET + FIT
```

```
###AND task
```

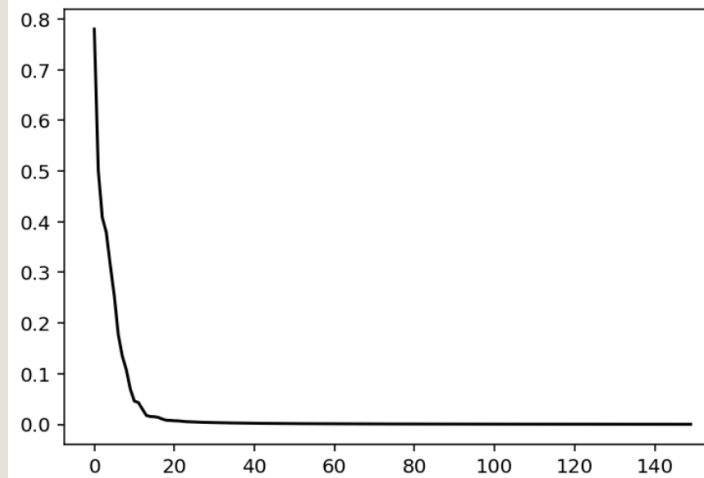
```
if model_name == 'AND_model':  
    train_y = np.array([0, 0, 0, 1])  
    train_y = train_y.reshape(4, 1)  
    history = model.fit(train_x, train_y, batch_size = 1, epochs=epochs)  
    loss_history = history.history["loss"]
```

```
###XOR task
```

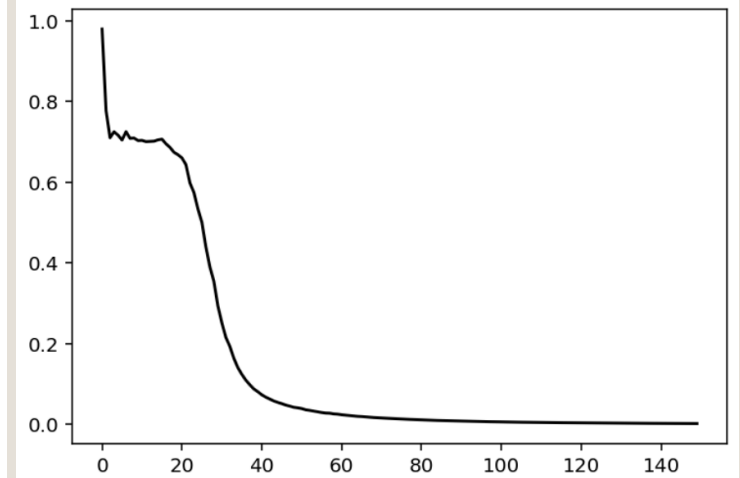
```
elif model_name == 'XOR_model':  
    train_y = np.array([0, 1, 1, 0])  
    train_y = train_y.reshape(4, 1)  
    history = model.fit(train_x, train_y, batch_size = 1, epochs=epochs)  
    loss_history = history.history["loss"]
```

```
###RM loop
```

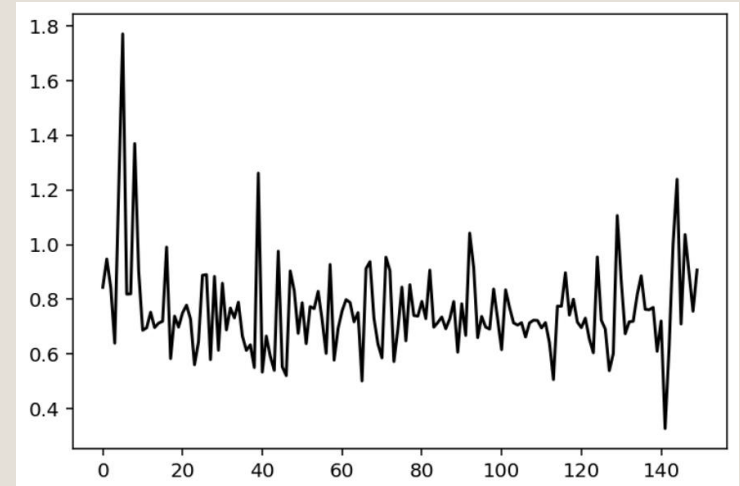
```
elif model_name == 'RM_model':  
    loss_history = []  
    accuracy_history = []  
  
    for epoch in range(epochs):  
        train_y = np.random.randint(0, 2, size=(4, 1))  
        print(f"Epoch {epoch+1}/{epochs}")  
        history = model.fit(train_x, train_y, batch_size = 1, epochs=1, verbose=1)  
        loss_history.append(history.history["loss"][0])  
        accuracy_history.append(history.history["accuracy"][0])
```



a) AND condition training (epochs/loss)



b) XOR condition training (epochs/loss)



c) random mapping condition training (epochs/loss)

QUESTIONS

QUESTION 1

- Q value network (value-based method) or policy network (policy-based methods)? And other variables...

QUESTION 2

- Which learning equation to use?

NOTES

- I need help to check my understanding of reinforcement learning (see next page)

Several variables affect the rules that shape reinforcement learning.

	Effect on reinforcement learning	Advantages and disadvantages
Value-based vs Policy-based	Value-based: algorithm stores and updates values through a value function Policy-based: algorithm learns and updates a policy	Value-based: efficient in discrete spaces, but struggles in more complex spaces, harder to adapt. Policy-based: complex spaces, stochastic (exploration and uncertainty), flexibility, requires more data.
Model-based vs Model-free	Model-based: models the environment through a transition and reward model. Model-free: uses estimates and updates based on prediction error.	Model-based: faster learning, better for limited data Model-free: better for rich data and uncertainty, flexibility
Sampling type	Bootstrapping: update estimates based on estimates of future values. Monte Carlo sampling: only updates after actual returns (full episodes).	Bootstrapping: faster learning but higher bias Monte Carlo sampling: slower learning but lower bias.
On-policy vs off-policy	...	

There are two steps in reinforcement learning: action choice and learning

Base: Rescorla-Wagner rule

$$V_t [\text{option}] = V_{t-1} [\text{option}] + \alpha (R_{t-1} - V_{t-1} [\text{option}])$$

	LEARNING	ACTION CHOICE
VALUE	Temporal Difference rule (Q) Incorporate temporal differences when multiple t exist. $Q_{\pi}(s, a) = E [\sum_{k=0}^{+\infty} \gamma^k R_{t+k+1} s, a, \pi]$ <i>π = policy</i> <i>γ = discount</i>	Highest Q-value?
POLICY	Temporal Difference rule (V) $V_{\pi}(s) = E [\sum_{k=0}^{+\infty} \gamma^k R_{t+k+1} s, \pi]$	Softmax function: calculating values based on expectations of value (probabilities) $P_t[\text{option} == A] = \frac{\exp(\beta v_t[\text{option} == A])}{\sum_i \exp(\beta v_t[\text{option} == i])}$

Solution: the Bellman equation

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

RL modelling concepts

NEXT WEEK'S OBJECTIVES



- ① **CREATE A RL MODEL**
Using a basic n-bandit task and a Rescorla-Wagner algorithm

INTEGRATE LVL 2
Using accuracy

②

INTEGRATE LVL 2
Using learning progress

④

...

...

③

**SUMMARISE
FINDINGS**

Write a short summary of findings, questions and plan for next week.

⑤