

CS287 Advanced Robotics
Lecture 4 (Fall 2019)

Function Approximation

Pieter Abbeel
UC Berkeley EECS

Value Iteration

Algorithm:

Start with $V_0^*(s) = 0$ for all s .

For $i = 1, \dots, H$

For all states s in S :

$$V_{i+1}^*(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i^*(s')]$$

$$\pi_{i+1}^*(s) \leftarrow \arg \max_{a \in A} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i^*(s')]$$

This is called a **value update** or **Bellman update/back-up**

Impractical for
large state spaces

$V_i^*(s)$ = expected sum of rewards accumulated starting from state s , acting optimally for i steps

$\pi_i^*(s)$ = optimal action when in state s and getting to act for i steps

Similar issue for policy iteration and linear programming

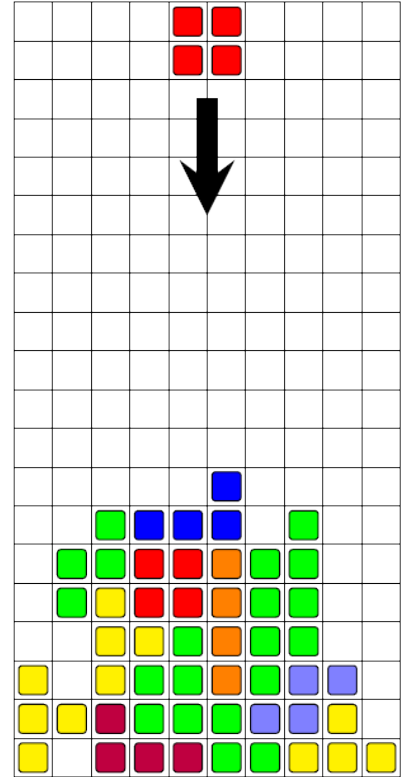
Outline

- Function approximation
- Value iteration with function approximation
- Policy iteration with function approximation
- Linear programming with function approximation

Function Approximation Example 1 : Tetris

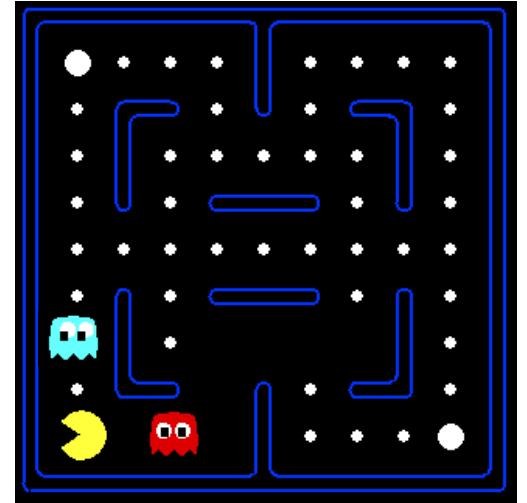
- state: board configuration + shape of the falling piece $\sim 2^{200}$ states!
- action: rotation and translation applied to the falling piece
- 22 features aka basis functions ϕ_i
 - Ten basis functions, $0, \dots, 9$, mapping the state to the height $h[k]$ of each column.
 - Nine basis functions, $10, \dots, 18$, each mapping the state to the absolute difference between heights of successive columns: $|h[k+1] - h[k]|$, $k = 1, \dots, 9$.
 - One basis function, 19, that maps state to the maximum column height: $\max_k h[k]$
 - One basis function, 20, that maps state to the number of 'holes' in the board.
 - One basis function, 21, that is equal to 1 in every state.

$$\hat{V}_\theta(s) = \sum_{i=0}^{21} \theta_i \phi_i(s) = \theta^\top \phi(s)$$



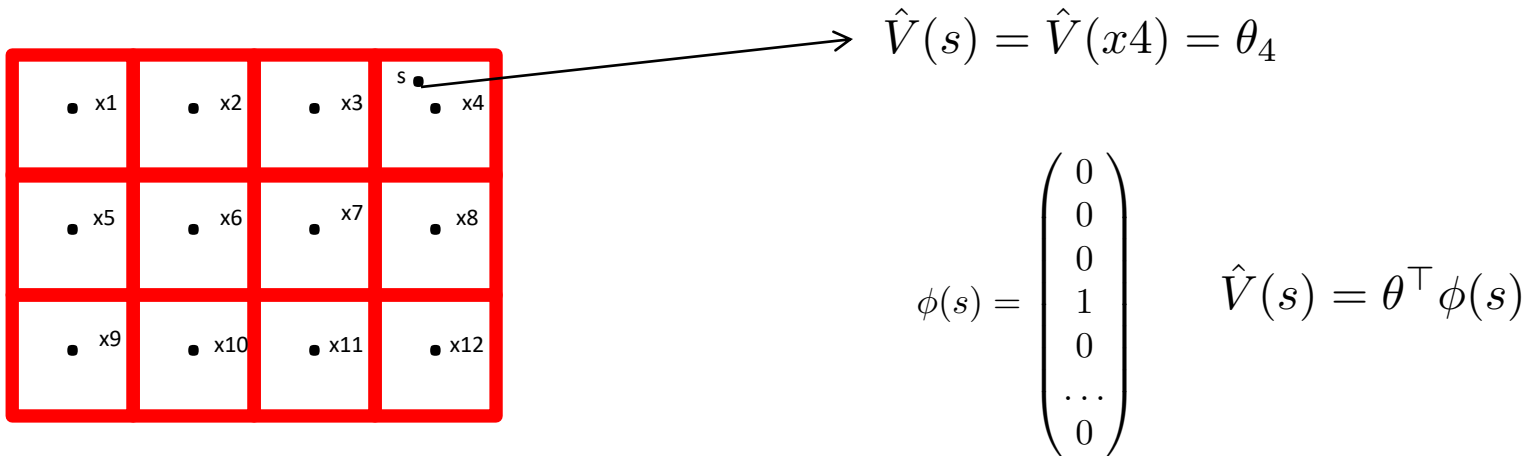
Function Approximation Example 2: Pacman

$$\begin{aligned} V(s) = & \theta_0 \\ & + \theta_1 \text{“distance to closest ghost”} \\ & + \theta_2 \text{“distance to closest power pellet”} \\ & + \theta_3 \text{“in dead-end”} \\ & + \theta_4 \text{“closer to power pellet than ghost”} \\ & + \dots \\ = & \sum_{i=0}^n \theta_i \phi_i(s) = \theta^\top \phi(s) \end{aligned}$$



Function Approximation Example 3: Nearest Neighbor

- 0'th order approximation (1-nearest neighbor):



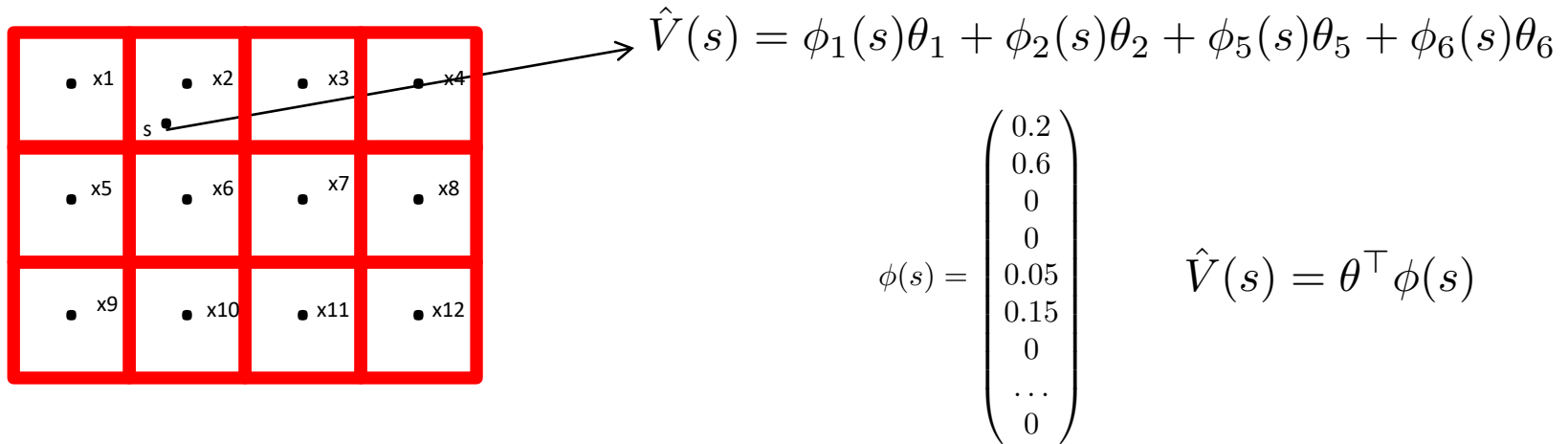
Only store values for x_1, x_2, \dots, x_{12}

– call these values $\theta_1, \theta_2, \dots, \theta_{12}$

Assign other states value of nearest “x” state

Function Approximation Example 4: k-Nearest Neighbor

- 1'th order approximation (k-nearest neighbor interpolation):



Only store values for x_1, x_2, \dots, x_{12}

– call these values $\theta_1, \theta_2, \dots, \theta_{12}$

Assign other states interpolated value of nearest 4 “x” states

More Function Approximation Examples

- Examples:

- $S = \mathbb{R}, \quad \hat{V}(s) = \theta_1 + \theta_2 s$

- $S = \mathbb{R}, \quad \hat{V}(s) = \theta_1 + \theta_2 s + \theta_3 s^2$

- $S = \mathbb{R}, \quad \hat{V}(s) = \sum_{i=0}^n \theta_i s^i$

- $S = \mathbb{R}^n \quad \hat{V}(s) = f_{\theta}(s) \quad (\text{e.g. neural net})$

Function Approximation

- Main idea:

- Use approximation \hat{V}_θ of the true value function V^*

- θ is a free parameter to be chosen from its domain Θ

- Representation size: $|S|$ down to $|\Theta|$

+ : less parameters to estimate

- : less expressiveness,

because typically there exist many V^* for which there is no θ such that $\hat{V}_\theta = V^*$

Supervised Learning

- Given:

- set of examples $(s^{(1)}, V(s^{(1)})), (s^{(2)}, V(s^{(2)})), \dots, (s^{(m)}, V(s^{(m)}))$,

- Asked for:

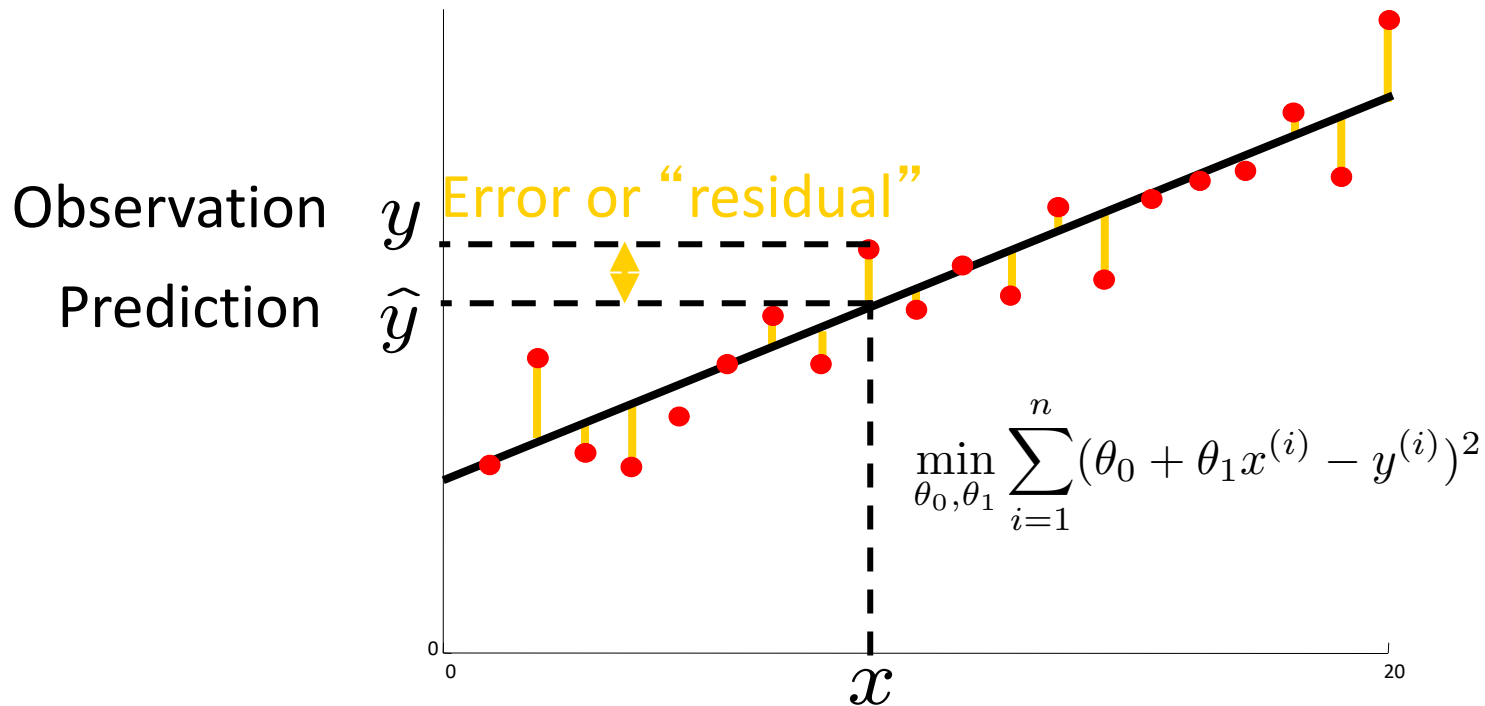
- “best” \hat{V}_θ

- Representative approach: find θ through least squares

$$\min_{\theta \in \Theta} \sum_{i=1}^m (\hat{V}_\theta(s^{(i)}) - V(s^{(i)}))^2$$

Supervised Learning Example

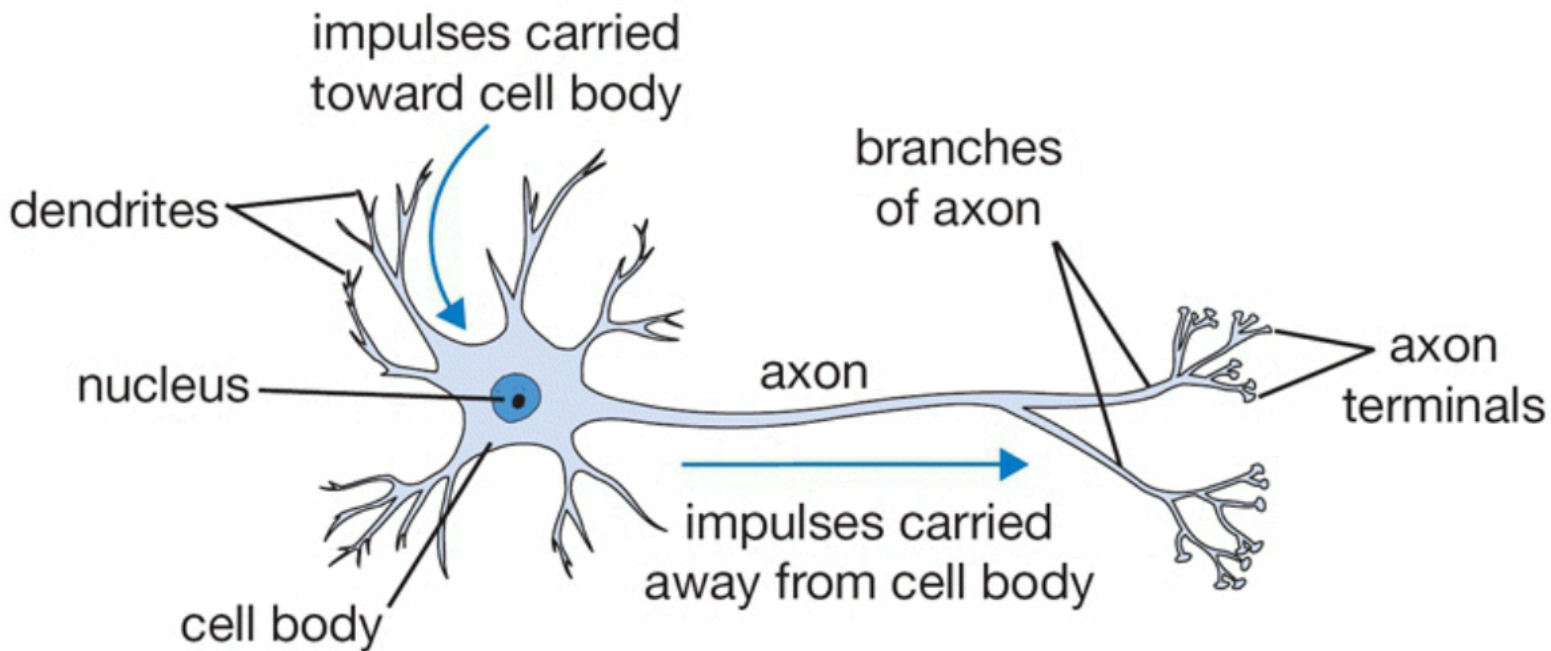
- Linear regression



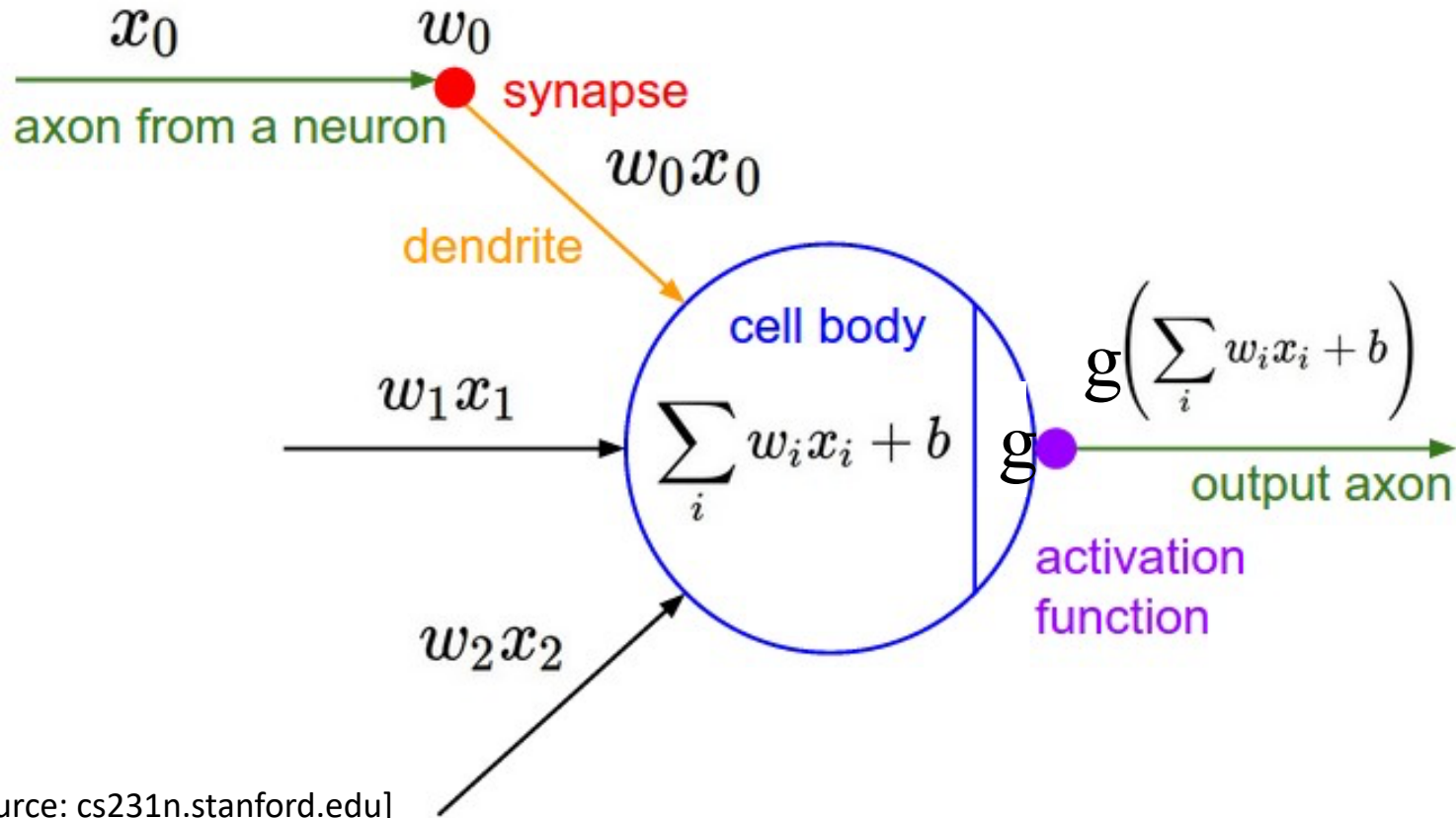
Supervised Learning Example

- Neural Nets

Single (Biological) Neuron

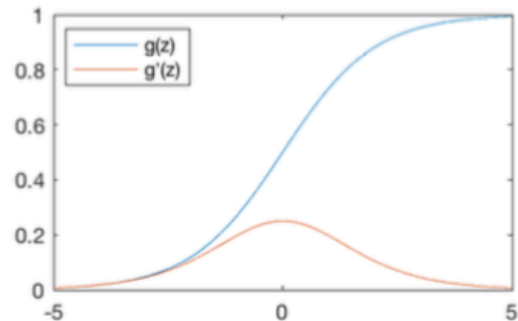


Single (Artificial) Neuron



Common Activation Functions

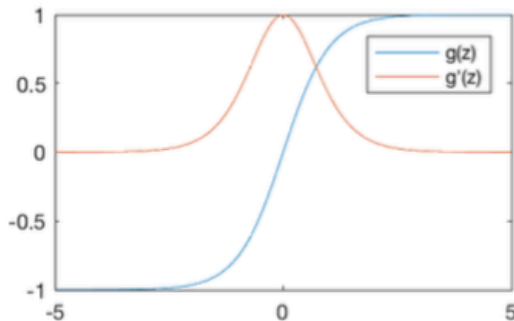
Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

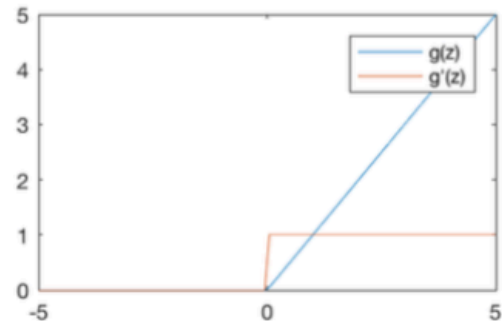
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

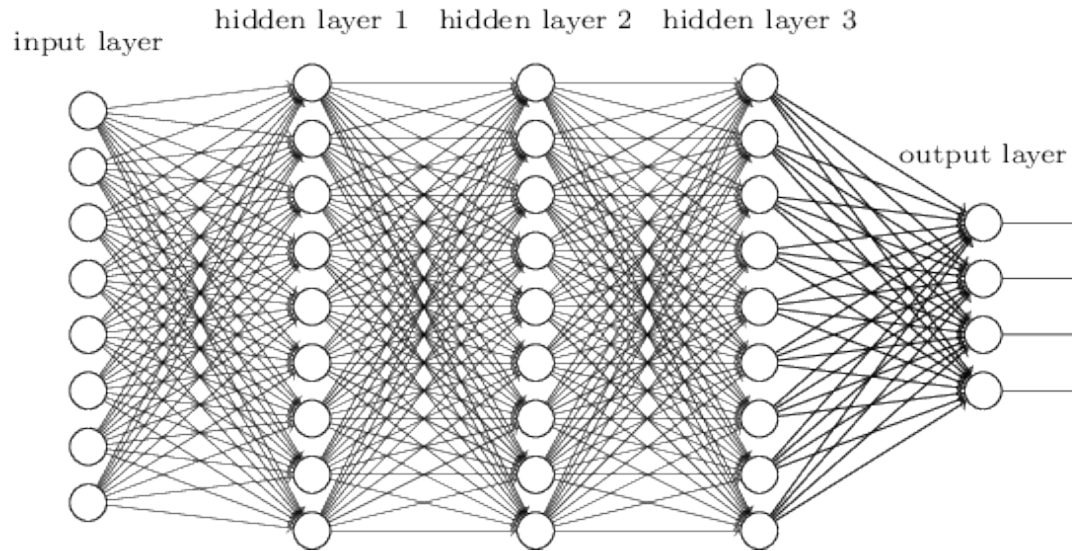
Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

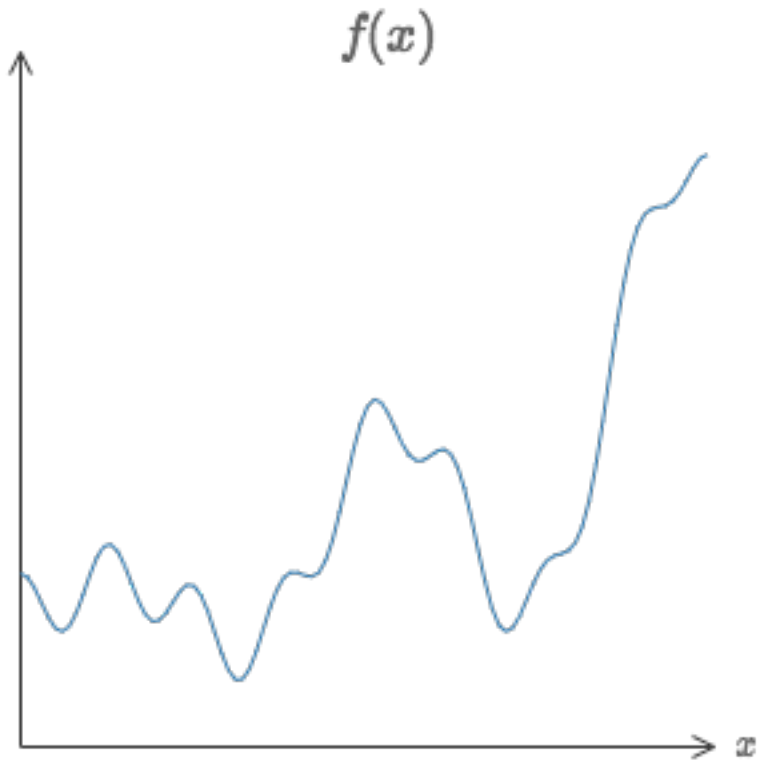
Neural Network



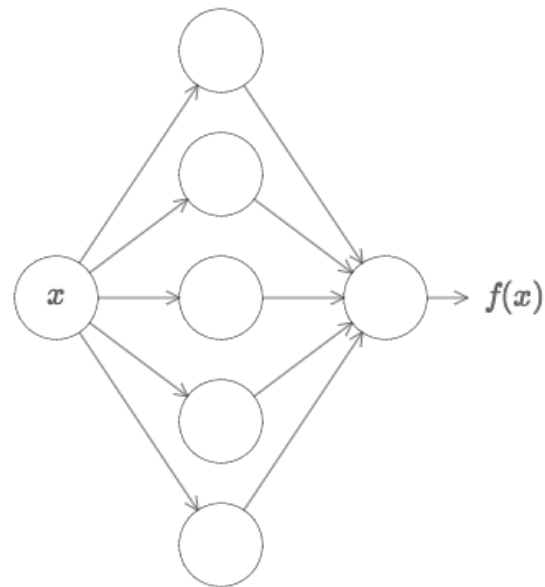
Notation: x $z^{(1)}$ $z^{(2)}$ $z^{(3)}$ $y = f(x, w)$

Choice of w determines the function from $x \rightarrow y$

What Functions Can a Neural Net Represent?



Does there exist a choice for w to make this work?



Universal Function Approximation Theorem

Hornik theorem 1: Whenever the activation function is *bounded and nonconstant*, then, for any finite measure μ , standard multilayer feedforward networks can approximate any function in $L^p(\mu)$ (the space of all functions on R^k such that $\int_{R^k} |f(x)|^p d\mu(x) < \infty$) arbitrarily well, provided that sufficiently many hidden units are available.

Hornik theorem 2: Whenever the activation function is *continuous, bounded and non-constant*, then, for arbitrary compact subsets $X \subseteq R^k$, standard multilayer feedforward networks can approximate any continuous function on X arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

- In words: Given any continuous function $f(x)$, if a 2-layer neural network has enough hidden units, then there is a choice of weights that allow it to closely approximate $f(x)$.

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"

Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"

Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

Universal Function Approximation Theorem

Math. Control Signals Systems (1989) 2: 303-314

Mathematics of Control,
Signals, and Systems
© 1989 Springer-Verlag New York Inc.

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functions can uniformly approximate any continuous function of a real variable with support in the unit hypercube, only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

Key words. Neural networks, Approximation, Completeness.

1. Introduction

A number of diverse application areas are concerned with the representation of general functions of an n -dimensional real variable, $x \in \mathbb{R}^n$, by finite linear combinations of the form

$$\sum_{j=1}^N \alpha_j \sigma(y_j^T x + \theta_j), \quad (1)$$

where $y_j \in \mathbb{R}^n$ and $\alpha_j, \theta_j \in \mathbb{R}$ are fixed, (y_j^T) is the transpose of y_j so that $y_j^T x$ is the inner product of y_j and x . Here the univariate function σ depends heavily on the context of the application. Our major concern is with so-called sigmoidal σ 's:

$$\sigma(t) \rightarrow \begin{cases} 1 & \text{as } t \rightarrow +\infty, \\ 0 & \text{as } t \rightarrow -\infty. \end{cases}$$

Such functions arise naturally in neural network theory as the activation function of a neural node (or unit as is becoming the preferred term) [L1], [RHM]. The main result of this paper is a demonstration of the fact that sums of the form (1) are dense in the space of continuous functions on the unit cube if σ is any continuous sigmoidal

* Date received: October 21, 1988. Date revised: February 17, 1989. This research was supported in part by NSF Grant DCR-8619103, ONR Contract N000-86-G-0202 and DOE Grant DE-FG02-85ER25001.

† Center for Supercomputing Research and Development and Department of Electrical and Computer Engineering, University of Illinois, Urbana, Illinois 61801, U.S.A.

Neural Networks, Vol. 4, pp. 251-257, 1991
Printed in the USA. All rights reserved.

0893-6003/91 \$3.00 + .00
Copyright © 1991 Pergamon Press plc

ORIGINAL CONTRIBUTION

Approximation Capabilities of Multilayer Feedforward Networks

KURT HORNIK

Technische Universität Wien, Vienna, Austria

(Received 30 January 1990; revised and accepted 25 October 1990)

Abstract. We show that standard multilayer feedforward networks with as few as a single hidden layer and arbitrary bounded and nonconstant activation function are universal approximators with respect to $L(p)$ performance criteria, for arbitrary finite input environment measures μ , provided only that sufficiently many hidden units are available. If the activation function is continuous, bounded and nonconstant, then continuous mappings can be learned uniformly over compact input sets. We also give very general conditions ensuring that networks with sufficiently smooth activation functions are capable of arbitrarily accurate approximation to a function and its derivatives.

Keywords. Multilayer feedforward networks, Activation function, Universal approximation capabilities, Input environment measure, $L(p)$ approximation, Uniform approximation, Sobolev spaces, Smooth approximation.

1. INTRODUCTION

The approximation capabilities of neural network architectures have recently been investigated by many authors, including Carroll and Dickinson (1989), Cybenko (1989), Funahashi (1989), Gallant and White (1988), Hecht-Nielsen (1989), Hornik, Stinchcombe, and White (1989, 1990), Irie and Miyake (1988), Lapedes and Farber (1988), Stinchcombe and White (1989, 1990). (This list is by no means complete.)

If we think of the network architecture as a rule for computing values at l output units given values at k input units, hence implementing a class of mappings from \mathbb{R}^k to \mathbb{R}^l , we can ask how well arbitrary mappings from \mathbb{R}^k to \mathbb{R}^l can be approximated by the network, in particular, if as many hidden units as required for internal representation and computation may be employed.

How to measure the accuracy of approximation depends on how we measure closeness between functions, which in turn varies significantly with the specific problem to be dealt with. In many applications, it is necessary to have the network perform simultaneously well on all input samples taken from some compact input set X in \mathbb{R}^k . In this case, closeness is

measured by the uniform distance between functions on X , that is,

$$\rho_\infty(f, g) = \sup |f(x) - g(x)|.$$

In other applications, we think of the inputs as random variables and are interested in the average performance where the average is taken with respect to the input environment measure μ , where $\mu(\mathbb{R}^k) = 1$. In this case, closeness is measured by the $L(p)$ distances

$$\rho_p(f, g) = \left[\int_X |f(x) - g(x)|^p d\mu(x) \right]^{1/p}.$$

$1 \leq p < \infty$, the most popular choice being $p = 2$, corresponding to mean square error.

Of course, there are many more ways of measuring closeness of functions. In particular, in many applications, it is also necessary that the derivatives of the approximating function implemented by the network closely resemble those of the function to be approximated, up to some order. This issue was first taken up in Hornik et al. (1990), who discuss the sources of need of smooth functional approximation in more detail. Typical examples arise in robotics (learning of smooth movements) and signal processing (analysis of chaotic time series); for a recent application to problems of nonparametric inference in statistics and econometrics, see Gallant and White (1989).

All papers establishing certain approximation ca-

MULTILAYER FEEDFORWARD NETWORKS WITH NON-POLYNOMIAL ACTIVATION FUNCTIONS CAN APPROXIMATE ANY FUNCTION

by

Moshe Leshno
Faculty of Management
Tel Aviv University
Tel Aviv, Israel 69978

and

Shimon Schocken
Leonard N. Stern School of Business
New York University
New York, NY 10003

September 1991

Center for Research on Information Systems
Information Systems Department
Leonard N. Stern School of Business
New York University

Working Paper Series

STERN IS-91-26

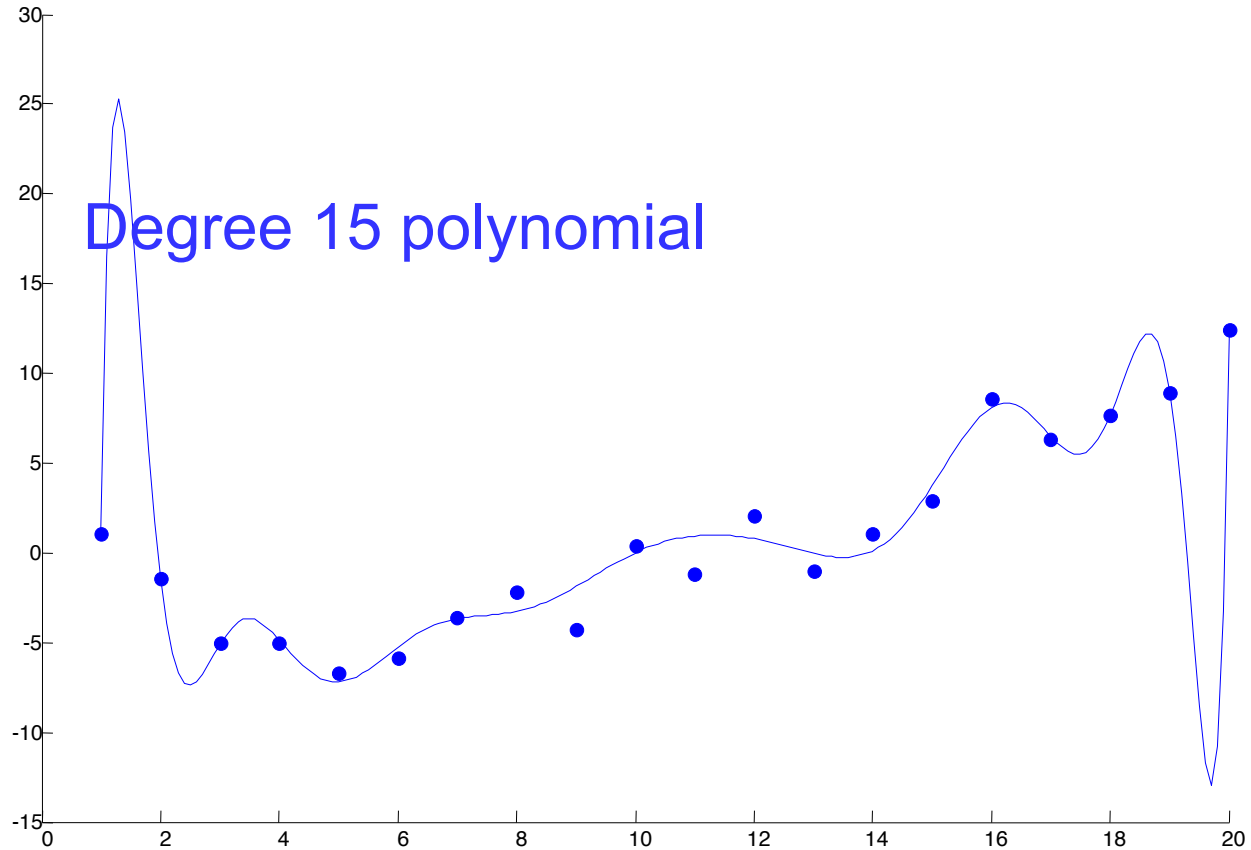
Appeared previously as Working Paper No. 21/91 at The Israel Institute Of Business Research

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"

Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"

Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

Overfitting



Avoiding Overfitting

- Reduce number of features or size of the network
- Regularize θ
- Early stopping: stop training updates once loss increases on hold-out data

Status

- Function approximation through supervised learning

BUT: where do the supervised examples come from?

Value Iteration with Function Approximation

- Initialize by choosing some setting for $\theta^{(0)}$
- Iterate for $i = 0, 1, 2, \dots, H$:
 - Step 0: Pick some $S' \subseteq S$ (typically $|S'| \ll |S|$)

- Step 1: Bellman back-ups

$$\forall s \in S' : \bar{V}_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma \hat{V}_{\theta^{(i)}}(s') \right]$$

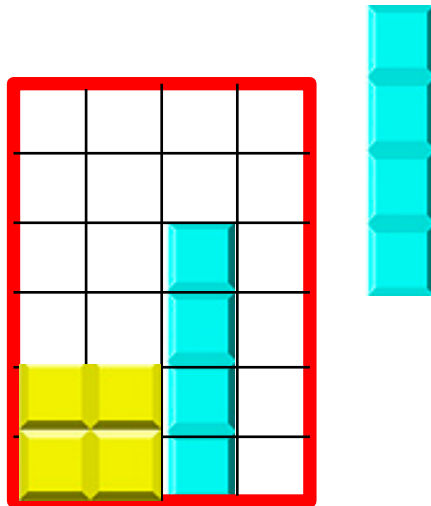
- Step 2: Supervised learning

find $\theta^{(i+1)}$ as the solution of:
$$\min_{\theta} \sum_{s \in S'} \left(\hat{V}_{\theta^{(i+1)}}(s) - \bar{V}_{i+1}(s) \right)^2$$

Value Iteration w/Function Approximation --- Example

- Mini-tetris: two types of blocks, can only choose translation (not rotation)

- Example state:



- Reward = 1 for placing a block
- Sink state / Game over is reached when block is placed such that part of it extends above the red rectangle
- If you have a complete row, it gets cleared

Value Iteration w/Function Approximation --- Example

$$S' = \left\{ \begin{array}{c} \text{Grid 1} \\ \text{Grid 2} \\ \text{Grid 3} \\ \text{Grid 4} \end{array} \right\}$$

The set S' contains four elements, each consisting of a 4x4 grid and an associated object:

- Element 1: A 4x4 grid with a red border. The bottom-left 2x2 area is yellow. The third column is cyan. A cyan vertical bar is positioned to the right of the grid.
- Element 2: A 4x4 grid with a red border. The left two columns are yellow. The third column is cyan. A cyan vertical bar is positioned to the right of the grid.
- Element 3: A 4x4 grid with a red border. The bottom-left 2x2 area is yellow. A yellow 2x2 square is positioned above and to the right of the grid.
- Element 4: A 4x4 grid with a red border. The first and third columns are cyan. A yellow 2x2 square is positioned above and to the right of the grid.

Value Iteration w/Function Approximation --- Example

$$S' = \left\{ \begin{array}{|c|c|c|c|} \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \end{array} \begin{array}{c} \color{cyan}{\square} \\ \color{cyan}{\square} \\ \color{cyan}{\square} \\ \color{cyan}{\square} \end{array} , \begin{array}{|c|c|c|c|} \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \end{array} \begin{array}{c} \color{cyan}{\square} \\ \color{cyan}{\square} \\ \color{cyan}{\square} \\ \color{cyan}{\square} \end{array} , \begin{array}{|c|c|c|c|} \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \end{array} \begin{array}{c} \color{yellow}{\square} \\ \color{yellow}{\square} \\ \color{yellow}{\square} \\ \color{yellow}{\square} \end{array} , \begin{array}{|c|c|c|c|} \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \hline \end{array} \begin{array}{c} \color{yellow}{\square} \\ \color{yellow}{\square} \\ \color{yellow}{\square} \\ \color{yellow}{\square} \end{array} \right\}$$

- 10 features (also called basis functions) φ_i
 - Four basis functions, $0, \dots, 3$, mapping the state to the height $h[k]$ of each of the four columns.
 - Three basis functions, $4, \dots, 6$, each mapping the state to the absolute difference between heights of successive columns: $|h[k+1] - h[k]|$, $k = 1, \dots, 3$.
 - One basis function, 7 , that maps state to the maximum column height: $\max_k h[k]$
 - One basis function, 8 , that maps state to the number of 'holes' in the board.
 - One basis function, 9 , that is equal to 1 in every state.
- Init with $\vartheta^{(0)} = (-1, -1, -1, -1, -2, -2, -2, -3, -2, 10)$

Value Iteration w/Function Approximation --- Example

- Bellman back-ups for the states in S' :

$$\begin{aligned}
 V(\text{state}) &= \max \{ 0.5 * (1 + \gamma V(\text{state}_1)) + 0.5 * (1 + \gamma V(\text{state}_2)), \\
 & 0.5 * (1 + \gamma V(\text{state}_3)) + 0.5 * (1 + \gamma V(\text{state}_4)), \\
 & 0.5 * (1 + \gamma V(\text{state}_5)) + 0.5 * (1 + \gamma V(\text{state}_6)), \\
 & 0.5 * (1 + \gamma V(\text{state}_7)) + 0.5 * (1 + \gamma V(\text{state}_8)) \}
 \end{aligned}$$

Value Iteration w/Function Approximation --- Example

- Bellman back-ups for the states in S' :

$$\begin{aligned}
 V(\text{state}) = \max \{ & 0.5 * (1 + \gamma V(\text{state}_1)) + 0.5 * (1 + \gamma V(\text{state}_2)), \\
 & 0.5 * (1 + \gamma V(\text{state}_3)) + 0.5 * (1 + \gamma V(\text{state}_4)), \\
 & 0.5 * (1 + \gamma V(\text{state}_5)) + 0.5 * (1 + \gamma V(\text{state}_6)), \\
 & 0.5 * (1 + \gamma V(\text{state}_7)) + 0.5 * (1 + \gamma V(\text{state}_8)) \}
 \end{aligned}$$

Value Iteration w/Function Approximation --- Example

$$S' = \left\{ \begin{array}{c} \text{4x4 grid with 1 yellow block at (1,1), 1 cyan block at (2,2), and 1 cyan block at (3,3)} \\ \text{4x4 grid with 1 yellow block at (1,1), 1 cyan block at (2,2), 1 cyan block at (3,3), and 1 cyan block at (4,3)} \\ \text{4x4 grid with 1 yellow block at (1,1), 1 cyan block at (2,2), 1 yellow block at (3,3), and 1 yellow block at (4,3)} \\ \text{4x4 grid with 1 cyan block at (1,1), 1 cyan block at (2,2), 1 yellow block at (3,3), and 1 yellow block at (4,3)} \end{array} \right\}$$

- 10 features aka basis functions φ_i
 - Four basis functions, $0, \dots, 3$, mapping the state to the height $h[k]$ of each of the four columns.
 - Three basis functions, $4, \dots, 6$, each mapping the state to the absolute difference between heights of successive columns: $|h[k+1] - h[k]|$, $k = 1, \dots, 3$.
 - One basis function, 7 , that maps state to the maximum column height: $\max_k h[k]$
 - One basis function, 8 , that maps state to the number of 'holes' in the board.
 - One basis function, 9 , that is equal to 1 in every state.
- Init with $\mathcal{V}^{(0)} = (-1, -1, -1, -1, -2, -2, -2, -3, -2, 10)$

Value Iteration w/Function Approximation --- Example

- Bellman back-ups for the states in S' :

$$\begin{aligned}
 v(\text{grid state}) &= \max \{ 0.5 * (1 + \gamma \theta^\top \phi(\text{grid state})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid state})), \\
 &\quad (6,2,4,0,4,2,4,6,0,1) \qquad (6,2,4,0,4,2,4,6,0,1) \\
 &0.5 * (1 + \gamma \theta^\top \phi(\text{grid state})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid state})), \\
 &\quad (2,6,4,0,4,2,4,6,0,1) \qquad (2,6,4,0,4,2,4,6,0,1) \\
 &0.5 * (1 + \gamma V(\text{sink-state, } V=0)) + 0.5 * (1 + \gamma V(\text{sink-state, } V=0)), \\
 &\quad (\text{sink-state, } V=0) \qquad (\text{sink-state, } V=0) \\
 &0.5 * (1 + \gamma \theta^\top \phi(\text{grid state})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid state})) \} \\
 &\quad (0,0,2,2,0,2,0,2,0,1) \qquad (0,0,2,2,0,2,0,2,0,1)
 \end{aligned}$$

Value Iteration w/Function Approximation --- Example

- Bellman back-ups for the states in S' :

$$\begin{aligned}
 v(\text{grid state}) &= \max \{ 0.5 * (1 + \gamma (-30)) + 0.5 * (1 + \gamma (-30)), \\
 &\quad 0.5 * (1 + \gamma (-30)) + 0.5 * (1 + \gamma (-30)), \\
 &\quad 0.5 * (1 + \gamma (0)) + 0.5 * (1 + \gamma (0)), \\
 &\quad 0.5 * (1 + \gamma (6)) + 0.5 * (1 + \gamma (6)) \} \\
 &= 6.4 \quad (\text{for } \gamma = 0.9)
 \end{aligned}$$

Value Iteration w/Function Approximation --- Example

$$\theta^{(0)} = (-1, -1, -1, -1, -2, -2, -2, -3, -2, 20)$$

- Bellman back-ups for the second state in S' :

$$\begin{aligned}
 V(\text{state}) &= \max \{ 0.5 * (1 + \gamma V(\text{sink-state, } V=0)) + 0.5 * (1 + \gamma V(\text{sink-state, } V=0)), \\
 & 0.5 * (1 + \gamma V(\text{sink-state, } V=0)) + 0.5 * (1 + \gamma V(\text{sink-state, } V=0)), \\
 & 0.5 * (1 + \gamma V(\text{sink-state, } V=0)) + 0.5 * (1 + \gamma V(\text{sink-state, } V=0)), \\
 & 0.5 * (1 + \gamma \theta^\top \phi(\text{state})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{state})) \} \\
 &= 19
 \end{aligned}$$

$(0,0,0,0, 0,0,0, 0, 0, 1) \rightarrow V = 20$

Value Iteration w/Function Approximation --- Example

$$\theta^{(0)} = (-1, -1, -1, -1, -2, -2, -2, -3, -2, 20)$$

- Bellman back-ups for the third state in S' :

$$\begin{aligned}
 V(\text{grid with 2 yellow squares}) &= \max \{ 0.5 * (1 + \gamma \theta^\top \phi(\text{grid with 3 yellow squares})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid with 4 yellow squares})), \\
 &\quad (4,4,0,0, 0,4,0, 4, 0, 1) \qquad \qquad \qquad (4,4,0,0, 0,4,0, 4, 0, 1) \\
 &\quad \rightarrow V = -8 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \rightarrow V = -8 \\
 &0.5 * (1 + \gamma \theta^\top \phi(\text{grid with 4 yellow squares})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid with 5 yellow squares})), \\
 &\quad (2,4,4,0, 2,0,4, 4, 0, 1) \qquad \qquad \qquad (2,4,4,0, 2,0,4, 4, 0, 1) \\
 &\quad \rightarrow V = -14 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \rightarrow V = -14 \\
 &0.5 * (1 + \gamma \theta^\top \phi(\text{grid with 6 yellow squares})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid with 7 yellow squares})) \} \\
 &\quad (0,0,0,0, 0,0,0, 0, 0, 1) \qquad \qquad \qquad (0,0,0,0, 0,0,0, 0, 0, 1) \\
 &\quad \rightarrow V = 20 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \rightarrow V = 20
 \end{aligned}$$

Value Iteration w/Function Approximation --- Example

$$\theta^{(0)} = (-1, -1, -1, -1, -2, -2, -2, -3, -2, 20)$$

- Bellman back-ups for the fourth state in S' :

$$\begin{aligned}
 V(\text{grid}) &= \max \{ 0.5 * (1 + \gamma \theta^\top \phi(\text{grid}_{(6,6,4,0,0,2,4,6,4,1)})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid}_{(6,6,4,0,0,2,4,6,4,1)})), \\
 &\quad \rightarrow V = -34 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \rightarrow V = -34 \\
 &0.5 * (1 + \gamma \theta^\top \phi(\text{grid}_{(4,6,6,0,2,0,6,6,4,1)})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid}_{(4,6,6,0,2,0,6,6,4,1)})), \\
 &\quad \rightarrow V = -38 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \rightarrow V = -38 \\
 &0.5 * (1 + \gamma \theta^\top \phi(\text{grid}_{(4,0,6,6,4,6,0,6,4,1)})) + 0.5 * (1 + \gamma \theta^\top \phi(\text{grid}_{(4,0,6,6,4,6,0,6,4,1)})) \} \\
 &\quad \rightarrow V = -42 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \rightarrow V = -42
 \end{aligned}$$

$$= -29.6$$

Value Iteration w/Function Approximation --- Example

- After running the Bellman back-ups for all 4 states in S' we have:

$$V(\text{state 1}) = 6.4$$

(2,2,4,0, 0,2,4, 4, 0, 1)

$$V(\text{state 2}) = 19$$

(4,4,4,0, 0,0,4, 4, 0, 1)

$$V(\text{state 3}) = 19$$

(2,2,0,0, 0,2,0, 2, 0, 1)

$$V(\text{state 4}) = -29.6$$

(4,0,4,0, 4,4,4, 4, 0, 1)

- We now run supervised learning on these 4 examples to find a new θ :

$$\begin{aligned} \min_{\theta} & (6.4 - \theta^T \phi(\text{state 1}))^2 \\ & + (19 - \theta^T \phi(\text{state 2}))^2 \\ & + (19 - \theta^T \phi(\text{state 3}))^2 \\ & + ((-29.6) - \theta^T \phi(\text{state 4}))^2 \end{aligned}$$

Running least squares gives:

$$\theta^{(1)} = (0.195, 6.24, -2.11, 0, -6.05, 0.13, -2.11, 2.13, 0, 1.59)$$

Value Iteration with **Neural Net** Function Approximation

- Initialize by choosing some setting for $\theta^{(0)}$
- Iterate for $i = 0, 1, 2, \dots, H$:
 - Step 0: Pick some $S' \subseteq S$ (typically $|S'| \ll |S|$)

- Step 1: Bellman back-ups

$$\forall s \in S' : \bar{V}_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma \hat{V}_{\theta^{(i)}}(s') \right]$$

- Step 2: Supervised learning

find $\theta^{(i+1)}$ as the solution of:
$$\min_{\theta} \sum_{s \in S'} \left(\hat{V}_{\theta^{(i+1)}}(s) - \bar{V}_{i+1}(s) \right)^2$$

To avoid overfitting: only small number of gradient updates on objective or early stopping based on hold-out set

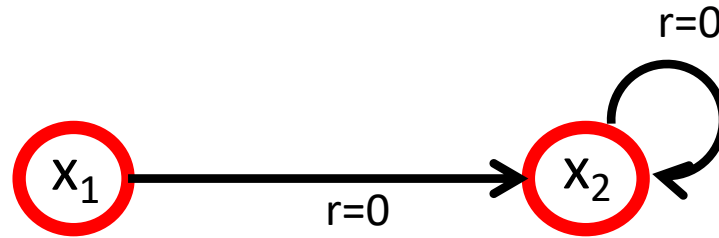
Potential Guarantees?

Theoretical Analysis of Value Iteration + Function Approximation

- We'll consider the following variation on the algorithm:
 - Assume we iterate over:
 - VI back-up for ALL states
 - Function approximation

Note: For ALL states is not practical (that's why we do function approximation). But (i) it's helpful to theoretically think through things; (ii) if we have a negative result, it's an even stronger negative result

Simple Example



θ

2θ

Function approximator: $[1 \ 2] * \theta$

Simple Example

$$\bar{J}_\theta = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \theta$$

$$\begin{aligned} \bar{J}^{(1)}(x_1) &= 0 + \gamma \hat{J}_{\theta^{(0)}}(x_2) = 2\gamma\theta^{(0)} \\ \bar{J}^{(1)}(x_2) &= 0 + \gamma \hat{J}_{\theta^{(0)}}(x_2) = 2\gamma\theta^{(0)} \end{aligned}$$

Function approximation with least squares fit:

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \theta^{(1)} \approx \begin{bmatrix} 2\gamma\theta^{(0)} \\ 2\gamma\theta^{(0)} \end{bmatrix}$$

Least squares fit results in:

$$\theta^{(1)} = \frac{6}{5}\gamma\theta^{(0)}$$

Repeated back-ups and function approximations result in:

$$\theta^{(i)} = \left(\frac{6}{5}\gamma\right)^i \theta^{(0)}$$

which diverges if $\gamma > \frac{5}{6}$ even though the function approximation class can represent the true value function.]

Composing Operators

- **Definition.** An operator G is a *non-expansion* with respect to a norm $\| \cdot \|$ if $\|GJ_1 - GJ_2\| \leq \|J_1 - J_2\|$
- **Fact.** If the operator F is a γ -contraction with respect to a norm $\| \cdot \|$ and the operator G is a non-expansion with respect to the same norm, then the sequential application of the operators G and F is a γ -contraction, i.e., $\|GFJ_1 - GFJ_2\| \leq \gamma\|J_1 - J_2\|$
- **Corollary.** If the supervised learning step is a non-expansion, then value iteration with function approximation is a γ -contraction, and in this case we have a convergence guarantee.

Averager Function Approximators Are Non-Expansions

DEFINITION: A real-valued function approximation scheme is an *averager* if every fitted value is the weighted average of zero or more target values and possibly some predetermined constants. The weights involved in calculating the fitted value Y_i may depend on the sample vector X_0 , but may not depend on the target values Y . More precisely, for a fixed X_0 , if Y has n elements, there must exist n real numbers k_i , n^2 nonnegative real numbers β_{ij} , and n nonnegative real numbers β_i , so that for each i we have $\beta_i + \sum_j \beta_{ij} = 1$ and $\hat{Y}_i = \beta_i k_i + \sum_j \beta_{ij} Y_j$.

- Examples:
 - nearest neighbor (aka state aggregation)
 - linear interpolation over triangles (tetrahedrons, ...)

Averager Function Approximators Are Non-Expansions

Proof: Let J_1 and J_2 be two vectors in \mathfrak{R}^n . Consider a particular entry s of ΠJ_1 and ΠJ_2 :

$$\begin{aligned} |(\Pi J_1)(s) - (\Pi J_2)(s)| &= \left| \beta_{s0} + \sum_{s'} \beta_{ss'} J_1(s') - \beta_{s0} + \sum_{s'} \beta_{ss'} J_2(s') \right| \\ &= \left| \sum_{s'} \beta_{ss'} (J_1(s') - J_2(s')) \right| \\ &\leq \max_{s'} |J_1(s') - J_2(s')| \\ &= \|J_1 - J_2\|_\infty \end{aligned}$$

This holds true for all s , hence we have

$$\|\Pi J_1 - \Pi J_2\|_\infty \leq \|J_1 - J_2\|_\infty$$

Guarantees for Fixed Point

Theorem. Let J^* be the optimal value function for a finite MDP with discount factor γ . Let the projection operator Π be a non-expansion w.r.t. the infinity norm and let \tilde{J} be any fixed point of Π . Suppose $\|\tilde{J} - J^*\|_\infty \leq \epsilon$. Then ΠT converges to a value function \bar{J} such that:

$$\|\bar{J} - J^*\| \leq 2\epsilon + \frac{2\gamma\epsilon}{1 - \gamma}$$

- I.e., if we pick a non-expansion function approximator which can approximate J^* well, then we obtain a good value function estimate.
- To apply to discretization: use continuity assumptions to show that J^* can be approximated well by chosen discretization scheme

Linear Regression ☹️

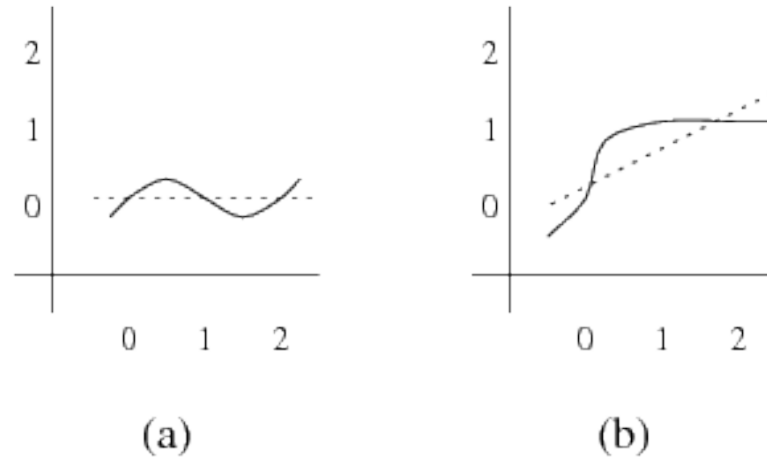


Figure 2: The mapping associated with linear regression when samples are taken at the points $x = 0, 1, 2$. In (a) we see a target value function (solid line) and its corresponding fitted value function (dotted line). In (b) we see another target function and another fitted function. The first target function has values $y = 0, 0, 0$ at the sample points; the second has values $y = 0, 1, 1$. Regression exaggerates the difference between the two functions: the largest difference between the two target functions at a sample point is 1 (at $x = 1$ and $x = 2$), but the largest difference between the two fitted functions at a sample point is $\frac{7}{6}$ (at $x = 2$).

Outline

- ✓ ■ Function approximation
- ✓ ■ Value iteration with function approximation
 - Policy iteration with function approximation
 - Linear programming with function approximation

Policy Iteration

One iteration of policy iteration:

- Policy evaluation: with fixed current policy π , find values with simplified Bellman updates:
 - Iterate until values converge

$$V_{i+1}^{\pi_k}(s) \leftarrow \sum_{s'} T(s, \pi_k(s), s') [R(s, \pi_k(s), s') + \gamma V_i^{\pi_k}(s')]$$

- Policy improvement: with fixed utilities, find the best action according to one-step look-ahead

$$\pi_{k+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_k}(s')]$$

Insert Function
Approximation Here

- Repeat until policy converges
- At convergence: optimal policy; and converges faster under some conditions

Approximate Policy Evaluation *is* a Contraction!

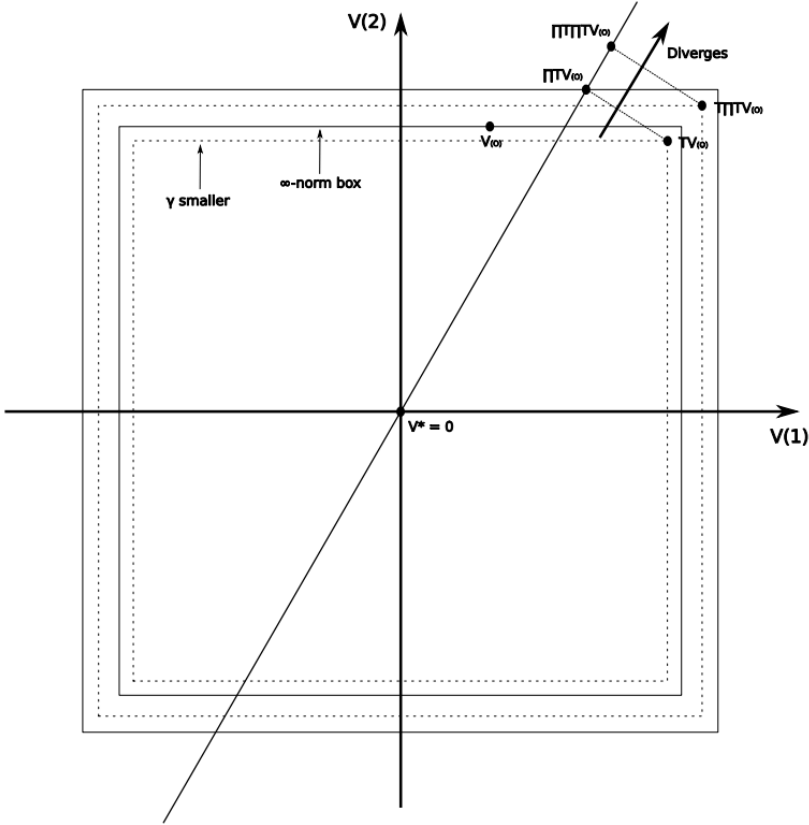
- IF we do weighted linear regression, weighted by the state visitation frequencies under the current policy
- THEN the resulting projection is a contraction w.r.t. the weighted 2-norm
- Policy Evaluation Bellman update is a contraction w.r.t. the same norm

→ Guaranteed convergence 😊 😊 😊

- Want to see the math:

<https://inst.eecs.berkeley.edu/~cs294-40/fa08/scribes/lecture5.pdf>

Extra Intermezzo on Incompatible Norms



Recent Related Paper**

- **Towards Characterizing Divergence in Deep Q-Learning,**
Joshua Achiam, Ethan Knight, Pieter Abbeel.
[arXiv 1903.08894](https://arxiv.org/abs/1903.08894)

Outline

- ✓ ■ Function approximation
- ✓ ■ Value iteration with function approximation
- ✓ ■ Policy iteration with function approximation
 - Linear programming with function approximation

Infinite Horizon Linear Program**

$$\begin{aligned} & \min_V \sum_{s \in S} \mu_0(s) V(s) \\ \text{s.t. } & V(s) \geq \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \quad \forall s \in S, a \in A \end{aligned}$$

Theorem. V^* is the solution to the above LP.

μ_0 is a probability distribution over S , with $\mu_0(s) > 0$ for all s in S .

Infinite Horizon Linear Program**

$$\begin{aligned} \min_V \quad & \sum_{s \in S} \mu_0(s) V(s) \\ \text{s.t.} \quad & V(s) \geq \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')], \quad \forall s \in S, a \in A \end{aligned}$$

Let $V(s) = \theta^\top \phi(s)$, and consider S' rather than S :

$$\begin{aligned} \min_\theta \quad & \sum_{s \in S'} \mu_0(s) \theta^\top \phi(s) \\ \text{s.t.} \quad & \theta^\top \phi(s) \geq \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \theta^\top \phi(s')], \quad \forall s \in S', a \in A \end{aligned}$$

We find approximate value function $\hat{V}_\theta(s) = \theta^\top \phi(s)$

Approximate Linear Program – Guarantees**

$$\begin{aligned} \min_{\theta} \quad & \sum_{s \in S'} \mu_0(s) \theta^\top \phi(s) \\ \text{s.t.} \quad & \theta^\top \phi(s) \geq \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \theta^\top \phi(s')] , \quad \forall s \in S', a \in A \end{aligned}$$

- LP solver will converge
- Solution quality: [de Farias and Van Roy, 2002]

Assuming one of the features is the feature that is equal to one for all states, and assuming $S'=S$ we have that:

$$\|V^* - \Phi\theta\|_{1, \mu_0} \leq \frac{2}{1 - \gamma} \min_{\theta} \|V^* - \Phi\theta\|_{\infty}$$

(slightly weaker, probabilistic guarantees hold for S' not equal to S , these guarantees require size of S' to grow as the number of features grows)