

# Deep learning and reinforcement learning

---

Jesús Fernández-Villaverde<sup>1</sup> and Galo Nuño<sup>2</sup>

October 15, 2021

<sup>1</sup>University of Pennsylvania

<sup>2</sup>Banco de España

## **A short introduction**

---

# The problem

- Let us suppose we want to approximate an unknown function:

$$y = f(\mathbf{x})$$

where  $y$  is a scalar and  $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$  a vector.

- We care about the case when  $N$  is large.
- Easy to generalize to the case where  $y$  is a vector (or a probability distribution), but notation becomes cumbersome.
- In economics,  $f(\mathbf{x})$  can be a value function, a policy function, a pricing kernel, a conditional expectation, a classifier, ...

# A neural network

- An artificial neural network (a.k.a. ANN or connectionist system) is an approximation to  $f(\mathbf{x})$  built as a linear combination of  $M$  generalized linear models of  $\mathbf{x}$  of the form:

$$y \cong g^{NN}(\mathbf{x}; \theta) = \theta_0 + \sum_{m=1}^M \theta_m \phi(z_m)$$

where  $\phi(\cdot)$  is an arbitrary activation function and:

$$z_m = \theta_{0,m} + \sum_{n=1}^N \theta_{n,m} x_n$$

- $M$  is known as the width of the model.
- We can select  $\theta$  such that  $g^{NN}(\mathbf{x}; \theta)$  is as close to  $f(\mathbf{x})$  as possible given some relevant metric (e.g.,  $L^2$  norm).
- This is known as “training” the network.

## Comparison with other approximations

- Compare:

$$y \cong g^{NN}(\mathbf{x}; \theta) = \theta_0 + \sum_{m=1}^M \theta_m \phi \left( \theta_{0,m} + \sum_{n=1}^N \theta_{n,m} x_n \right)$$

with a standard projection:

$$y \cong g^{CP}(\mathbf{x}; \theta) = \theta_0 + \sum_{m=1}^M \theta_m \phi_m(\mathbf{x})$$

where  $\phi_m$  is, for example, a Chebyshev polynomial.

- We exchange the rich parameterization of coefficients for the parsimony of basis functions.
- Later, we will explain why this is often a good idea.
- How we determine the coefficients will also be different, but this is somewhat less important.

- A deep learning network is an acyclic *multilayer* composition of  $J > 1$  neural networks:

$$y \cong g^{DL}(\mathbf{x}; \theta) = g^{NN(1)}\left(g^{NN(2)}\left(\dots; \theta^{(2)}\right); \theta^{(1)}\right)$$

where the  $M^{(1)}, M^{(2)}, \dots$  and  $\phi^1(\cdot), \phi^2(\cdot), \dots$  are possibly different across each layer of the network.

- Sometimes known as deep feedforward neural networks or multilayer perceptrons.
- “Feedforward” comes from the fact that the composition of neural networks can be represented as a directed acyclic graph, which lacks feedback. We can have more general recurrent structures.
- $J$  is known as the depth of the network. The case  $J = 1$  is a standard neural network.
- As before, we can select  $\theta$  such that  $g^{DL}(\mathbf{x}; \theta)$  approximates a target function  $f(\mathbf{x})$  as closely as possible under some relevant metric.

# Why are neural networks a good solution method in economics?

- From now on, I will refer to neural networks as including both single and multilayer networks.
- With suitable choices of activation functions, neural networks can efficiently approximate extremely complex functions.
- In particular, under certain (relatively weak) conditions:
  1. Neural networks are universal approximators.
  2. Neural networks break the “curse of dimensionality.”
- Furthermore, neural networks are easy to code, stable, and scalable for multiprocessing.
- Thus, neural networks have considerable option value as solution methods in economics.

## Current interest

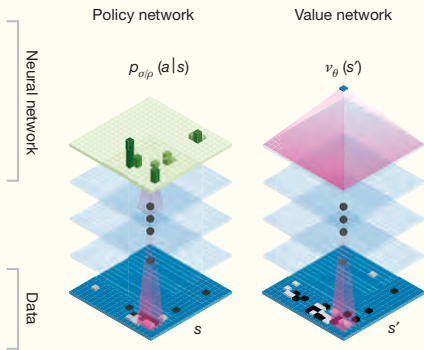
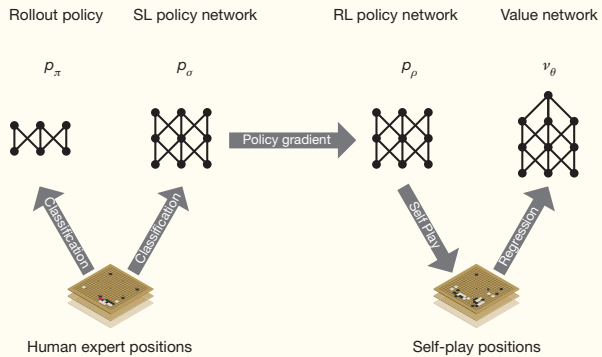
- Currently, neural networks are among the most active areas of research in computer science and applied math.
- While original idea goes back to the 1940s, neural networks were rediscovered in the second half of the 2000s.
- Why?
  1. Suddenly, the large computational and data requirements required to train the networks efficiently became available at a reasonable cost.
  2. New algorithms such as *back propagation* through gradient descent became popular.
- Some well-known successes and industrial applications.

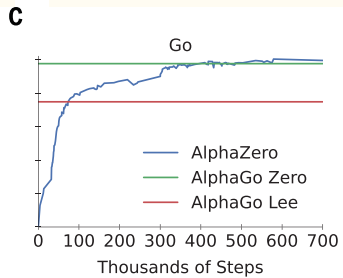
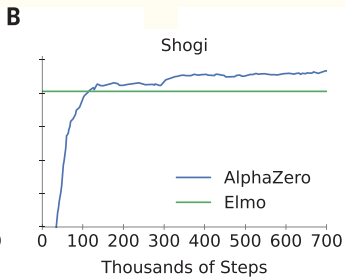
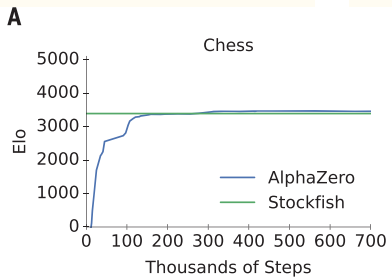




# AlphaGo

- Big splash: AlphaGo vs. Lee Sedol in March 2016.
- *Silver et al. (2018)*: now applied to chess, shogi, Go, and StarCraft II.
- Check also:
  1. <https://deepmind.com/research/alphago/>.
  2. <https://www.alphagomovie.com/>
  3. <https://deepmind.com/blog/article/alphastar-mastering-real-time-strategy-game-starcraft-ii>
- Very different than Deep Blue against Kasparov.
- New and surprising strategies.
- However, you need to keep this accomplishment in perspective.





## Further advantages

- Neural networks and deep learning often require less “inside knowledge” by experts on the area.
- Results can be highly counter-intuitive and yet, deliver excellent performance.
- Outstanding open source libraries: Tensorflow, Pytorch, Flux.
- More recently, development of dedicated hardware (TPUs, AI accelerators, FPGAs) are likely to maintain a hedge for the area.
- The width of an ecosystem is key for its long-run success.



TensorFlow

## Limitations of neural networks and deep learning

- While neural networks and deep learning can work extremely well, there is no such a thing as a silver bullet.
- Clear and serious trade-offs in real-life applications.
- Rule-of-thumb in the industry is that one needs around  $10^7$  labeled observations to properly train a complex ANN with around  $10^4$  observations in each relevant group.
- Of course, sometimes “observations” are endogenous (we can simulate them), but if your goal is to forecast GDP next quarter, it is unlikely a neural network will beat an ARIMA(n,p,q) (at least only with macro variables).
- Issues of interpretation.



# DEEP LEARNING

Ian Goodfellow, Yoshua Bengio,  
and Aaron Courville



## Digging deeper

---

## More details on neural networks

- Non-linear functional approximation method.
- Much hype around them and over-emphasis of biological interpretation.
- We will follow a much sober formal treatment (which, in any case, agrees with state-of-art researchers approach).
- In particular, we will highlight connections with econometrics (e.g., NOLS, semiparametric regression, and sieves).
- We will start describing the simplest possible neural network.

## A neuron

- $N$  observables:  $x_1, x_2, \dots, x_N$ . We stack them in  $\mathbf{x}$ .
- Coefficients (or weights):  $\theta_0$  (a constant),  $\theta_1, \theta_2, \dots, \theta_N$ . We stack them in  $\theta$ .
- We build a linear combination of observations:

$$z = \theta_0 + \sum_{n=1}^N \theta_n x_n$$

Theoretically, we could build non-linear combinations, but unlikely to be a fruitful idea in general.

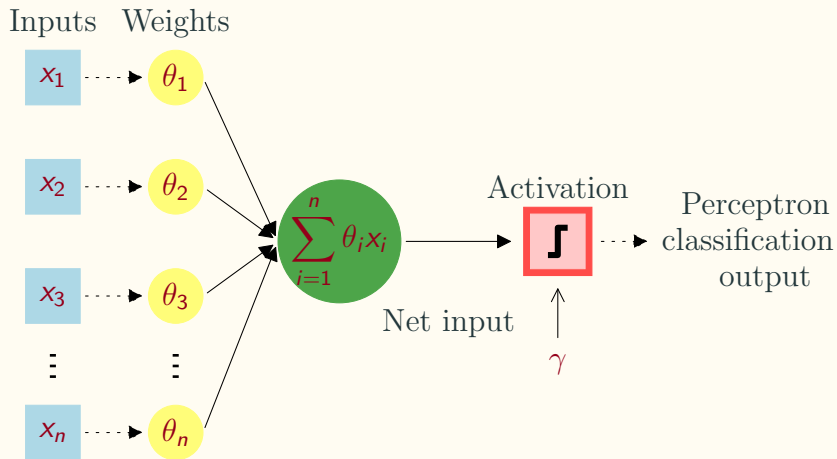
- We transform such linear combination with an activation function:

$$y = g(\mathbf{x}; \theta) = \phi(z)$$

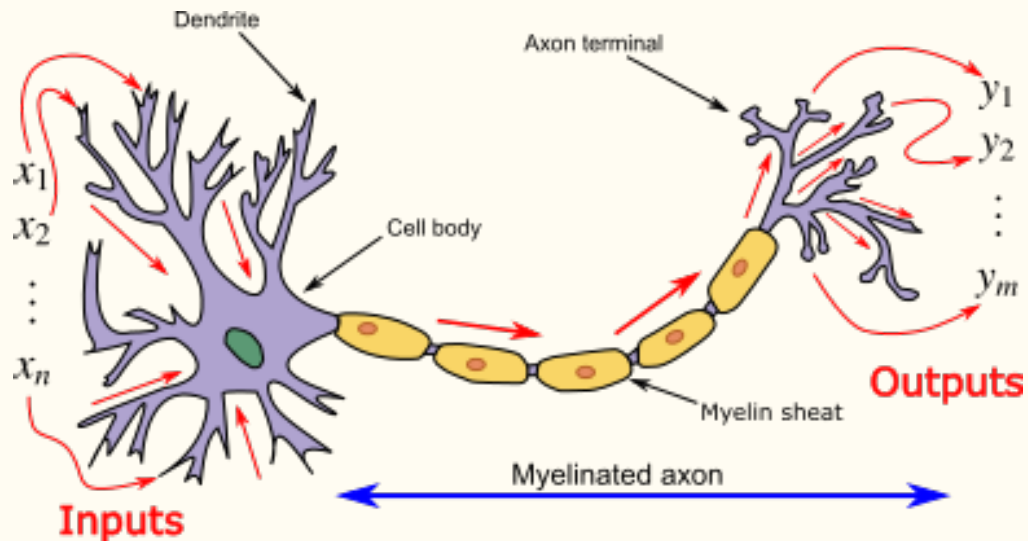
The activation function might have some coefficients  $\gamma$  on its own.

- Why do we need an activation function?

# Flow representation



## The biological analog



# Activation functions I

- Traditionally:

1. Identity function:

$$\phi(z) = z$$

Used in linear regression.

2. A sigmoidal function:

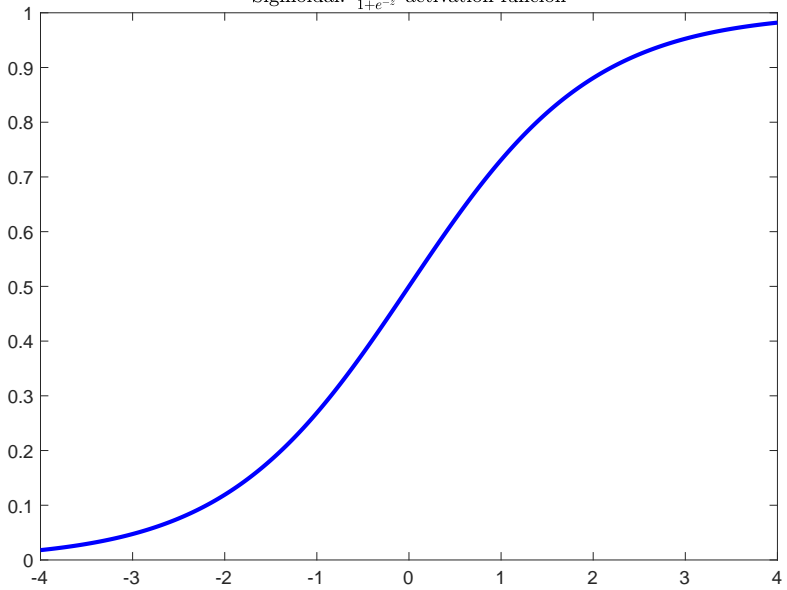
$$\phi(z) = \frac{1}{1 + e^{-z}}$$

A particular limiting case as  $z$  grows quickly: step function.

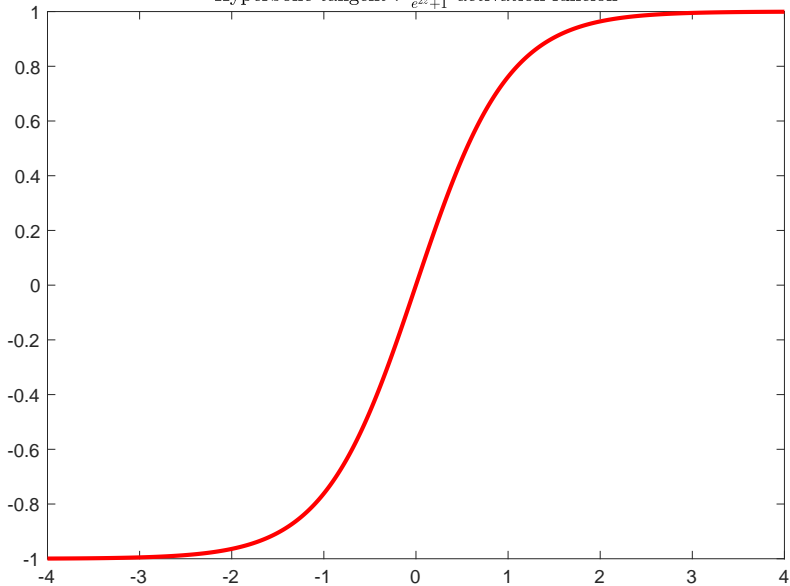
3. Hyperbolic tangent:

$$\phi(z) = \frac{e^{2z} - 1}{e^{2z} + 1}$$

Sigmoidal:  $\frac{1}{1+e^{-z}}$  activation function



Hyperbolic tangent :  $\frac{e^{2z}-1}{e^{2z}+1}$  activation function





# Activation functions II

- Some activation functions that have gained popularity recently:

1. Rectified linear unit (ReLU):

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

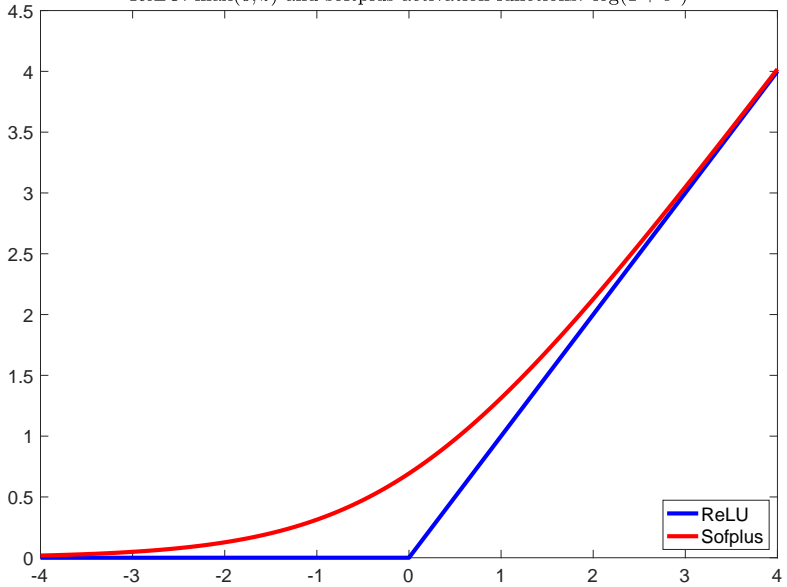
2. Parametric ReLU:

$$\phi(z) = \max(z, az)$$

3. Softplus:

$$\phi(z) = \log(1 + e^z)$$

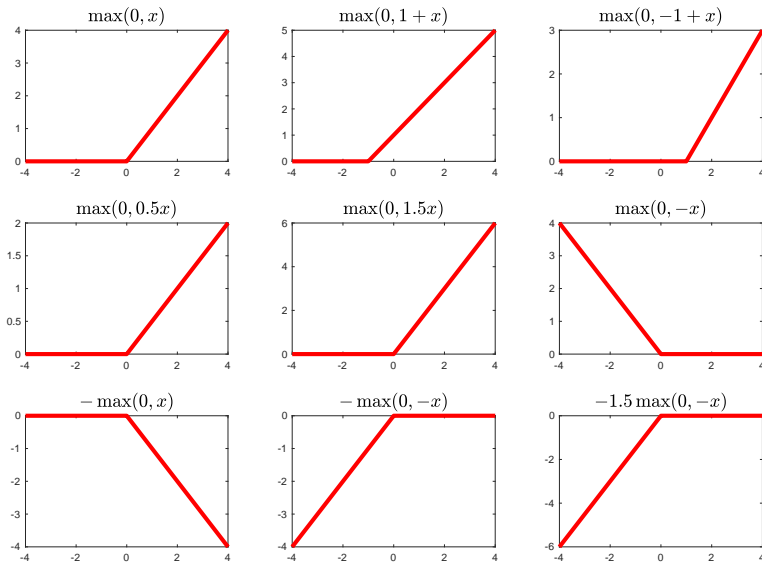
ReLU:  $\max(0, z)$  and softplus activation functions:  $\log(1 + e^z)$



# Interpretation

- $\theta_0$  controls the activation threshold.
- The level of the  $\theta_i$ 's for  $i > 0$  control the activation rate (the higher the  $\theta_i$ 's, the harder the activation).
- Some textbooks separate the activation threshold and scaling coefficients from  $\theta$  as different coefficients in  $\phi$ , but such separation moves notation farther away from standard econometrics.
- Potential identification problem between  $\theta$  and more general activation functions with their own parameters.
- But in practice  $\theta$  does not have a structural interpretation, so the identification problem is of secondary importance.
- As mentioned in the introduction, a neuron closely resembles a generalized linear model in econometrics.

Different ReLUs:  $\theta_i \max(0, \theta_{i,0} + \theta_{i,1}x)$



## Combining neurons into a neural network

- As before, we have  $N$  observables:  $x_1, x_2, \dots, x_N$ .
- Coefficients (or weights):  $\theta_{0,m}$  (a constant),  $\theta_{1,m}, \theta_{2,m}, \dots, \theta_{N,m}$ .

- We build  $M$  linear combinations of observations:

$$z_m = \theta_{0,m} + \sum_{n=1}^N \theta_{n,m} x_n$$

- We transform and add such linear combinations with an activation function:

$$y \cong g(\mathbf{x}; \theta) = \theta_0 + \sum_{m=1}^M \theta_m \phi(z_m)$$

- Also, quasi-linear structure in terms of vectors of observables and coefficients.
- This is known as a single layer network.

# Two classic (yet remarkable) results I

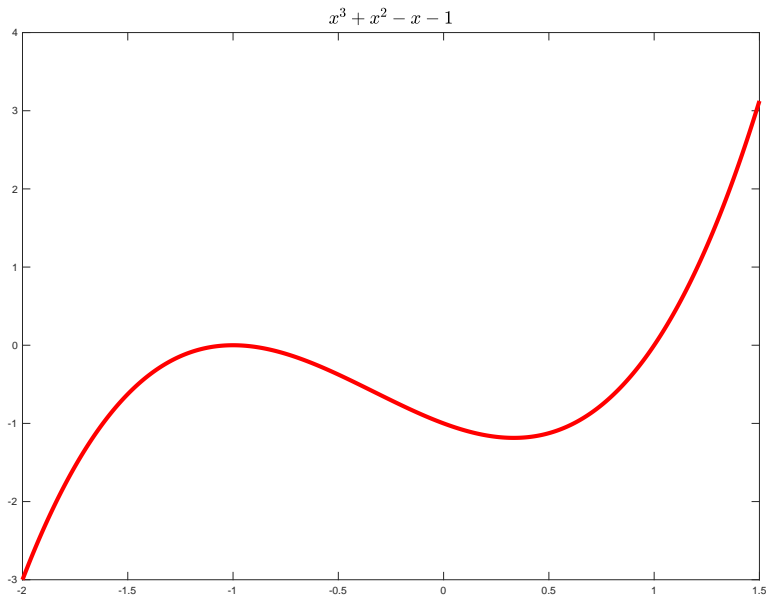
## Borel measurable function

A map  $f : X \rightarrow Y$  between two topological spaces is called Borel measurable if  $f^{-1}(A)$  is a Borel set for any open set  $A$  on  $Y$  (the Borel sets are all the open sets built through the operations of countable union, countable intersection, and relative complement).

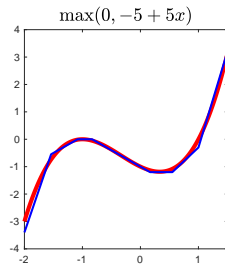
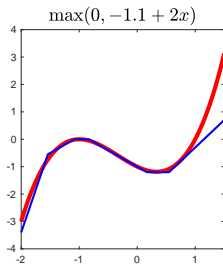
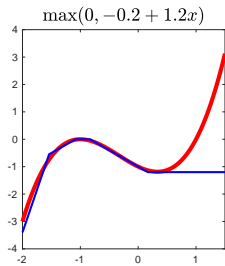
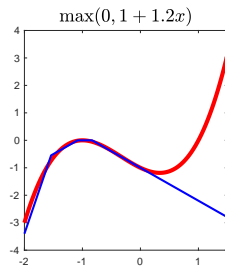
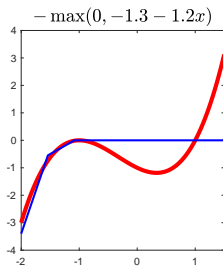
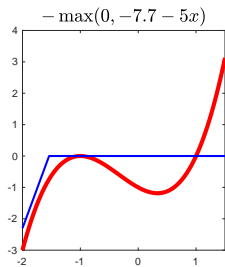
## Universal approximation theorem: Hornik, Stinchcombe, and White (1989)

A neural network with at least one hidden layer can approximate any Borel measurable function mapping finite-dimensional spaces to any desired degree of accuracy.

- Intuition of the result.
- Comparison with other results in series approximations.



### A six ReLUs approximation





## Two classic (yet remarkable) results II

- Assume, as well, that we are dealing with the class of functions for which the Fourier transform of their gradient is integrable.

### Breaking the curse of dimensionality: Barron (1993)

A one-layer NN achieves integrated square errors of order  $\mathcal{O}(1/M)$ , where  $M$  is the number of nodes. In comparison, for series approximations, the integrated square error is of order  $\mathcal{O}(1/(M^{2/N}))$  where  $N$  is the dimensions of the function to be approximated.

- More general theorems by [Leshno et al. \(1993\)](#) and [Bach \(2017\)](#).
- What about Chebyshev polynomials? Splines? Problems of convergence and extrapolation.
- There is another, yet more subtle curse of dimensionality.

## Training the network

- $\theta$  is selected to minimize the quadratic error function  $\mathcal{E}(\theta; \mathbf{Y}, \hat{\mathbf{y}})$ :

$$\begin{aligned}\theta^* &= \arg \min_{\theta} \mathcal{E}(\theta; \mathbf{Y}, \hat{\mathbf{y}}) \\ &= \arg \min_{\theta} \sum_{j=1}^J \mathcal{E}(\theta; y_j, \hat{y}_j) \\ &= \arg \min_{\theta} \frac{1}{2} \sum_{j=1}^J \|y_j - g(\mathbf{x}_j; \theta)\|^2\end{aligned}$$

- Where from do the observations  $\mathbf{Y}$  come? Observed data vs. simulated epochs.
- How do we solve this minimization problem?
- Other objective functions are possible.

# Back propagation

- In general, we can easily calculate  $\mathcal{E}(\theta^*; Y, \hat{\mathbf{y}})$  and  $\nabla \mathcal{E}(\theta^*; Y, \hat{\mathbf{y}})$  for a given  $\theta^*$ .
- In particular, for the gradient, we can use *back propagation* (Rumelhart et al., 1986):

$$\frac{\partial \mathcal{E}(\theta; y_j, \hat{y}_j)}{\partial \theta_0} = y_j - g(\mathbf{x}_j; \theta)$$

$$\frac{\partial \mathcal{E}(\theta; y_j, \hat{y}_j)}{\partial \theta_m} = (y_j - g(\mathbf{x}_j; \theta)) \phi(z_m), \text{ for } \forall m$$

$$\frac{\partial \mathcal{E}(\theta; y_j, \hat{y}_j)}{\partial \theta_{0,m}} = (y_j - g(\mathbf{x}_j; \theta)) \theta_m \phi'(z_m), \text{ for } \forall m$$

$$\frac{\partial \mathcal{E}(\theta; y_j, \hat{y}_j)}{\partial \theta_{n,m}} = (y_j - g(\mathbf{x}_j; \theta)) \theta_m x_n \phi'(z_m), \text{ for } \forall n, m$$

where  $\phi'(z)$  is the derivative of the activation function.

- The derivative  $\phi'(z)$  will be trivial to evaluate if we use a ReLU.
- Back propagation will be particularly important below when we introduce multiple layers.

## An approach to minimization

- One approach to optimization is to minimize a local model that approximates the true objective function.
- The local model can be a first- or second-order Taylor approximation of the objective function.
- For example, suppose a function  $\mathcal{E}$  is roughly approximated as a quadratic form:

$$\mathcal{E}(\theta) \approx \frac{1}{2}\theta^T \mathbf{A}\theta - \mathbf{b}^T \theta + c$$

where  $\mathbf{A}$  is a square, symmetric, positive-definite matrix.

- Then  $\mathcal{E}(\theta)$  is minimized by the solution to:

$$\mathbf{A}\theta = \mathbf{b}$$

- We can use this result to build a descent direction iteration if we know  $\mathbf{A}$  and  $\mathbf{b}$  (or we have approximations to them).

# Descent direction iteration

- Starting at point  $\theta^{(1)}$ , a descent direction algorithm generates sequence of steps (called iterates) that converge to a local minimum.
- The descent direction iteration algorithm:
  1. At iteration  $k$ , check whether  $\theta^{(k)}$  satisfies termination condition. If so stop; otherwise go to step 2.
  2. Determine the descent direction  $\mathbf{d}^{(k)}$  using local information such as gradient or Hessian.
  3. Compute step size  $\alpha^{(k)}$ .
  4. Compute the next candidate point:  $\theta^{(k+1)} \leftarrow \theta^{(k)} + \alpha^{(k)} \mathbf{d}^{(k)}$ .
- Choice of  $\alpha$  and  $\mathbf{d}$  determines the flavor of the algorithm.

# Gradient descent method

- A natural choice for  $\mathbf{d}$  is the direction of steepest descent (first proposed by Cauchy).
- The direction of steepest descent is given by the direction opposite the gradient  $\nabla\mathcal{E}(\theta)$ . Thus, a.k.a. steepest descent.
- If function is smooth and the step size small, the method leads to improvement (as long as the gradient is not zero).
- The normalized direction of steepest descent is:

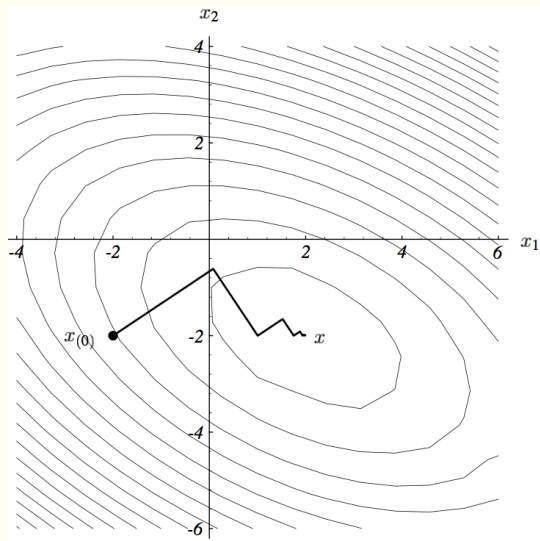
$$\mathbf{d}^{(k)} = -\frac{\nabla\mathcal{E}(\theta^{(k)})}{\|\nabla\mathcal{E}(\theta^{(k)})\|}$$

- One way to set the step size is to solve:

$$\alpha^k = \arg \min_{\alpha} \mathcal{E}(\theta^{(k)} + \alpha\mathbf{d}^{(k)})$$

- Under this step size choice, it can be shown  $\mathbf{d}^{(k+1)}$  and  $\mathbf{d}^{(k)}$  are orthogonal.

# Steepest descent method

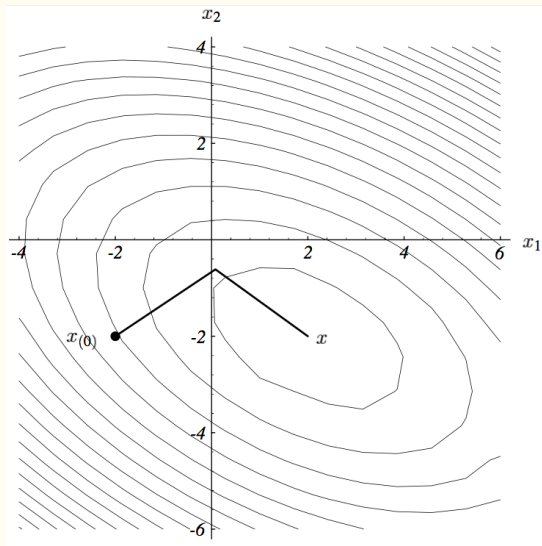


# Conjugate descent method

- Gradient descent can perform poorly in narrow valleys (it may require many steps to make progress).
- Famous example: Rosenbrock function.
- The *conjugate gradient* method overcomes this problem by somehow constructing to be conjugate to the old gradient, and to all previous directions traversed.
- Define  $\mathbf{g}(\theta) = \nabla \mathcal{E}(\theta)$ .
- In first iteration, set:  $\mathbf{d}^{(1)} = -\mathbf{g}(\theta^{(1)})$  and  $\theta^{(2)} = \theta^{(1)} + \alpha^{(1)}\mathbf{d}^{(1)}$ . Here,  $\alpha^{(1)}$  is arbitrary.
- Subsequent iterations set  $\mathbf{d}^{(k+1)} = -\mathbf{g}^{(k+1)} + \beta^{(k)}\mathbf{d}^{(k)}$ .



# Conjugate descent method



# Conjugate descent method

- There are two approaches to set  $\beta$ :

1. Fletcher-Reeves:

$$\beta^{(k)} = \frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k-1)T} \mathbf{g}^{(k-1)}}$$

2. Olak-Ribiere:

$$\beta^{(k)} = \frac{\mathbf{g}^{(k)T} (\mathbf{g}^{(k)} - \mathbf{g}^{(k-1)})}{\mathbf{g}^{(k-1)T} \mathbf{g}^{(k-1)}}$$

- The Olak-Ribiere requires an automatic reset at every iteration:  $\beta \leftarrow \max(\beta, 0)$ .
- If the function to minimize has flat areas, one can introduce a *momentum* update equation:

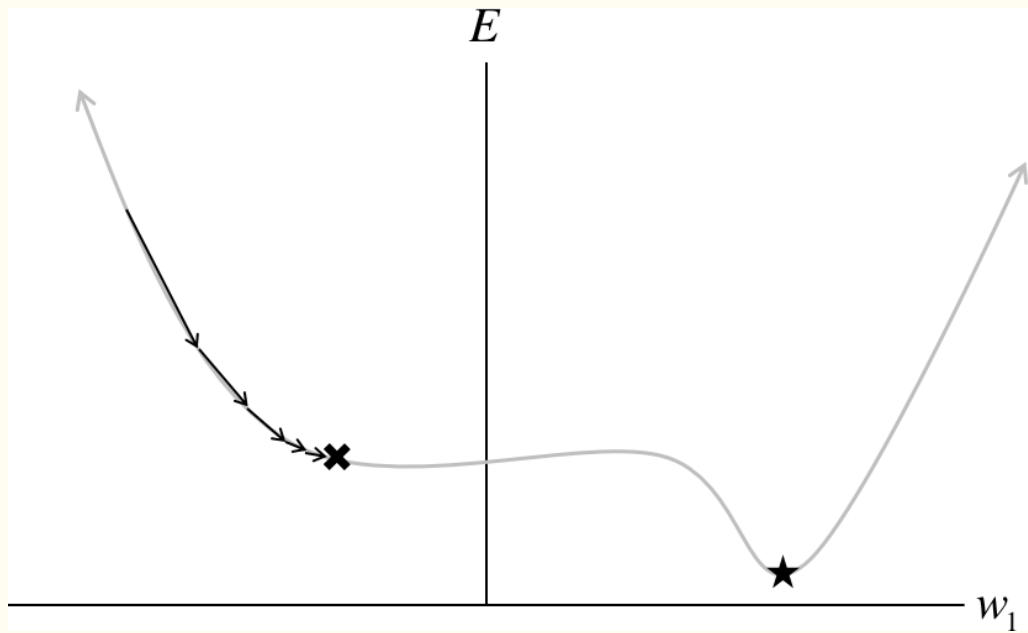
$$\mathbf{v}^{(k+1)} = \beta \mathbf{v}^{(k)} - \alpha \mathbf{g}^{(k)}$$

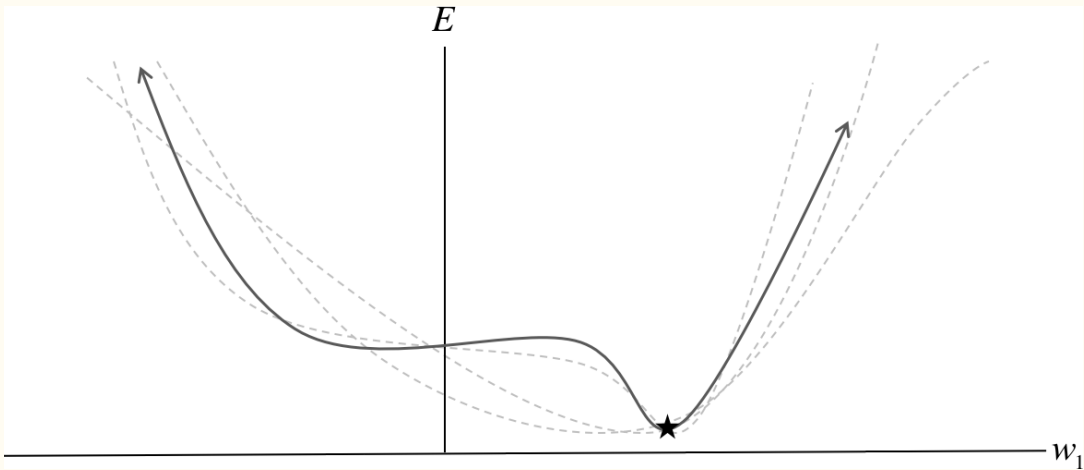
$$\theta^{(k+1)} = \theta^{(k)} + \mathbf{v}^{(k+1)}$$

- The modification reverts to the gradient descent version if  $\beta = 0$ .
- Intuitively, the momentum update is like a ball rolling down an almost horizontal surface.

# Stochastic gradient descent and minibatch

- Even with back propagation, evaluating the gradient for the whole training set can be costly.
- Stochastic gradient descent: Intuition from Monte Carlo.
- An additional advantage.
- A compromise between using the whole training set and pure stochastic gradient descent: minibatch gradient descent.
- This is the most popular algorithm to train neural networks.
- Intuition from GMM. Notice also resilience to scaling.
- In practice, we do not need a global min ( $\neq$  likelihood).
- You can flush the algorithm to a graphics processing unit (GPU) or a tensor processing unit (TPU) instead of a standard CPU.

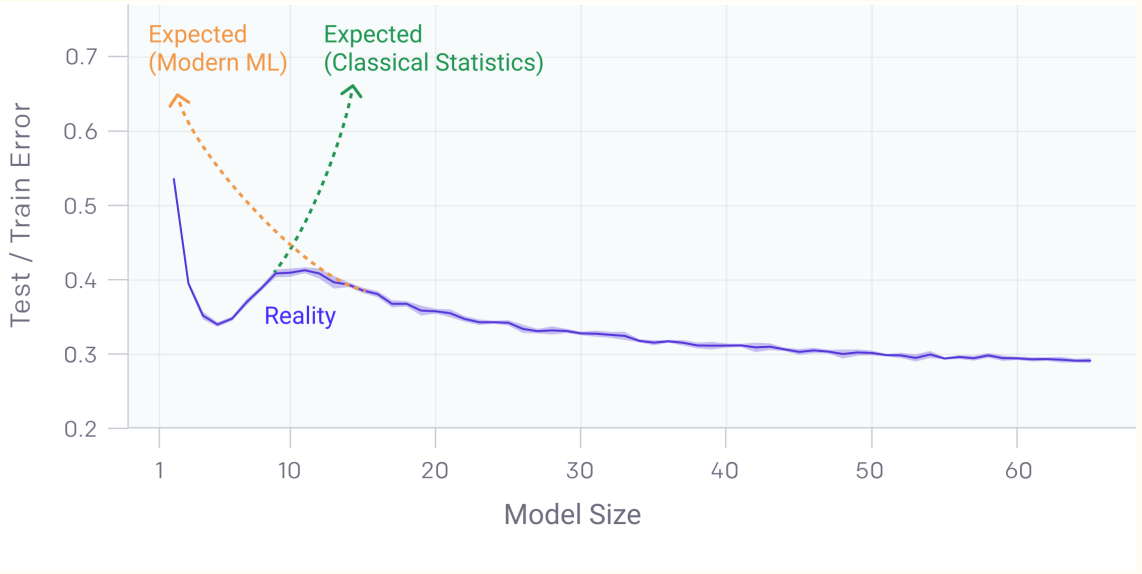




# Alternative minimization algorithms

1. More sophisticated stochastic gradient descent: Adam (Adaptive Moment Estimation). It uses running averages of both the gradients and the second moments of the gradients.
2. Newton and Quasi-Newton methods are unlikely to be of much use in practice. Why?
3. McMc/Simulated annealing.
4. Genetic algorithms:
  - In fact, much of the research in deep learning incorporates some flavor of genetic selection.
  - Basic idea.

- Design of the network architecture:
  1. Trade-off error/computational time.
  2. Better to err on the side of too many  $M$ .
- Double descent phenomenon.





# Multiple layers I

- The hidden layers can be multiplied without limit in a feed-forward ANN.
- We build  $K$  layers:

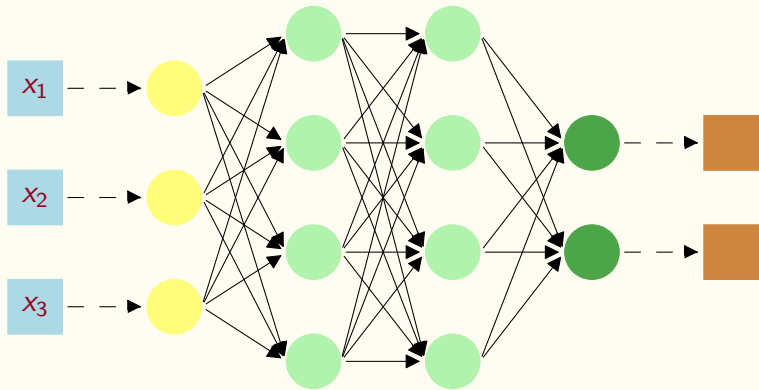
$$z_m^1 = \theta_{0,m}^1 + \sum_{n=1}^N \theta_{n,m}^1 x_n$$

and

$$z_m^2 = \theta_{0,m}^2 + \sum_{m=1}^M \theta_m^2 \phi(z_m^1)$$

...

$$y \cong g(x; \theta) = \theta_0^K + \sum_{m=1}^M \theta_m^K \phi(z_m^{K-1})$$



Input Values

Hidden Layer 1

Output Layer

Input Layer

Hidden Layer 2

## Multiple layers II

- Why do we want to introduce hidden layers?
  1. It works! Our brains have six layers. AlphaGo has 12 layers with ReLUs.
  2. Hidden layers induce highly nonlinear behavior.
  3. Allow for clustering of variables.
- We can have different  $M$ 's in each layer  $\Rightarrow$  fewer neurons in higher layers allow for compression of learning into fewer features.
- We can also add multidimensional outputs.
- Or even to produce, as output, a probability distribution, for example, using a softmax layer:

$$y_m = \frac{e^{z_m^{K-1}}}{\sum_{m=1}^M e^{z_m^{K-1}}}$$

## **Application to Economics**

---

# Solving high-dimensional dynamic programming problems using Deep Learning

- Our goal is to solve the recursive continuous-time Hamilton-Jacobi-Bellman (HJB) equation globally:

$$\rho V(\mathbf{x}) = \max_{\alpha} r(\mathbf{x}, \alpha) + \nabla_{\mathbf{x}} V(\mathbf{x}) f(\mathbf{x}, \alpha) + \frac{1}{2} \text{tr}(\sigma(\mathbf{x}))^T \Delta_{\mathbf{x}} V(\mathbf{x}) \sigma(\mathbf{x}))$$

$$\text{s.t. } G(\mathbf{x}, \alpha) \leq \mathbf{0} \quad \text{and} \quad H(\mathbf{x}, \alpha) = \mathbf{0}$$

- Think about the cases where we have many state variables.
- Alternatives for this solution?

- We define four neural networks:
  1.  $\tilde{V}(\mathbf{x}; \Theta^V) : \mathbb{R}^N \rightarrow \mathbb{R}$  to approximate the value function  $V(\mathbf{x})$ .
  2.  $\tilde{\alpha}(\mathbf{x}; \Theta^\alpha) : \mathbb{R}^N \rightarrow \mathbb{R}^M$  to approximate the policy function  $\alpha$ .
  3.  $\tilde{\mu}(\mathbf{x}; \Theta^\mu) : \mathbb{R}^N \rightarrow \mathbb{R}^{L_1}$ , and  $\tilde{\lambda}(\mathbf{x}; \Theta^\lambda) : \mathbb{R}^N \rightarrow \mathbb{R}^{L_2}$  to approximate the Karush-Kuhn-Tucker (KKT) multipliers  $\mu$  and  $\lambda$ .
- To simplify notation, we accumulate all weights in the matrix  $\Theta = (\Theta^V, \Theta^\alpha, \Theta^\mu, \Theta^\lambda)$ .
- We could think about the approach as just one large neural network with multiple outputs.

# Error criterion I

- The HJB error:

$$\begin{aligned} err_{HJB}(\mathbf{x}; \Theta) \equiv & r(\mathbf{x}, \tilde{\alpha}(\mathbf{s}; \Theta^\alpha)) + \nabla_x \tilde{V}(\mathbf{x}; \Theta^V) f(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha)) + \\ & + \frac{1}{2} tr[\sigma(\mathbf{x})^T \Delta_x \tilde{V}(\mathbf{x}; \Theta^V) \sigma(\mathbf{x})] - \rho \tilde{V}(\mathbf{x}; \Theta^V) \end{aligned}$$

- The policy function error:

$$\begin{aligned} err_\alpha(\mathbf{x}; \Theta) \equiv & \frac{\partial r(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha))}{\partial \alpha} + D_\alpha f(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha))^T \nabla_x \tilde{V}(\mathbf{x}; \Theta^V) \\ & - D_\alpha G(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha))^T \tilde{\mu}(\mathbf{x}; \Theta^\mu) - D_\alpha H(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha)) \tilde{\lambda}(\mathbf{x}; \Theta^\lambda), \end{aligned}$$

where  $D_\alpha G \in \mathbb{R}^{L_1 \times M}$ ,  $D_\alpha H \in \mathbb{R}^{L_2 \times M}$ , and  $D_\alpha f \in \mathbb{R}^{N \times M}$  are the submatrices of the Jacobian matrices of  $G$ ,  $H$  and  $f$  respectively containing the derivatives with respect to  $\alpha$ .

## Error criterion II

- The constraint error is itself composed of the primal feasibility errors:

$$err_{PF_1}(\mathbf{x}; \Theta) \equiv \max\{0, G(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha))\}$$

$$err_{PF_2}(\mathbf{x}; \Theta) \equiv H(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha))$$

the dual feasibility error:

$$err_{DF}(\mathbf{x}; \Theta) = \max\{0, -\tilde{\mu}(\mathbf{x}; \Theta^\mu)\}$$

and the complementary slackness error:

$$err_{CS}(\mathbf{x}; \Theta) = \tilde{\mu}(\mathbf{x}; \Theta)^\top G(\mathbf{x}, \tilde{\alpha}(\mathbf{x}; \Theta^\alpha))$$

- We combine these four errors by using the squared error as our loss criterion:

$$\begin{aligned} \mathcal{E}(\mathbf{x}; \Theta) \equiv & \|err_{HJB}(\mathbf{x}; \Theta)\|_2^2 + \|err_\alpha(\mathbf{x}; \Theta)\|_2^2 + \|err_{PF_1}(\mathbf{x}; \Theta)\|_2^2 + \\ & + \|err_{PF_2}(\mathbf{x}; \Theta)\|_2^2 + \|err_{DF}(\mathbf{x}; \Theta)\|_2^2 + \|err_{CS}(\mathbf{x}; \Theta)\|_2^2 \end{aligned}$$



# Training

- We train our neural networks by minimizing the above error criterion through mini-batch gradient descent over points drawn from the ergodic distribution of the state vector.
- The efficient implementation of this last step is the key to the success of our algorithm.
- We start by initializing our network weights and we perform  $K$  learning steps called epochs, where  $K$  can be chosen in a variety of ways.
- For each epoch, we draw  $I$  points from the state space by simulating from the ergodic distribution.
- Then, we randomly split this sample into  $B$  mini-batches of size  $S$ . For each mini-batch, we define the mini-batch error, by averaging the loss function over the batch.
- Finally, we perform mini-batch gradient descent for all network weights, with  $\eta_k$  being the learning rate in the  $k$ -th epoch.

## An Example

---

# The continuous-time neoclassical growth model I

- We start with the continuous-time neoclassical growth model because it has closed-form solutions for the policy functions, which allows us to focus our attention on the analysis of the value function approximation.
- We can then back out the policy function from this approach and compare it to the results of the next step in which we approximate the policy functions themselves with a neural net.
- A single agent deciding to either save in capital or consume with a HJB equation :

$$\rho V(k) = \max_c U(c) + V'(k)[F(k) - \delta * k - c]$$

- Notice that  $c = (U')^{-1}(V'(k))$ . With CRRA utility, this simplifies further to  $c = (V'(k))^{-\frac{1}{\gamma}}$ .
- We set  $\gamma = 2$ ,  $\rho = 0.04$ ,  $F(k) = 0.5 * k^{0.36}$ ,  $\delta = 0.05$ .

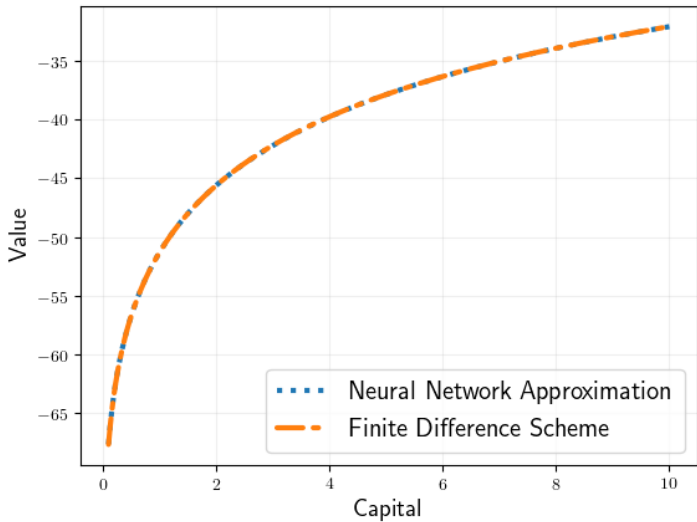
## The continuous-time neoclassical growth model II

- We approximate the value function  $V(k)$  with a neural network,  $\tilde{V}(k; \Theta)$  with an “HJB error”:

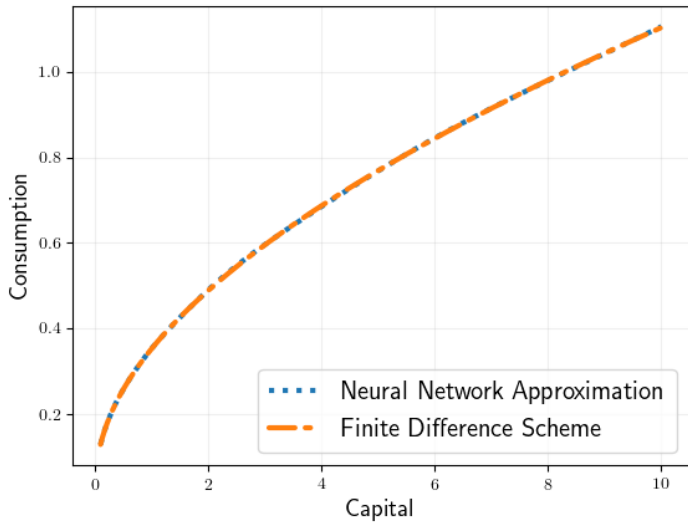
$$\begin{aligned} err_{HJB} = & \rho \tilde{V}(k; \Theta) - U \left( (U')^{-1} \left( \frac{\partial \tilde{V}(k; \Theta)}{\partial k} \right) \right) \\ & - \frac{\partial \tilde{V}(k; \Theta)}{\partial k} \left[ F(k) - \delta * k - (U')^{-1} \left( \frac{\partial \tilde{V}(k; \Theta)}{\partial k} \right) \right] \end{aligned}$$

- Details:

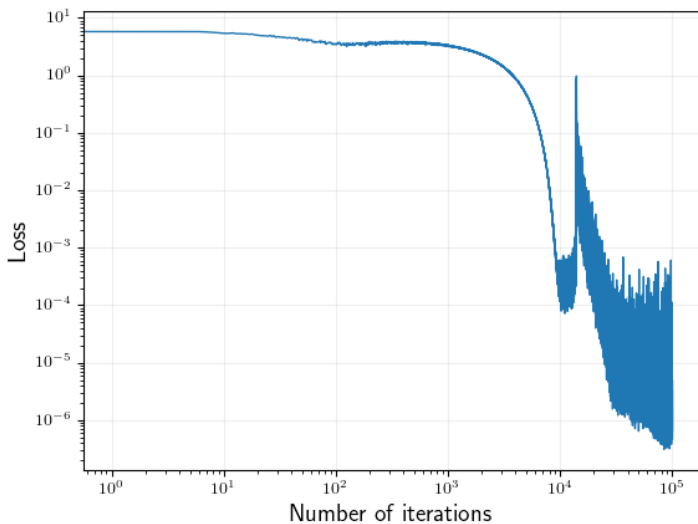
1. 3 layers.
2. 8 neurons per layers.
3.  $\tanh(x)$  activation.
4. Normal initialization  $\mathcal{N} \left( 0, 4 \sqrt{\frac{2}{n_{input} + n_{output}}} \right)$  with input normalization.



(a) Value with closed-form policy



(c) Consumption with closed-form policy



(e) HJB error with closed-form policy

## Approximating the policy function

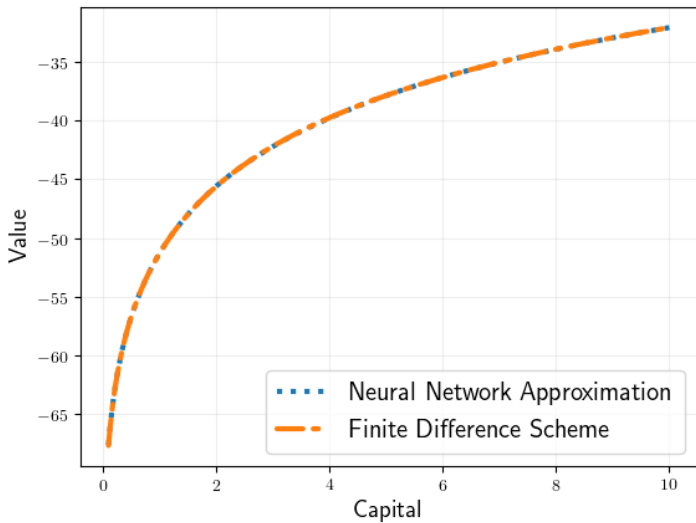
- Let us not use the closed-form consumption policy function but rather approximate said policy function directly with a policy neural network  $\tilde{C}(k; \Theta^C)$ .
- The new HJB error:

$$err_{HJB} = \rho \tilde{V}(k; \Theta^V) - U(\tilde{C}(k; \Theta^C)) - \frac{\partial \tilde{V}(k; \Theta^V)}{\partial k} [F(k) - \delta * k - \tilde{C}(k; \Theta^C)]$$

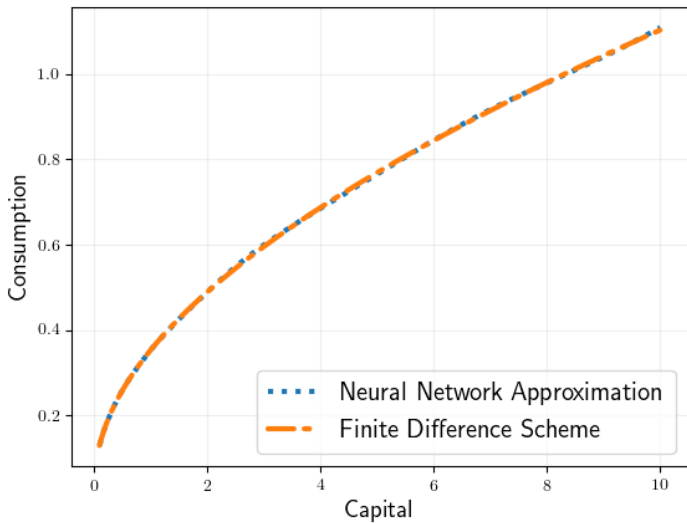
- Now we have a policy function error:

$$err_C = (U')^{-1} \left( \frac{\partial \tilde{V}(k; \Theta^V)}{\partial k} \right) - \tilde{C}(k; \Theta^C)$$

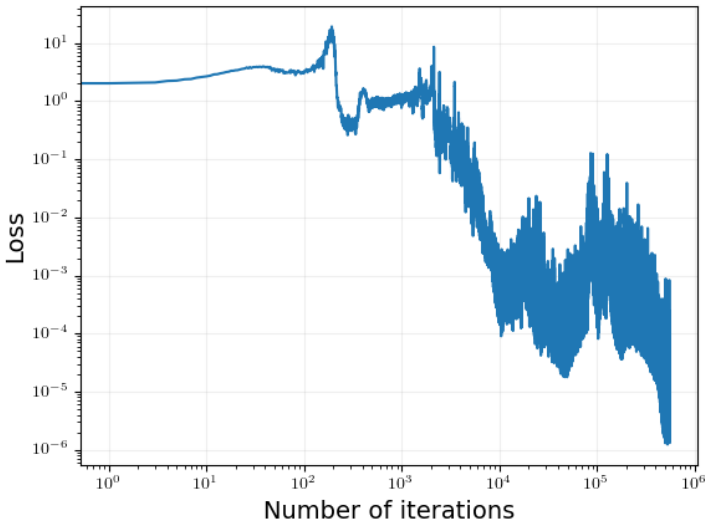




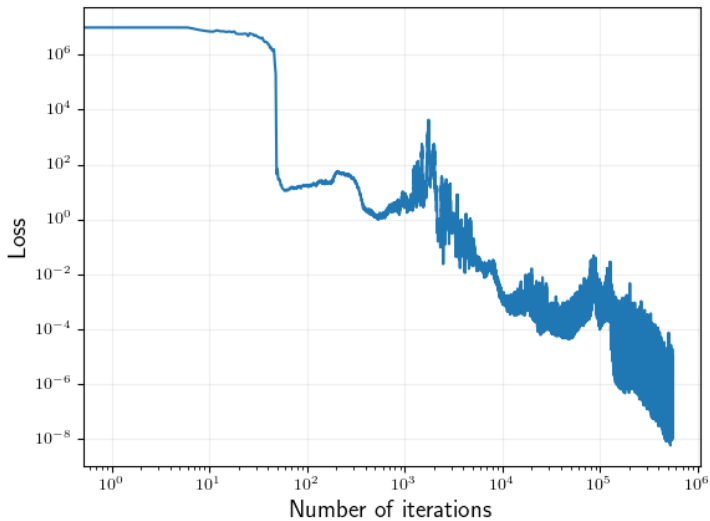
(b) Value with policy approximation



(d) Consumption with policy approximation



(f) HJB error with policy approximation

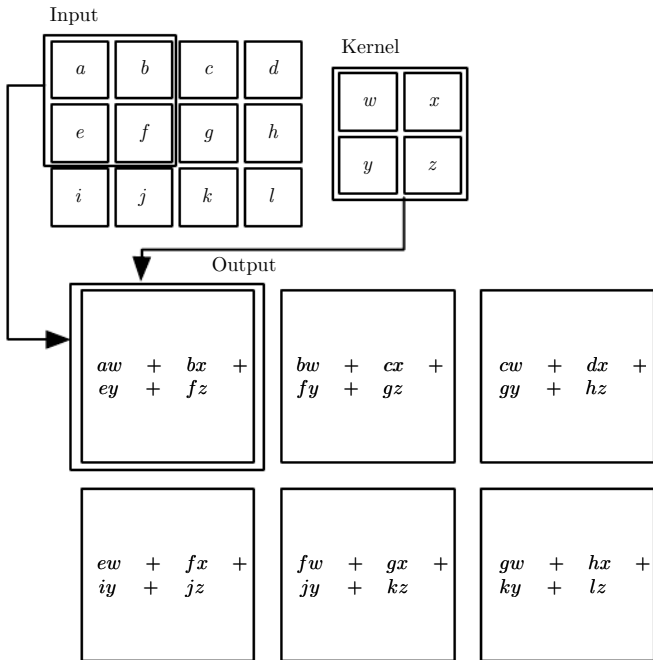


(g) Policy error with policy approximation

## Alternative ANNs

---

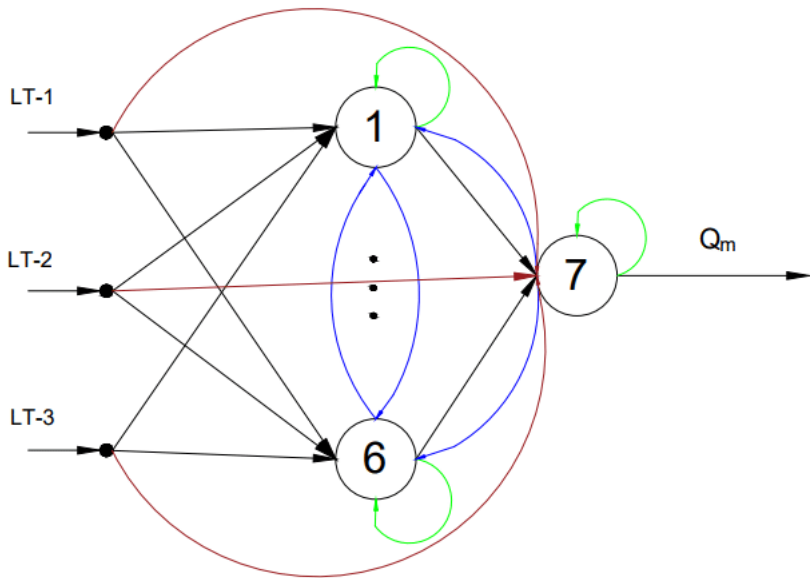
- Convolutional neural networks.
- Feedback ANN such as the Hopfield network.
- Self-organizing maps (SOM).
- ANN and reinforcement learning.



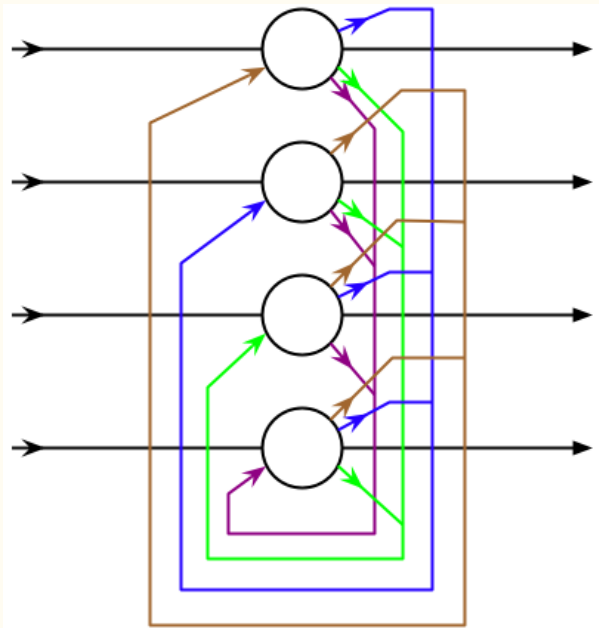
Input Layer

Hidden Layer

Output Layer







# Reinforcement learning

---

# Reinforcement learning

- Main idea: Algorithms that use training information that evaluates the actions taken instead of deciding whether the action was correct.
- Purely *evaluative feedback* to assess how good the action taken was, but not whether it was the best feasible action.
- Useful when:
  1. The dynamics of the state is unknown but simulation is easy: model-free vs. model-based reinforcement learning.
  2. Or the dimensionality is so high that we cannot store the information about the DP in a table.
- Work surprisingly well in a wide range of situations, although no methods that are guaranteed to work.
- Key for success in economic applications: ability to simulate fast (link with massive parallelization). Also, it complements very well with neural networks.

## Comparison with alternative methods

- Similar (same?) ideas are called approximate dynamic programming or neuro-dynamic programming.
- Traditional dynamic programming: we optimize over best feasible actions.
- Supervised learning: purely *instructive feedback* that indicates best feasible action regardless of action actually taken.
- Unsupervised learning: hard to use for optimal control problems.
- In practice, we mix different methods.
- Current research challenge: how do we handle associate behavior effectively?



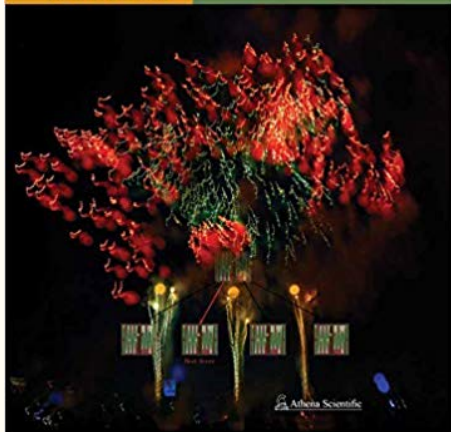
# Reinforcement Learning

An Introduction  
second edition

Richard S. Sutton and Andrew G. Barto

# Reinforcement Learning and Optimal Control

Dimitri P. Bertsekas



## Example: Multi-armed bandit problem

- You need to choose action  $a$  among  $k$  available options.
- Each option is associated with a probability distribution of payoffs.
- You want to maximize the expected (discounted) payoffs.
- But you do not know which action is best, you only have estimates of your value function (dual control problem of identification and optimization).
- You can observe actions and period payoffs.
- Go back to the study of “sequential design of experiments” by Thompson (1933, 1934) and Bellman (1956).





# Theory vs. practice

- You can follow two pure strategies:
  1. Follow *greedy* actions: actions with highest expected value. This is known as *exploiting*.
  2. Follow *non-greedy* actions: actions with dominated expected value. This is known as *exploring*.
- This should remind you of a basic dynamic programming problem: what is the optimal mix of pure strategies?
- If we impose enough structure on the problem (i.e., distributions of payoffs belong to some family, stationarity, etc.), we can solve (either theoretically or applying standard solution techniques) the optimal strategy (at least, up to some upper bound on computational capabilities).
- But these structures are too restrictive for practical purposes outside the pages of *Econometrica*.

# A policy-based method I

- Proposed by **Thathachar and Sastry (1985)**.
- A very simple method that uses the averages  $Q_n(a)$  of rewards  $R_i(a), i = \{1, \dots, n\}$ , actually received:

$$Q_n(a) = \frac{1}{n} \sum_{i=1}^{n-1} R_i(a)$$

- We start with  $Q_0(a) = 0$  for all  $k$ . Here (and later), we randomize among ties.
- We update  $Q_n(a)$  thanks to the nice recursive update based on linearity of means:

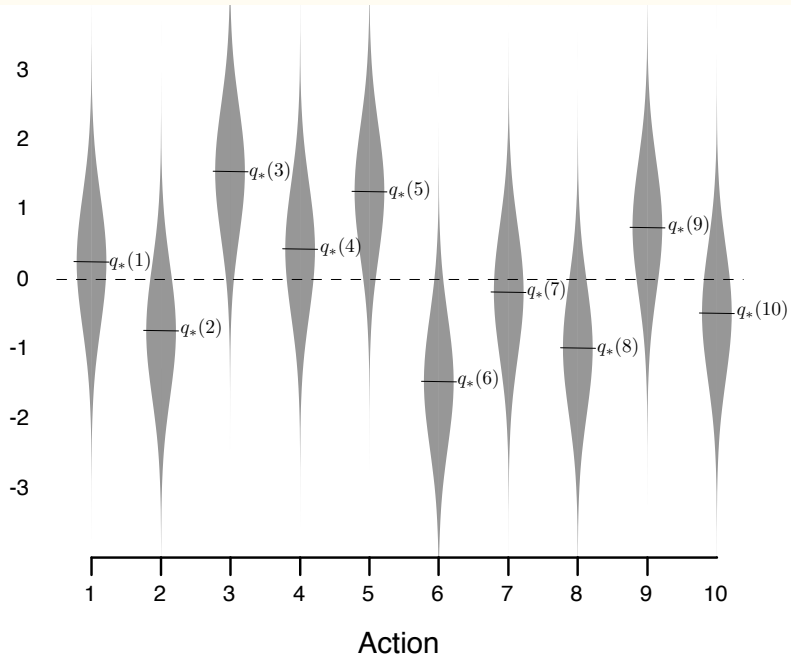
$$Q_{n+1}(a) = Q_n(a) + \frac{1}{n} [R_n(a) - Q_n(a)]$$

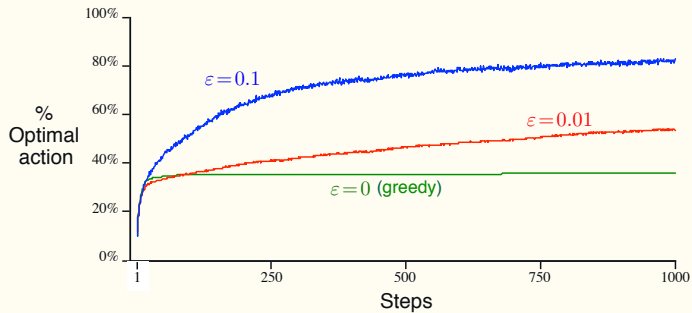
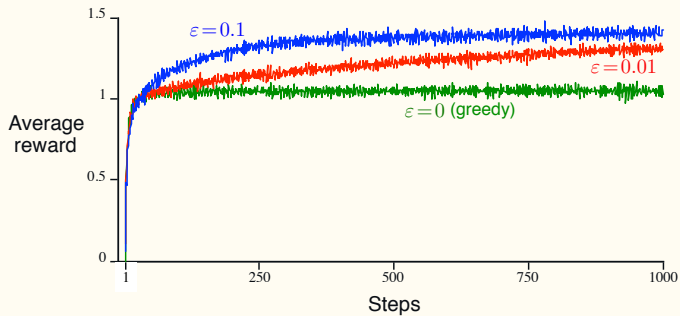
Averages of actions not picked are not updated.

## A policy-based method II

- How do we pick actions?
  1. Pure greedy method:  $\arg \max_a Q_t(a)$ .
  2.  $\epsilon$ -greedy method. Mixed best action with a random trembling.
- Easy to generalize to more sophisticated strategies.
- In particular, we can connect with genetic algorithms (AlphaGo).

Reward  
distribution





## A more general update rule

- Let's think about a modified update rule:

$$Q_{n+1}(a) = Q_n(a) + \alpha [R_n(a) - Q_n(a)]$$

for  $\alpha \in (0, 1]$ .

- This is equivalent, by recursive substitution, to:

$$Q_{n+1}(a) = (1 - \alpha)^n Q_1(a) + \alpha \sum_{i=1}^{n-1} \alpha (1 - \alpha)^{n-i} R_i(a)$$

- We can also have a time-varying  $\alpha_n(a)$ , but, to ensure convergence with probability 1 as long as:

$$\sum_{i=1}^{\infty} \alpha_n(a) = \infty$$

$$\sum_{i=1}^{\infty} \alpha_n^2(a) = \infty$$

## Improving the algorithm

- We can start with “optimistic”  $Q_0$  to induce exploration.

- We can implement an upper-confidence-bound action selection

$$\arg \max_a \left[ Q_n(a) + c \sqrt{\frac{\log n}{N_n(a)}} \right]$$

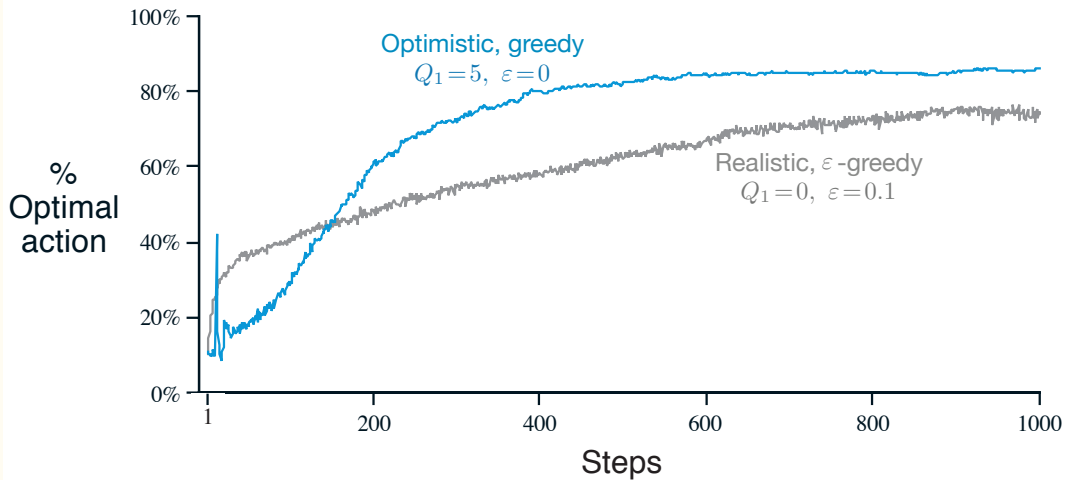
- We can have a gradient bandit algorithms based on a softmax choice:

$$\pi_n(a) = P(A_n = a) = \frac{e^{H_n(a)}}{\sum_{b=1}^k e^{H_n(b)}}$$

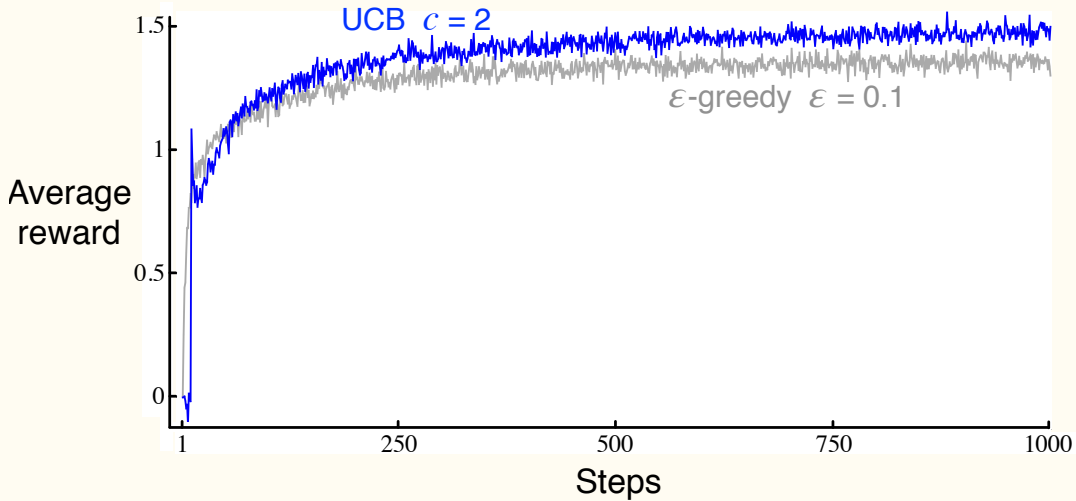
where

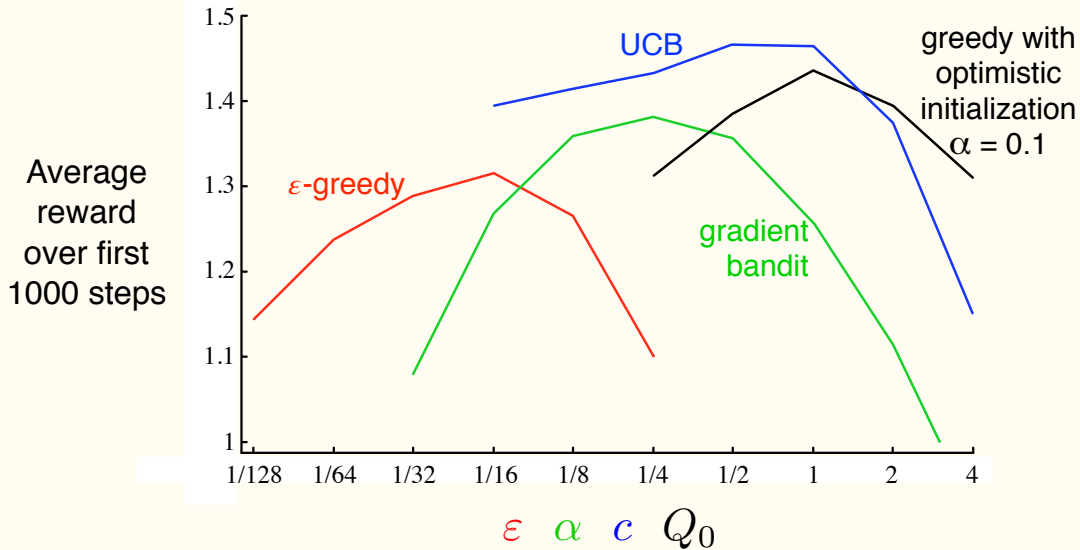
$$\begin{aligned} H_{n+1}(A_n) &= H_n(A_n) + \alpha(1 - \pi_n(A_n))(R_n(a) - \bar{R}_n) \\ H_{n+1}(a) &= H_n(a) - \alpha\pi_n(a)(R_n(a) - \bar{R}_n) \text{ for all } a \neq A_n \end{aligned}$$

This is a slightly hidden version of a stochastic gradient algorithm that we will see soon when we talk about deep learning.









## Other algorithms

- Monte Carlo prediction.
- Temporal-difference (TD) learning:

$$V^{n+1}(s_t) = V^n(s_t) + \alpha(r_{t+1} + \beta V^n(s_{t+1}) - V^n(s_t))$$

- SARSA  $\Rightarrow$  On-policy TD control:

$$Q^{n+1}(a_t, s_t) = Q^n(a_t, s_t) + \alpha(r_{t+1} + \beta Q^n(a_{t+1}, s_{t+1}) - Q^n(a_t, s_t))$$

- Q-learning  $\Rightarrow$  Off-Policy TD Control:

$$Q^{n+1}(a_t, s_t) = Q^n(a_t, s_t) + \alpha \left( r_{t+1} + \beta \max_{a_{t+1}} Q^n(a_{t+1}, s_{t+1}) - Q^n(a_t, s_t) \right)$$

- Value-based methods.
- Actor-critic methods.