

Reminder I

You should be enrolled in the private session we created in Coursera for CMPUT 365.

I cannot use marks from the public repository for your course marks.

You **need** to **check**, **every time**, if you are in the private session and if you are submitting quizzes and assignments to the private section.

The deadlines in the public session do not align with the deadlines in Coursera.

If you have any questions or concerns, **talk with the TAs** or email us cmput365@ualberta.ca.

Reminders and Notes

- The practice quiz and programming assignment for "Constructing features for prediction" are due on Friday.
- Exam viewing for Midterm 2:
 - Tuesday (Nov 18): 15:00 17:00 @ UCOMM-3-480
 - Wednesday (Nov 19): 10:00 12:00 @ UCOMM-3-162
 - Thursday (Nov 20): 14:00 16:00 @ UCOMM-3-162
- The remaining worksheets are now all available on Canvas.
- The answers for the worksheet exercises that are not from the textbook are now all available on Canvas as well.
- I'll be in CMPUT 495 Honors Seminar today at 3 p.m.
 - o I'll talk about agent-state construction and neuroscience

Reminders and Notes

- The practice quiz and programming assignment for "Constructing features for prediction" are due on Friday.
- Exam viewing for Midterm 2:
 - Tuesday (Nov 18): 15:00 17:00 @ UCOMM-3-480
 - → Wednesday (Nov 19): 10:00 12:00 @ UCOMM-3-162
 - Thursday (Nov 20): 14:00 16:00 @ UCOMM-3-162
- The remaining worksheets and the answers for all worksheet exercises that are not from the textbook are now available on Canvas.
- Unfortunately, I can't stay after class today.

Reminders and Notes

- The practice quiz and programming assignment for "Constructing features for prediction" are due on Friday.
- Exam viewing for Midterm 2:
 - Tuesday (Nov 18): 15:00 17:00 @ UCOMM 3-480
 - → Wednesday (Nov 19): 10:00 12:00 @ UCOMM-3-162
 - Thursday (Nov 20): 14:00 16:00 @ UCOMM-3-162
- The remaining worksheets and the answers for all worksheet exercises that are not from the textbook are now available on Canvas.

Please, interrupt me at any time!



Two weeks ago: Semi-gradient TD with Function Approximation

```
Semi-gradient TD(0) for estimating \hat{v} \approx v_{\pi}
Input: the policy \pi to be evaluated What approximation should we use?!
Input: a differentiable function \hat{v}: \mathbb{S}^+ \times \mathbb{R}^d \to \mathbb{R} such that \hat{v}(\text{terminal},\cdot) = 0
Algorithm parameter: step size \alpha > 0
Initialize value-function weights \mathbf{w} \in \mathbb{R}^d arbitrarily (e.g., \mathbf{w} = \mathbf{0})
Loop for each episode:
    Initialize S
    Loop for each step of episode:
        Choose A \sim \pi(\cdot|S)
        Take action A, observe R, S'
        \mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})
        S \leftarrow S'
    until S is terminal
```

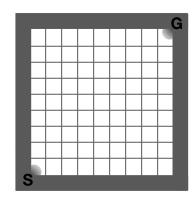


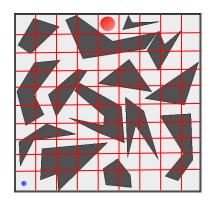
Feature Construction for Linear Methods

- Linear methods can be effective, but they heavily rely on how states are represented in terms of features.
- Feature construction is a way of adding domain knowledge; but at the same time, it went out of fashion because of deep reinforcement learning.
- Naïve linear function approximation methods do not take into consideration the interaction between features.

State Aggregation

- Simplest form of representation
- States are grouped together (one component of the vector w) for each group.
- State aggregation is a special case of SGD in which the gradient, ∇v̂(S_t,w_t), is 1 for S_t's group's component and 0 for the other components.





• Doesn't work so well, but they are one of the simplest families of features.

- Doesn't work so well, but they are one of the simplest families of features.
- Suppose an RL problem has states with two numerical dimensions.

$$\mathbf{x}(s) = (s_1, s_2)^{\top}$$

- Doesn't work so well, but they are one of the simplest families of features.
- Suppose an RL problem has states with two numerical dimensions.

$$\mathbf{x}(s) = (s_1, s_2)^{\top}$$

But what about interactions? What if both features were zero?

$$\mathbf{x}(s) = (1, s_1, s_2, s_1 s_2)^{\top}$$

- Doesn't work so well, but they are one of the simplest families of features.
- Suppose an RL problem has states with two numerical dimensions.

$$\mathbf{x}(s) = (s_1, s_2)^{\top}$$

But what about interactions? What if both features were zero?

$$\mathbf{x}(s) = (1, s_1, s_2, s_1 s_2)^{\top}$$

And we can keep going...

$$\mathbf{x}(s) = (1, s_1, s_2, s_1 s_2, s_1^2, s_2^2, s_1 s_2^2, s_1^2 s_2, s_1^2 s_2^2)^{\top}$$

Suppose each state s corresponds to k numbers, $s_1, s_2, ..., s_k$, with each $s_i \in \mathbb{R}$. For this k-dimensional state space, each order-n polynomial-basis feature x_i can be written as

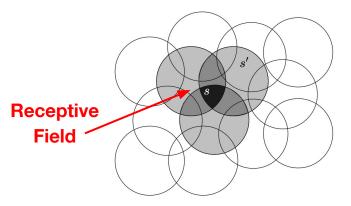
$$x_i(s) = \prod_{j=1}^k s_j^{c_{i,j}}, \tag{9.17}$$

where each $c_{i,j}$ is an integer in the set $\{0,1,\ldots,n\}$ for an integer $n \geq 0$. These features make up the order-n polynomial basis for dimension k, which contains $(n+1)^k$ different features.

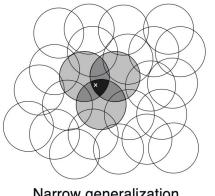


Coarse Coding

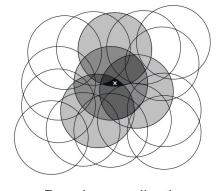
- Consider a task in which the natural representation of the state set is a continuous two-dimensional space.
- We define binary features indicating whether a state is present or not in a specific circle.



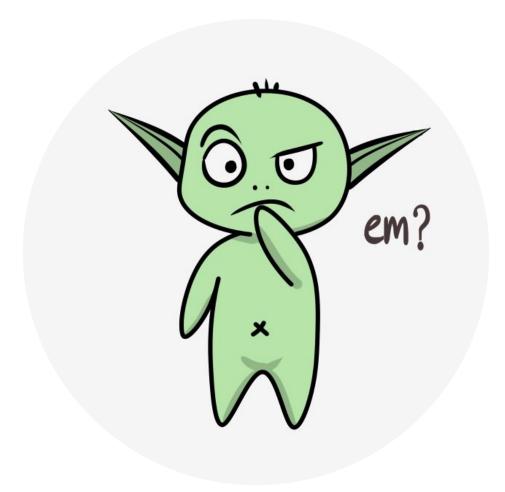
The shape defines generalization





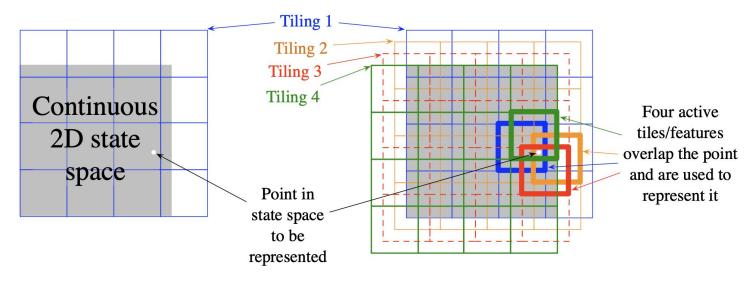


Broad generalization

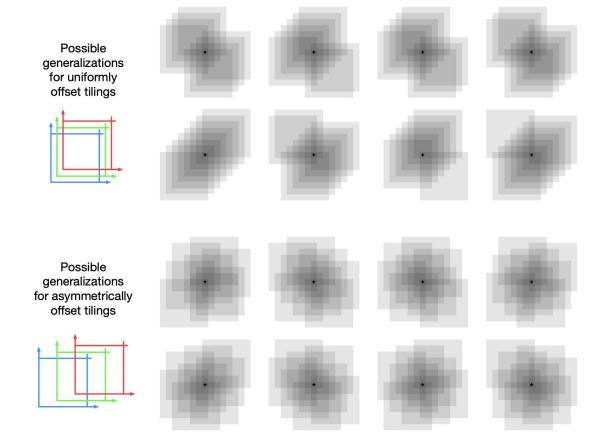


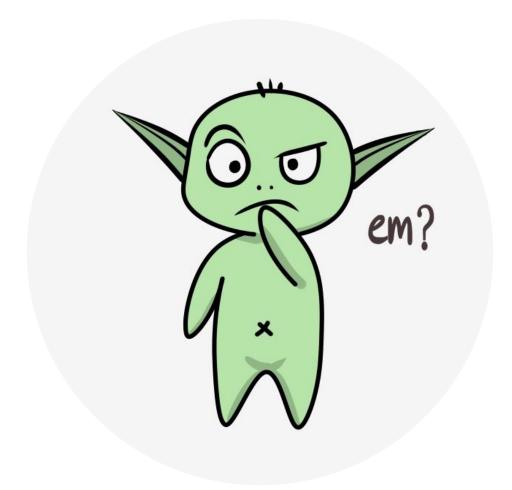
Tile Coding

• Tile coding is a form of coarse coding for multi-dimensional continuous spaces (with a fixed number of active features per timestep).



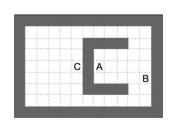
Tile Coding

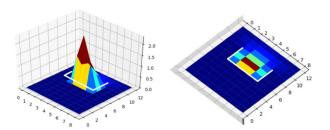




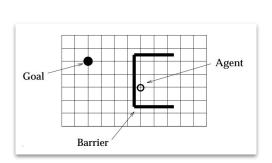
It Isn't that We do Function Approximation Because We Cannot do Tabular Reinforcement Learning

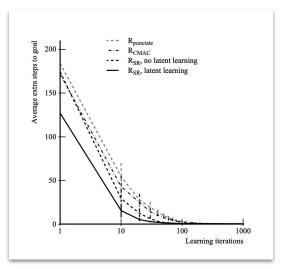
Successor Representation [Dayan, Neural Computation 1993].

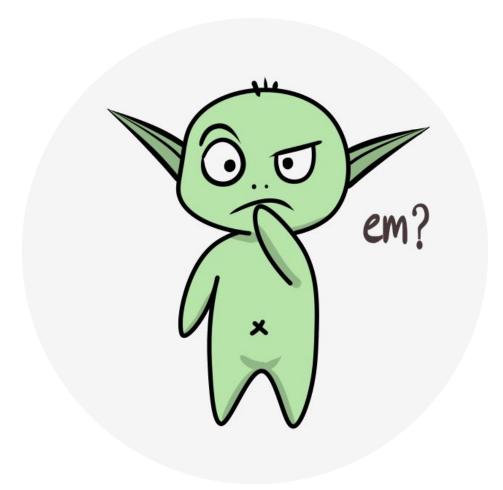








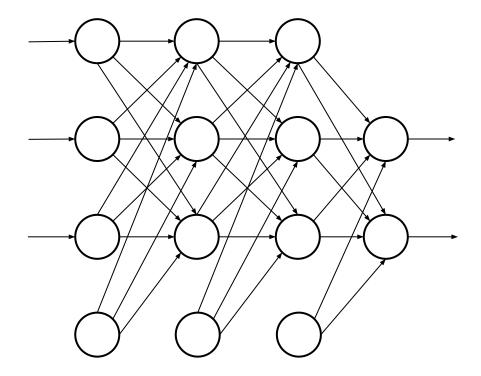


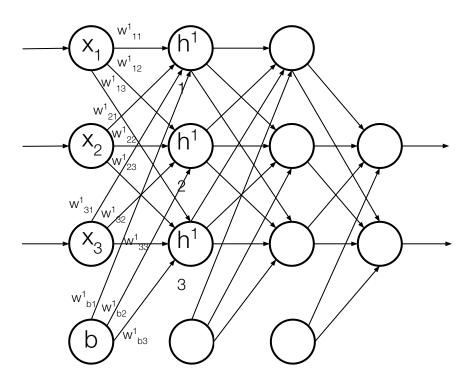


Nonlinear Function Approximation: Artificial Neural Networks

- The basics of deep reinforcement learning.
- Idea: Instead of using linear features, we feed the "raw" input to a neural network and ask it to predict the state (or state-action) value function.

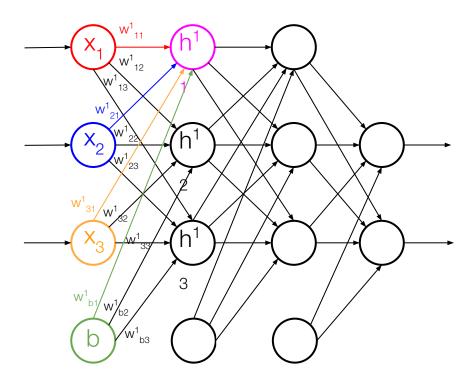






$$\mathbf{h^1} = \operatorname{act}(\mathbf{xW^1})$$

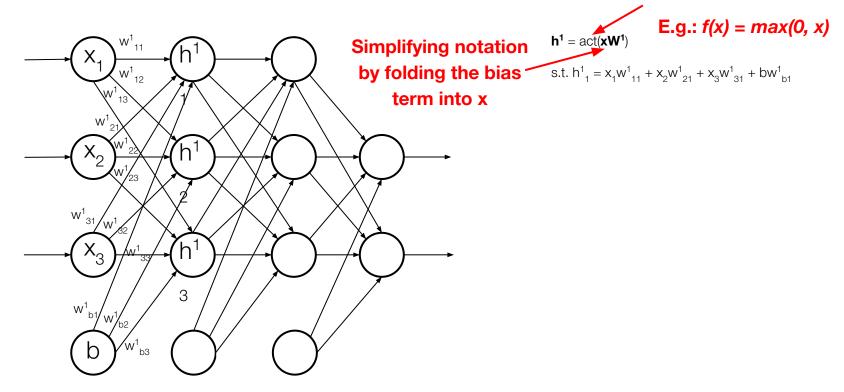
s.t.
$$h_1^1 = x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b w_{b1}^1$$



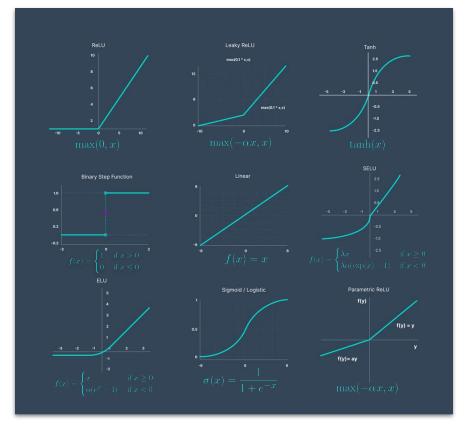
$$\boldsymbol{h^1} = \text{act}(\boldsymbol{xW^1})$$

s.t.
$$h_1^1 = X_1 w_{11}^1 + X_2 w_{21}^1 + X_3 w_{31}^1 + bw_{b1}^1$$

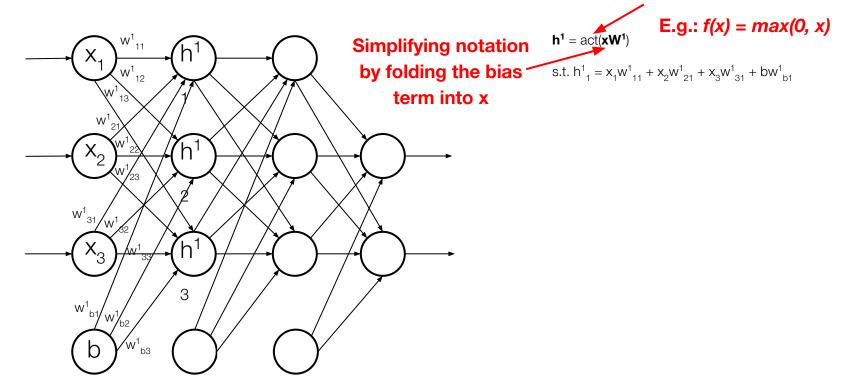
The activation function introduces non-linearity

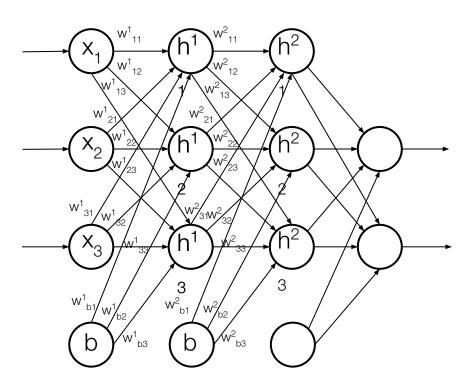


Main Types of Activation Functions



The activation function introduces non-linearity



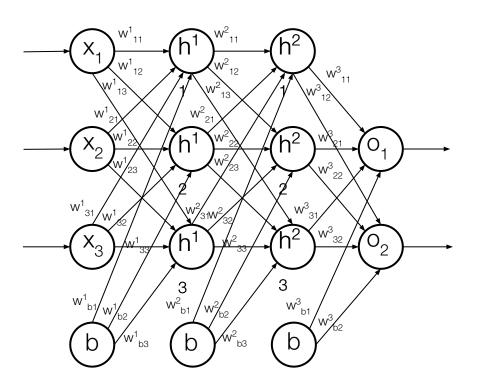


$$\mathbf{h^1} = \operatorname{act}(\mathbf{xW^1})$$

s.t.
$$h_1^1 = x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b w_{b1}^1$$

$$h^2 = act(h^1W^2)$$

s.t.
$$h_1^2 = h_1^1 w_{11}^2 + h_2^1 w_{21}^2 + h_3^1 w_{31}^2 + b w_{b1}^2$$



$$\mathbf{h^1} = \operatorname{act}(\mathbf{xW^1})$$

s.t.
$$h_1^1 = x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b w_{b1}^1$$

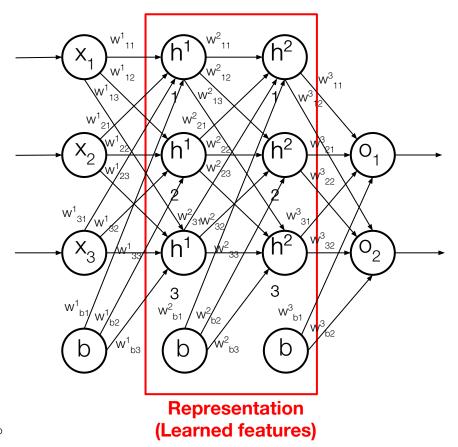
$$\mathbf{h^2} = \operatorname{act}(\mathbf{h^1W^2})$$

s.t.
$$h_{1}^{2} = h_{1}^{1}w_{11}^{2} + h_{2}^{1}w_{21}^{2} + h_{3}^{1}w_{31}^{2} + bw_{b1}^{2}$$

$$\mathbf{o} = \operatorname{act}(\mathbf{h}^2\mathbf{W}^3)$$

s.t.
$$o_1 = h_{11}^2 w_{11}^3 + h_{22}^2 w_{21}^3 + h_{33}^2 w_{31}^3 + b w_{b1}^3$$

 $o = act(act(act(xW^1)W^2)W^3)$



 $\mathbf{h^1} = \operatorname{act}(\mathbf{xW^1})$

s.t. $h_1^1 = x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b w_{b1}^1$

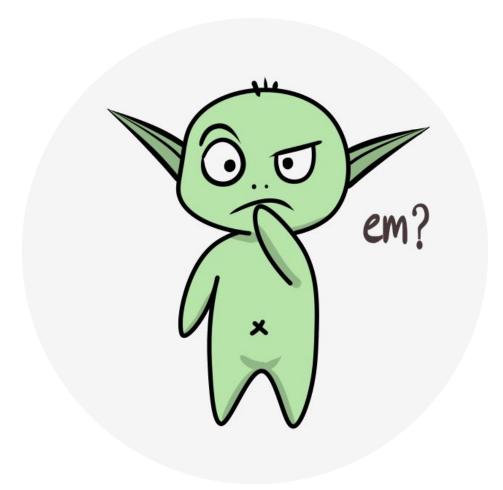
 $\mathbf{h^2} = \operatorname{act}(\mathbf{h^1W^2})$

s.t. $h_{1}^{2} = h_{1}^{1}w_{11}^{2} + h_{2}^{1}w_{21}^{2} + h_{3}^{1}w_{31}^{2} + bw_{b1}^{2}$

 $\mathbf{o} = \operatorname{act}(\mathbf{h^2W^3})$

s.t. $o_1 = h_1^2 w_{11}^3 + h_2^2 w_{21}^3 + h_3^2 w_{31}^3 + b w_{b1}^3$

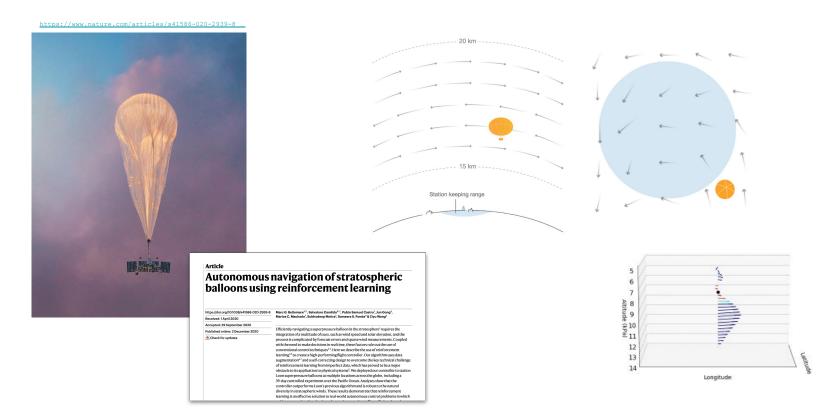
 $o = act(act(act(xW^1)W^2)W^3)$



A Note from the Textbook

The backpropagation algorithm can produce good results for shallow networks having 1 or 2 hidden layers, but it may not work well for deeper ANNs. In fact, training a network with k+1 hidden layers can actually result in poorer performance than training a network with k hidden layers, even though the deeper network can represent all the functions that the shallower network can (Bengio, 2009). Explaining results like these is not easy, but several factors are important. First, the large number of weights in a typical deep ANN makes it difficult to avoid the problem of overfitting, that is, the problem of failing to generalize correctly to cases on which the network has not been trained. Second, backpropagation does not work well for deep ANNs because the partial derivatives computed by its backward passes either decay rapidly toward the input side of the network, making learning by deep layers extremely slow, or the partial derivatives grow rapidly toward the input side of the network, making learning unstable. Methods

What we feed to a neural network still matters (a lot!)





Deep Convolutional Network

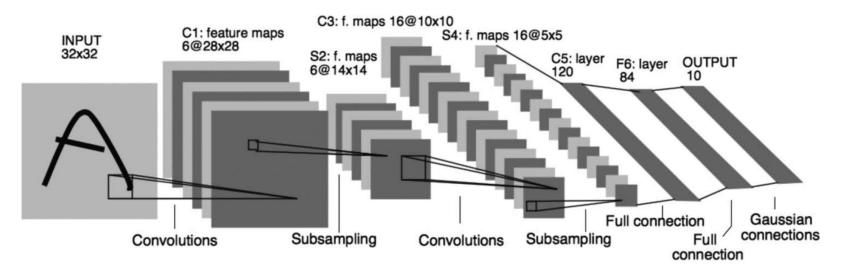
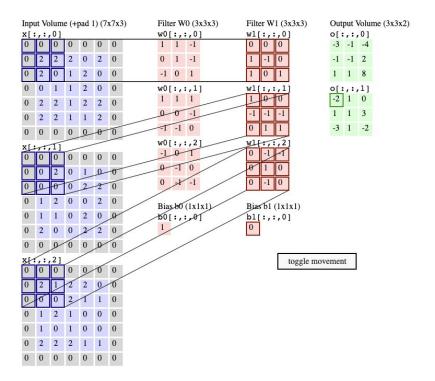
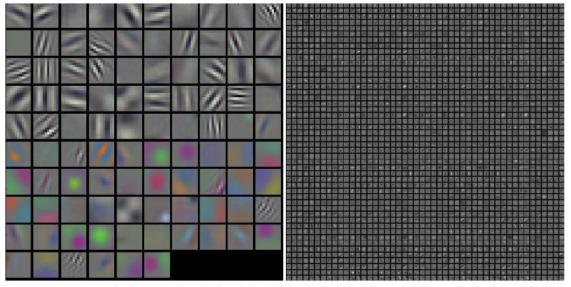


Figure 9.15: Deep Convolutional Network. Republished with permission of Proceedings of the IEEE, from Gradient-based learning applied to document recognition, LeCun, Bottou, Bengio, and Haffner, volume 86, 1998; permission conveyed through Copyright Clearance Center, Inc.

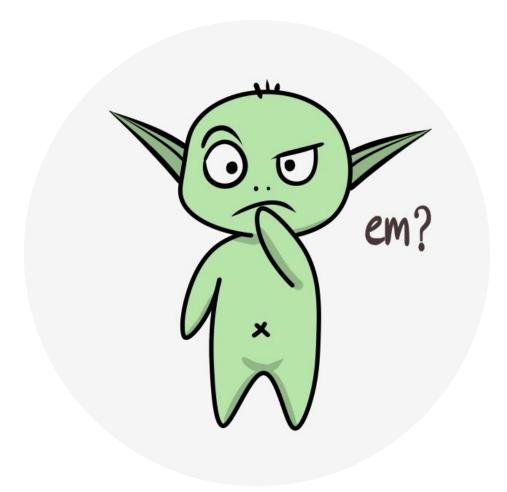
Deep Convolutional Network



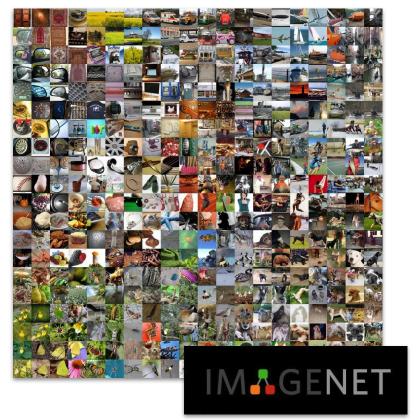
Learned Representations



Typical-looking filters on the first CONV layer (left), and the 2nd CONV layer (right) of a trained AlexNet. Notice that the first-layer weights are very nice and smooth, indicating nicely converged network. The color/grayscale features are clustered because the AlexNet contains two separate streams of processing, and an apparent consequence of this architecture is that one stream develops high-frequency grayscale features and the other low-frequency color features. The 2nd CONV layer weights are not as interpretable, but it is apparent that they are still smooth, well-formed, and absent of noisy patterns.



Are we improving anything? From features to architecture...



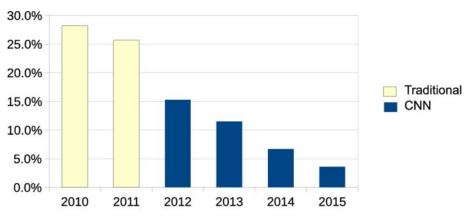


Image by Bottou, Curtis, and Nocedal (2016)



Journal of Artificial Intelligence Research 47 (2013) 253–279

Submitted 02/13; published 06/13

The Arcade Learning Environment: An Evaluation Platform for General Agents

Marc G. Bellemare

MG17@CS.UALBERTA.CA

University of Alberta, Edmonton, Alberta, Canada

Yavar Naddaf

YAVAR@EMPIRICALRESULTS.CA

Empirical Results Inc., Vancouver, British Columbia, Canada

Joel Veness Michael Bowling VENESS@CS.UALBERTA.CA BOWLING@CS.UALBERTA.CA

University of Alberta, Edmonton, Alberta, Canada

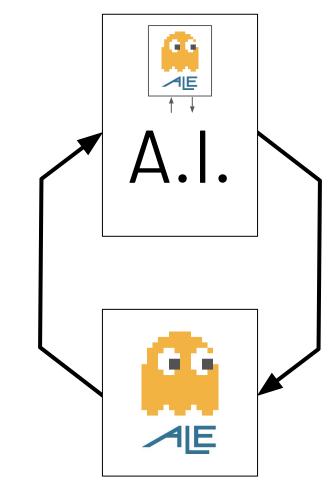
Abstract

In this article we introduce the Arcade Learning Environment (ALE): both a challenge problem and a platform and methodology for evaluating the development of general, domain-independent AI technology. ALE provides an interface to hundreds of Atari 2600 game environments, each one different, interesting, and designed to be a challenge for human players. ALE presents significant research challenges for reinforcement learning, model learning, model-based planning, imitation learning, transfer learning, and intrinsic motivation. Most importantly, it provides a rigorous testbed for evaluating and comparing approaches to these problems. We illustrate the promise of ALE by developing and benchmarking domain-independent agents designed using well-established AI techniques for both reinforcement learning and planning. In doing so, we also propose an evaluation

Arcade Learning Environment

Over 50 Domains in 8 Minutes 23 Seconds

Reinforcement Learning 2000 Copyright 1982



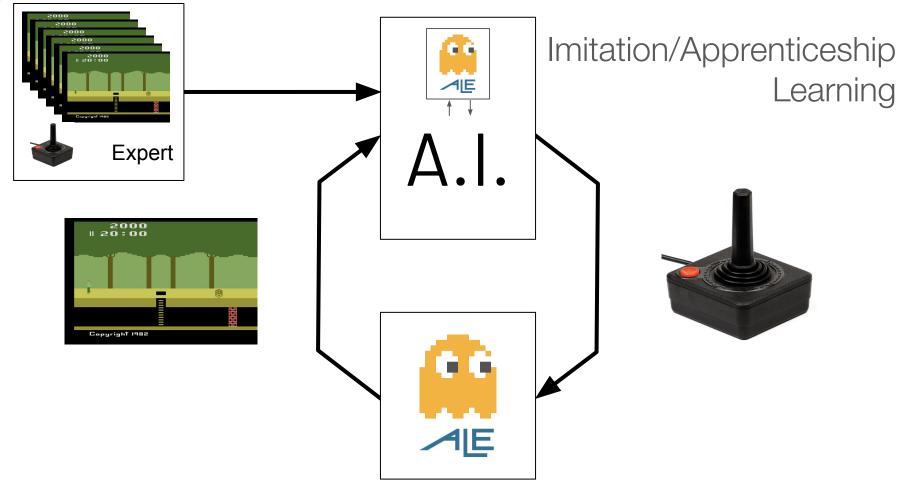
Model Learning



00680

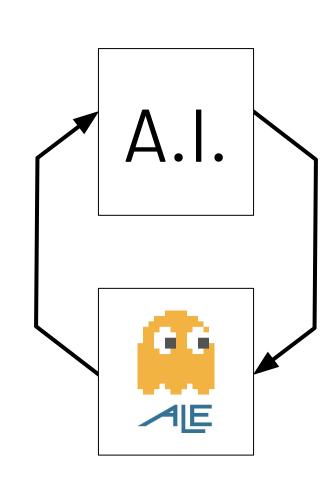
2000

Copyright 1982



5854 || 18:03 ACTIVISTON 2000 Copyright 1982 1899 || 19:38

Ç. AŒIVSION







2000 Copyright 1982 1692

Transfer Learning

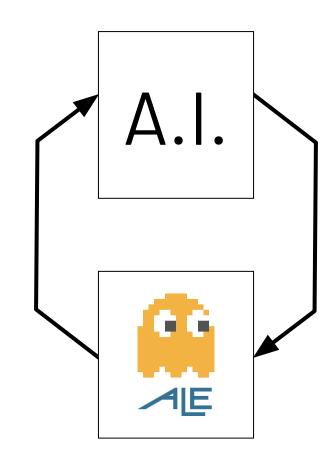


Pitfall II

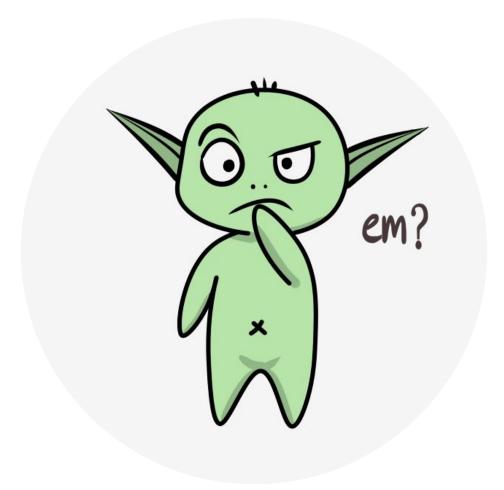
Pitfall!

Intrinsic Motivation



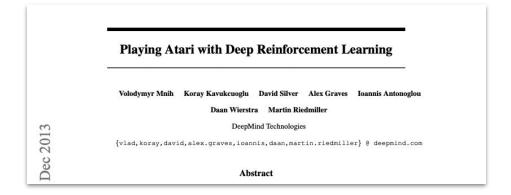


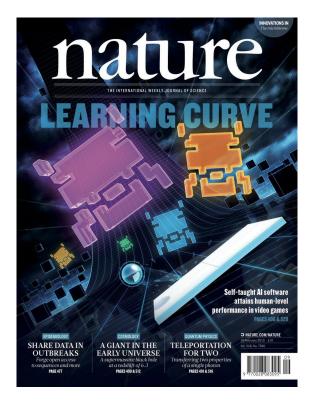




Deep Q-Network (and Deep RL)

[Mnih et al., 2013, 2015]

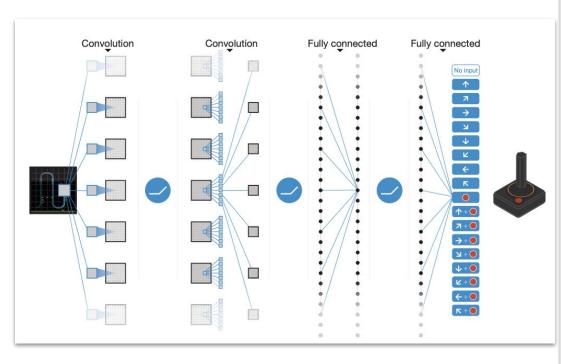


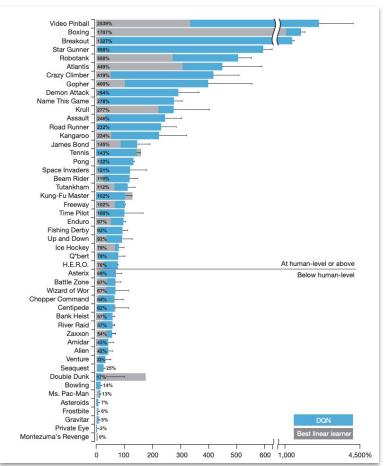




Deep Q-Network (and Deep RL)

[Mnih et al., 2013, 2015]







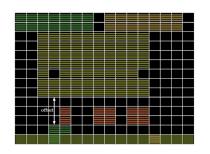
What if we were to design features instead?

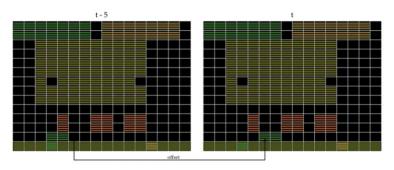
What if we were to design features instead?

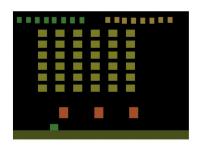
State of the Art Control of Atari Games Using Shallow Reinforcement Learning

Yitao Liang[†], Marlos C. Machado[‡], Erik Talvitie[†], and Michael Bowling[‡] †Franklin & Marshall College Lancaster, PA, USA {yliang, erik.talvitie}@fandm.edu

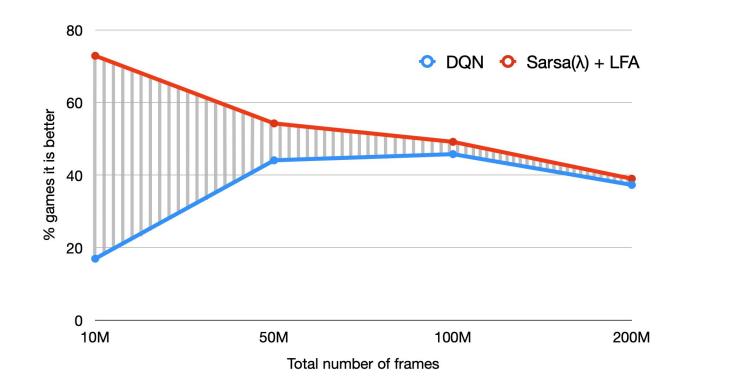
[‡]University of Alberta Edmonton, AB, Canada {machado, mbowling}@ualberta.ca

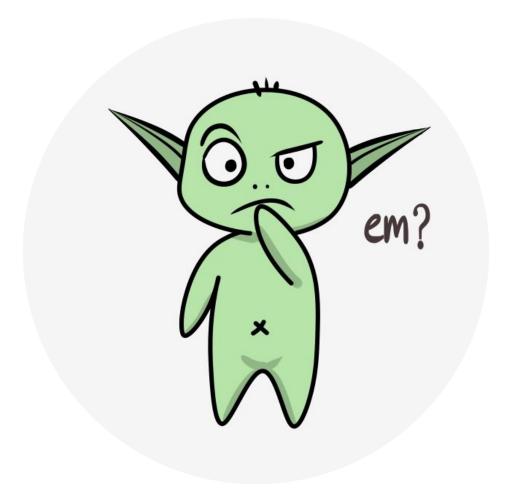




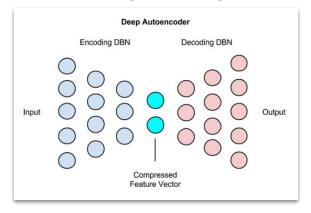


There are many trade-offs, we need to understand them [Liang et al., 2016; Machado et al. 2018]





There are many many inductive biases one can use



A number of convolutional layers

Fully connected layers

https://wiki.pathmind.com/deep-autoencoder

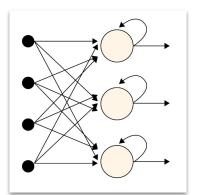


Image by Yu, Miao, and Wang (2022)

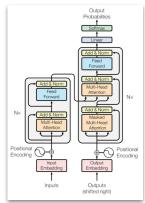


Image by Vaswani et al. (2017)

Marlos C. Machado

https://www.scaler.com/topics/deep-learning/rnn/