# Summary on Value Function Approximation

**Christopher Mutschler**

# Value Function Approximation

- In practice we often need to approximate the value function, because:
  - (Explicit) tabular representations require too much space
  - We want to generalize information across state (see also: POMDPs!)
- For linear function approximation almost all convergence guarantees hold
  - For non-linear function approximation such guarantees cannot be given
  - But careful scheduling and several tricks help to stabilize training
- But:
  - Non-linear function approximation is very sensitive to hyper-parameter tuning!
  - See also: https://www.youtube.com/watch?v=Vh4H0gOwdlg (not directly related but definitely worth watching!)
  - And also: https://www.alexirpan.com/2018/02/14/rl-hard.html (but please read with humor)

# The Deadly Triad

- Stability in RL is a very serious thing!

- Instability and divergence in RL mainly stem from

  1. Function Approximation.

  2. Bootstrapping.

  3. Off-policy training.

- Unfortunately, in most of the case we really use the full combination.

# Fuzzy Tiling Activations



(a) TA, $k = 4$     (b) FTA, $k = 4, \eta = 0.1$     (c) FTA, $k = 4, \eta = 0.25$

- DQNs need target networks to reduce the chance of divergence

- Main reason:
  - The Q-Targets are non-stationary and moving
  - over subsequent gradient descent steps

- Idea:
  - Introduce a special activation function that produces sparse representations! i.e., updates on a particular Q-value does no affect nearby Q-values that much.
  - FTA layers stack $k$-dimensional sparse encodings for each element $h_1 W_2$ ($h_1 = x W_1$)
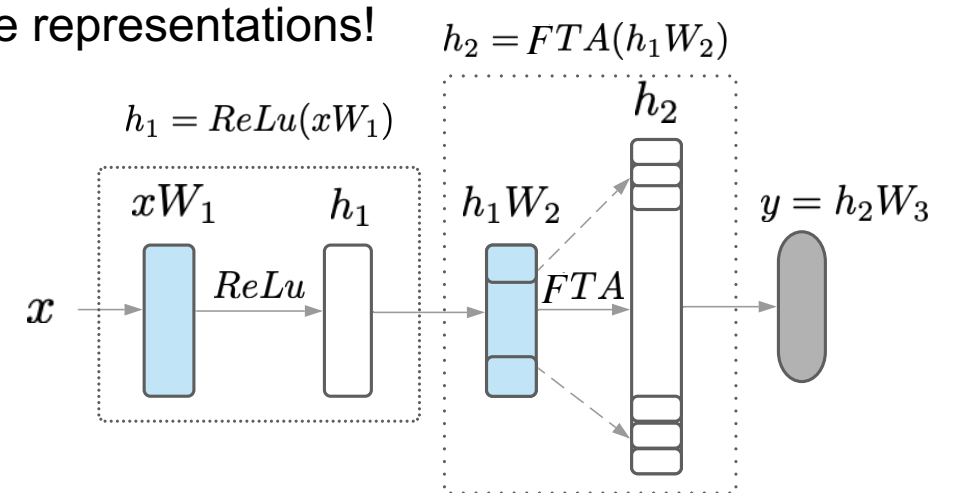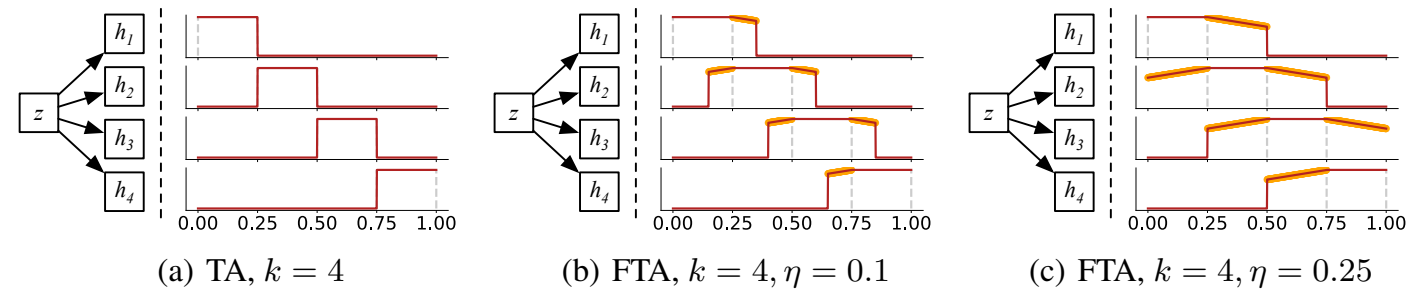


Figure 2: A visualization of an FTA layer

*Pan et al.: Fuzzy Tiling Activations: A Simple Approach to Learning Sparse Representations Online. ICLR 2021.*

# Lesson of today

**"Be careful with (non-linear) function approximation"**

# References

**Books:**

- Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.

**Lectures:**

- Pieter Abbeel: CS 188 Introduction to Artificial Intelligence. Fall 2018
- UCL Course on RL. http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html

**Web:**

- https://medium.com/init27-labs/understanding-q-learning-the-cliff-walking-problem-80198921abbc
- https://deepmind.com/blog/alphastar-mastering-real-time-strategy-game-starcraft-ii/
- https://danieltakeshi.github.io/2016/10/31/going-deeper-into-reinforcement-learning-understanding-q-learning-and-linear-function-approximation/
- http://www0.cs.ucl.ac.uk/staff/d.silver/web/Resources_files/deep_rl.pdf