Introductory remarks for RL seminar

The aim of this series of seminars is to understand AlphaGo, which is an important advance not only in the field of artificial intelligence but in science more broadly. To understand AlphaGo we need to understand its three components:

 Monte	e - Carlo t	tree sec	irch	(Le	ctures	7,8)	
 Deep	learning			(Lea	ctures	ر4 را	5,6)
•	0	N 1		1		- (~	\sim)

- Reinforcement learning (Lectures 2,6,8,9)

Why is AlphaGo important?

That depends who you are. It is not clear AlphaGo is important in pure mathematics, except perhaps for logic (see the interview with Szegedy linked on the seminar webpage). However, for applied mathematics, physics and other natural sciences it seems obviously important, for the same reason that calculus and the study of polynomial approximation has been important for the last several centuries.

- <u>Theorem</u> (Stone-Weierstrass) If $X \subseteq \mathbb{R}^n$ is compact, then polynomial functions are dense in $Ct_s(X, \mathbb{R})$. (compact-open to pology, 1-e. sup-metric)
 - <u>Banach + Picard</u>: solve ODEs by iteration of a contraction mapping to find a fixed point, gives polynomial approximations.
 - This is <u>actually use ful</u> in many cases!

But there are functions in nature for which this approach completely fails:

Example Consider the state space of the game
$$G_0$$
 on a 19×19 board:
 $S = \{empty, while, black\}^{[1,...,19] \times \{1,...,19\}} |S| = 3^{361}$
 $S_{legal} \subseteq S |S_{legal}| \sim 2.082 \times 10^{170}$
Let V^* : $S_{legol} \rightarrow \mathbb{R}$ dende the optimal value function, given
rewards + 1 for win, -1 for a loss, and a discount factor $T < 1$
(1.e. the expected value of time discounted rewavel for a "perfect" player).
(that this makes sense will be justified in Lecture 2), in a game of Go.
To find V^* by fixed point iteration V_0, V_1, \ldots would need to store

$$2 \times 10^{10} \times 16 \text{ bits} = 4 \times 10^{100} \text{ bytes}$$

$$16 \text{ bit real number} = 4 \times 10^{152} \times 10^{18} \text{ bytes}$$

$$= 4 \times 10^{152} \text{ exabytes}$$

In 2011 world disk storage amounted to 295 exabytes. So maybe you don't want to do that.

We can extend V to a continuous function (say by linear interpolation)

$$\begin{bmatrix} 0,1 \end{bmatrix}^{3\cdot361} \cong \left(\begin{bmatrix} 0,1 \end{bmatrix}^3 \right)^{\{1,\dots,19\}\times\{1,\dots,19\}} \longrightarrow \mathbb{R}$$

$$(1,0,0), (0,1,0), (0,0,1)$$

$$(1,0,0), (0,1,0), (0,0,1)$$

$$(0,0,0), (0,0,0)$$

$$(0,0,0), (0,0,0)$$

Although this optimal value function <u>exists</u> as a mathematical object, it was far from clear that there existed any reasonable algorithm for producing a sequence of simple functions $(f_n)_n \gg 0$ converging to V^*



The work of the research group at DeepMind led by David Silver shows that cleep neural networks and stochastic gradient descent + RL + Monte-Carlo tree search + self-play suffices.

<u>Theorem</u> If $X \subseteq \mathbb{R}^n$ is compact then feedforward ReLU neural networks of width n+1 are dense in $Ct_3(X, \mathbb{R})$.