

Spooky Boundaries at a Distance: Inductive Bias, Dynamic Models, and Behavioral Macro

Mahdi Ebrahimi Kahou¹ Jesús Fernández-Villaverde² Sebastián Gómez-Cardona³
Jesse Perla⁴ Jan Rosa⁴

–

¹Bowdoin College

²University of Pennsylvania

³Morningstar, Inc.

⁴University of British Columbia

Motivation, Question, and Contribution

In the long run, we are all dead—*J.M. Keynes, A Tract on Monetary Reform (1923)*

- Numerical solutions to dynamical systems are central to many quantitative fields in economics.
- Many dynamical systems in economics are **boundary value** problems:
 1. The boundary is at **infinity**.
 2. The values at the boundary are potentially **unknown**.
- Resulting from **forward looking** behavior of agents.
- Examples include the transversality and the no-bubble condition.
- Without them, the problems are **ill-posed** and have infinitely many solutions:
 - These forward-looking boundary conditions are a key limitation on increasing dimensionality.

Question

Question:

*Can we (economists and agents) **ignore** these long-run boundary conditions and still have accurate short/medium-run dynamics disciplined by the long-run conditions?*

1. **Yes**, it is possible to meet long-run boundary conditions **without** strictly enforcing them as a constraint on the model's dynamics.
 - We show how using Machine Learning (ML) methods achieve this.
 - This is due to the **inductive bias** of ML methods.
 - In this paper focusing on deep neural networks
2. We argue the inductive bias provides a foundation for modeling forward-looking behavioral agents with self-consistent expectations.
 - Easy to compute.
 - Provides short-run accuracy.

Background: Economic Models, Deep learning and inductive bias

Economic Models: functional equations

Many theoretical models can be written as functional equations:

- Economic object of interest: f , where $f : \mathcal{X} \rightarrow \mathcal{R} \subseteq \mathbb{R}^N$
 - e.g., asset price, investment choice, best-response, etc.
- Domain of f : \mathcal{X}
 - e.g. space of dividends, capital, opponents state or time in sequential models.
- The “Economics model” error: $\ell(x, f)$
 - e.g., Euler and Bellman residuals, equilibrium FOCs.

Then a **solution** is $f^* \in \mathcal{F}$ where $\ell(x, f^*) = \mathbf{0}$ for all $x \in \mathcal{X}$.

Example: one formulation of neoclassical growth

An example of a recursive case:

- Domain: $x = [k]$ and $\mathcal{X} = \mathbb{R}_+$.
- Solve for the optimal policy $k'(\cdot)$ and consumption function $c(\cdot)$: So $\psi : \mathbb{R} \rightarrow \mathbb{R}^2$ and $\mathcal{Y} = \mathbb{R}_+^2$.
- Residuals are the Euler equation and feasibility condition, so $\mathcal{R} = \mathbb{R}^2$:

$$\ell(\underbrace{\begin{bmatrix} k'(\cdot) & c(\cdot) \end{bmatrix}}_{\equiv f}, \underbrace{k}_{\equiv x}) = \underbrace{\begin{bmatrix} u'(c(k)) - \beta u'(c(k'(k))) (f'(k'(k)) + 1 - \delta) \\ f(k) - c(k) - k'(k) + (1 - \delta)k \end{bmatrix}}_{\text{model}}$$

- Finally, $f^* = [k'(\cdot), c(\cdot)]$ is a solution if it has zero residuals on domain \mathcal{X} .

Approximate solution: deep neural networks

1. Sample \mathcal{X} : $\mathcal{D} = \{x_1, \dots, x_N\}$
2. Pick a deep neural network $f_\theta(\cdot) \in \mathcal{H}(\theta)$:
 - θ : parameters for optimization (i.e., weights and biases).
3. To find an approximation for f solve:

$$\min_{\theta} \frac{1}{N} \sum_{x \in \mathcal{D}} \underbrace{\| \ell(x, f_\theta) \|_2^2}_{\text{Econ model error}}$$

- Deep neural networks are highly over-parameterized.
- Formally, $|\theta| \gg N$

Deep learning is **highly-overparameterized** $\mathcal{H}(\Theta)$ ($M \gg D$) class of functions.

- Example: one layer neural network, $\hat{\psi} : \mathbb{R}^Q \rightarrow \mathbb{R}$:

$$\hat{\psi}(x; \theta) = W_2 \cdot \sigma(W_1 \cdot x + b_1) + b_2$$

- $W_1 \in \mathbb{R}^{P \times Q}$, $b_1 \in \mathbb{R}^{P \times 1}$, $W_2 \in \mathbb{R}^{1 \times P}$, and $b_2 \in \mathbb{R}$.
- $\sigma(\cdot)$ is a nonlinear function applied element-wise (e.g., $\max\{\cdot, 0\}$).
- $\Theta \equiv \{b_1, W_1, b_2, W_2\}$ are the coefficients, in this example $M = PQ + P + P + 1$.
- Making it “deeper” by adding another “layer”: $\hat{\psi}(x; \theta) \equiv W_3 \cdot \sigma(W_2 \cdot \sigma(W_1 \cdot x + b_1) + b_2) + b_3$.

Over-parameterized interpolation

- Being over-parameterized ($|\theta| \gg N$), the optimization problem can have many solutions.
- Since individual θ are irrelevant it is helpful to think of optimization directly within \mathcal{H}

$$\min_{f_\theta \in \mathcal{H}} \frac{1}{N} \sum_{x \in \mathcal{D}} \|\ell(x, f_\theta)\|_2^2$$

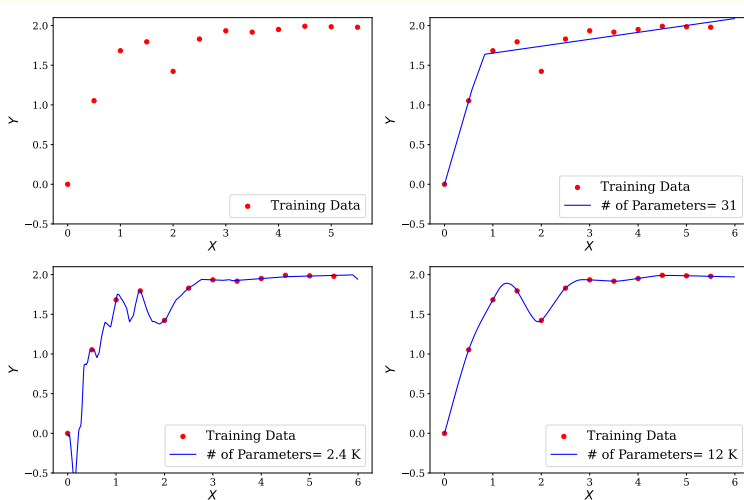
- But which f_θ ?
- **Mental model:** chooses min-norm interpolating solution for a (usually) unknown functional norm ψ

$$\begin{aligned} \min_{f_\theta \in \mathcal{H}} & \|f_\theta\|_\psi \\ \text{s.t. } & \ell(x, f_\theta) = 0, \quad \text{for all } x \in \mathcal{D} \end{aligned}$$

- That is what we mean by **inductive bias** (see Belkin, 2021 and Ma and Yang, 2021).
- Characterizing ψ (e.g., Sobolev norms or semi-norms?) is an active research area in ML.

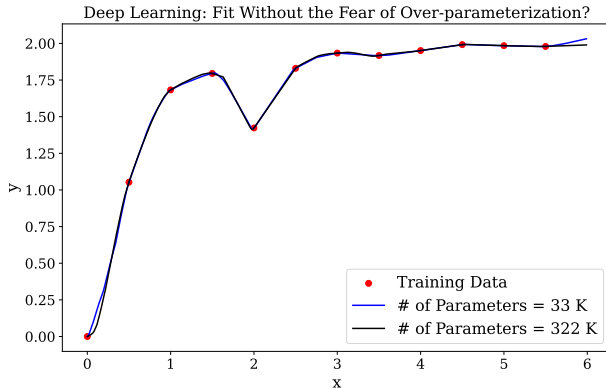
Over-parameterization and smooth interpolation

- Intuition: biased toward solutions which are flatter and have smaller derivatives

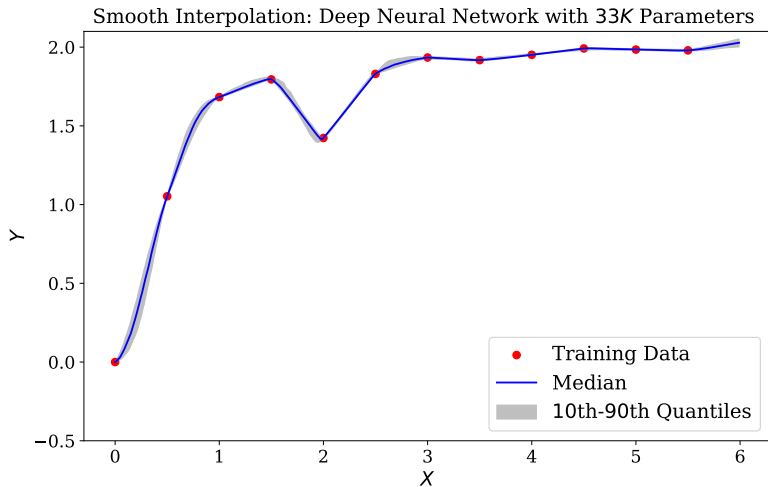


Deep Learning: “Fit Without Fear”?

- *“I remember my friend Johnny von Neumann used to say, with four parameters I can fit an elephant, and with five I can make him wiggle his trunk.”* Enrico Fermi
- *“The best way to solve the problem from practical standpoint is you build a very big system ... basically you want to make sure you hit the zero training error”* Ruslan Salakhutdinov



Deep Learning: random initialization and non-convex optimization



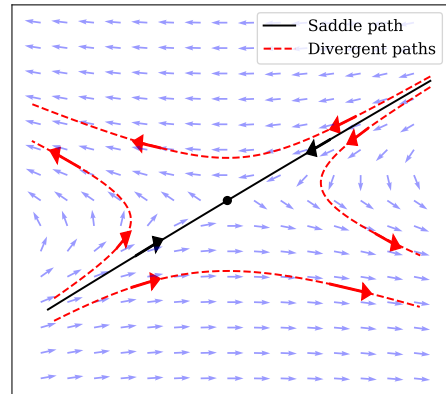
Intuition of the paper

- **Minimum-norm inductive bias:**

- Over-parameterized models (e.g., large neural networks) interpolate the train data.
- They are biased towards interpolating functions with smaller norms.
- So they don't like explosive functions.

- **Violation of economic boundary conditions:**

- Sub-optimal solutions diverge (explode) over time.
- This is due to the **saddle-path** nature of econ problems.
- The long-run boundary conditions rule out the explosive solutions.



Outline

To explore how we can ignore the long-run boundary conditions, we show deep learning solutions to

1. Classic linear-asset pricing model.
2. Sequential formulation of the neoclassical growth model.
3. Sequential formulation of the neoclassical growth model with non-concave production function.
4. Equivalent for a recursive formulation of the neoclassical growth model.

Linear asset pricing and the no-bubble condition

Linear asset pricing: setup

- The risk-neutral price, $p(t)$, of a claim to a stream of dividends, $y(t)$, is given by the recursive equation:

$$p(t) = y(t) + \beta p(t+1), \quad \text{for } t = 0, 1, \dots$$

- $\beta < 1$, and $y(t)$ is exogenous, $y(0)$ given.
- This is a two dimensional dynamical system with unknown initial condition $p(0)$. This problem is **ill-posed**.
- A family of solutions

$$p(t) = \underbrace{p_f(t)}_{\text{fundamentals}} + \underbrace{\zeta \left(\frac{1}{\beta} \right)^t}_{\text{explosive bubble}}$$

- $p_f(t) \equiv \sum_{\tau=0}^{\infty} \beta^{\tau} y(t+\tau)$. Each solution corresponds to a different $\zeta > 0$.

Linear asset pricing: the long-run boundary condition

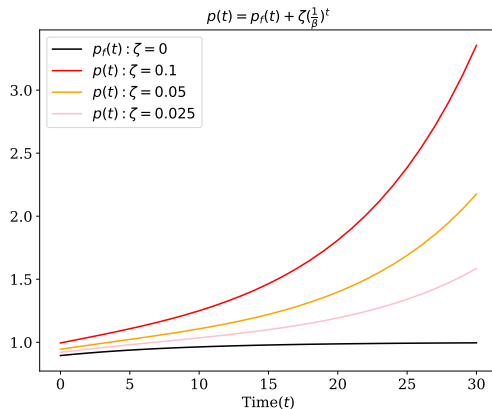
- Long-run boundary condition that rule out the explosive bubbles and chooses $\zeta = 0$

$$\lim_{t \rightarrow \infty} \beta^t p(t) = 0.$$

- Any norm that preserve monotonicity, like L_p and Sobolev (semi-)norms

$$\min_{\zeta \geq 0} \|p\|_{\psi} = \|p_f\|_{\psi}$$

- Ignoring the no-bubble condition and using a deep neural network provides an accurate approximation for $p_f(t)$.



Linear asset pricing: numerical method

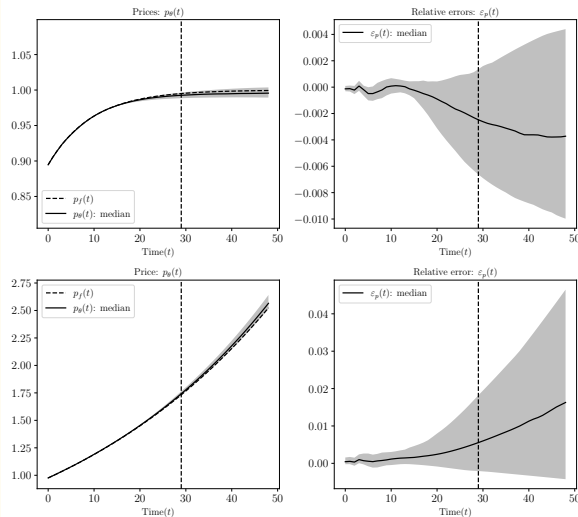
- Sample for time: $\mathcal{D} = \{t_1, \dots, t_N\}$.
- Generating the dividend process: $y(t+1) = c + (1+g)y(t)$, given $y(0)$.
- An over-parameterized neural network $p_\theta(t)$, **ignore** the non-bubble condition and solve

$$\min_{\theta} \frac{1}{N} \sum_{t \in \mathcal{D}} [p_\theta(t) - y(t) - \beta p_\theta(t+1)]^2$$

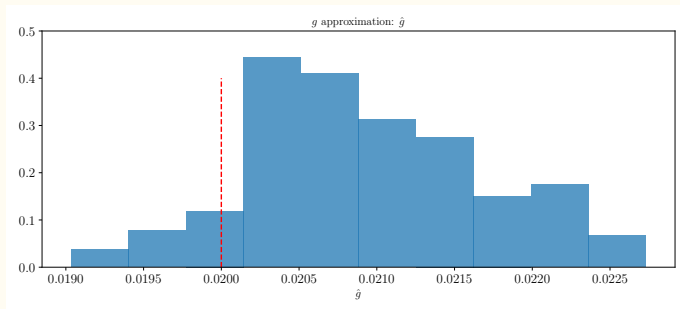
- This minimization should provide an accurate short- and medium-run approximation for price based on the fundamentals $p_f(t)$.

Linear asset pricing: results

- Two cases: $g < 0$ and $g > 0$.
- Relative errors: $\varepsilon_p(t) \equiv \frac{p_\theta(t) - p_f(t)}{p_f(t)}$.
- for $g > 0$: $p_\theta(t) = e^{\phi t} N N_\theta(t)$, ϕ is “learnable”.
- Results for 100 different seeds (initialization of the parameters):
 - important for non-convex optimizations.
- Very accurate short- and medium-run approximation.



Learning the growth rate



- $\hat{g} = e^{\phi} - 1$.
- Slightly biased due to small sample size, i.e., $\mathcal{D} = \{0, 1, \dots, 29\}$.

Sequential neoclassical growth model and the transversality condition

Neoclassical growth model: setup

- Total factor productivity $z(t)$ exogenously given, capital $k(t)$ with given $k(0)$, consumption $c(t)$, production function $f(\cdot)$, depreciation rate $\delta < 1$, discount factor β :

$$\underbrace{k(t+1) = z(t)^{1-\alpha} f(k(t)) + (1-\delta)k(t) - c(t)}_{\text{feasibility constraint}},$$

$$\underbrace{c(t+1) = \beta c(t) [z(t+1)^{1-\alpha} f'(k(t+1)) + 1 - \delta]}_{\text{Euler equation}}.$$

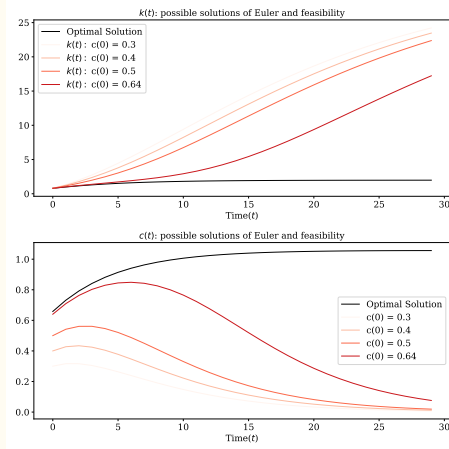
- This is a three dimensional dynamical system with unknown initial condition $c(0)$. This problem is **ill-posed**.
- A family of solutions, each solution corresponds to a different $c(0)$. Only one of them is the optimal solution.

Neoclassical growth model: the long-run boundary condition

- To rule out sub-optimal solutions, transversality condition

$$\lim_{t \rightarrow \infty} \beta^t \frac{k(t+1)}{c(t)} = 0.$$

- Using a deep neural network and ignoring the transversality condition provides a an accurate approximation for the optimal capital path.



Neoclassical growth model: numerical method

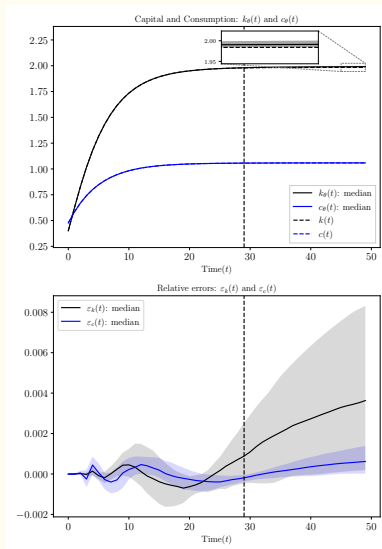
- Sample for time: $\mathcal{D} = \{t_1, \dots, t_N\}$.
- Generating the TFP process: $z(t+1) = (1+g)z(t)$, given $z(0)$.
- A over-parameterized neural network $k_\theta(t)$,
- Given $k_\theta(t)$, define the consumption function $c(t; k_\theta) = z(t)^{1-\alpha} f(k_\theta(t)) + (1-\delta)k_\theta(t) - k_\theta(t+1)$
- **Ignore** the transversality condition and solve

$$\min_{\theta \in \Theta} \left[\frac{1}{N} \sum_{t \in \mathcal{D}} \underbrace{\left(\frac{c(t+1; k_\theta)}{c(t; k_\theta)} - \beta [z(t+1)^{1-\alpha} f'(k_\theta(t+1)) + (1-\delta)] \right)^2}_{\text{Euler residuals}} + \underbrace{\left(k_\theta(0) - k_0 \right)^2}_{\text{Initial condition residual}} \right]$$

- This minimization should provide an accurate short- and medium-run approximation for the optimal capital and consumption path.

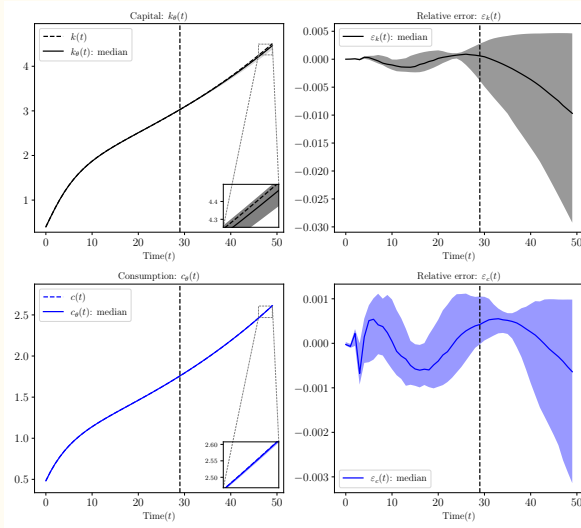
Neoclassical growth model, no TFP growth: results

- $g = 0$, $z(0) = 1$.
- $\varepsilon_k(t) \equiv \frac{k_\theta(t) - k(t)}{k(t)}$, and $\varepsilon_c(t) \equiv \frac{c(t; k_\theta) - c(t)}{c(t)}$
- Benchmark solution: value function iteration.
- Results for 100 different seeds.
- Very accurate short- and medium-run approximation.



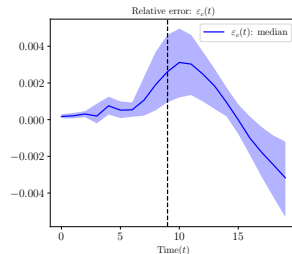
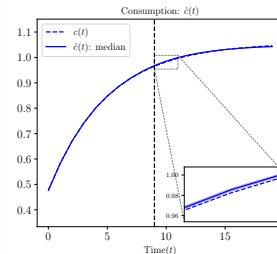
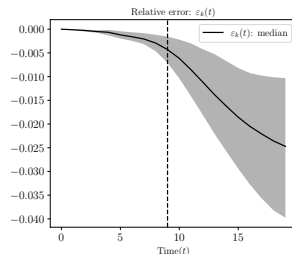
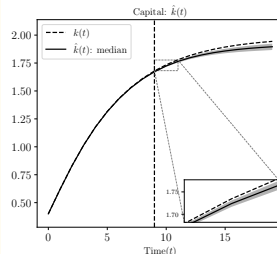
Neoclassical growth model with TFP growth: results

- $g > 0$ and $z(0) = 1$.
- $k_{\theta}(t) = e^{\phi t} N N_{\theta}(t)$, ϕ is "learnable".
- Results for 100 different seeds.
- Very accurate short- and medium-run approximation.



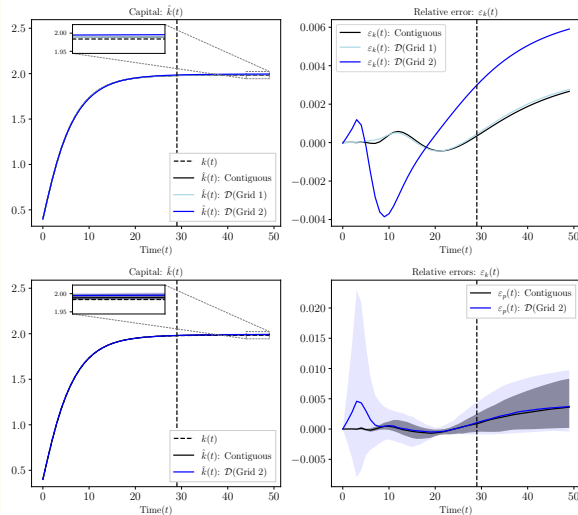
But, seriously “in the long run, we are all dead”

- So far, we have used long time-horizon $\mathcal{D} = \{0, 1, \dots, 29\}$.
- In other methods, choosing the time-horizon T is a challenge:
 - Too large \rightarrow accumulation of errors, and numerical instability. We don't have that problem.
 - Too small \rightarrow convergence to the steady state too quickly.
- An accurate short-run solution, even for a medium-sized T .



Do we need a dense and contiguous grid?

- We have used a dense $\mathcal{D} = \{0, 1, \dots, 29\}$.
- What if
 - $\mathcal{D}(\text{Grid 1}) = \{0, 1, 2, 4, 6, 8, 12, 16, 20, 24, 29\}$
 - $\mathcal{D}(\text{Grid 2}) = \{0, 1, 4, 8, 12, 18, 24, 29\}$
- An accurate short-run solution, even for a sparse grid.

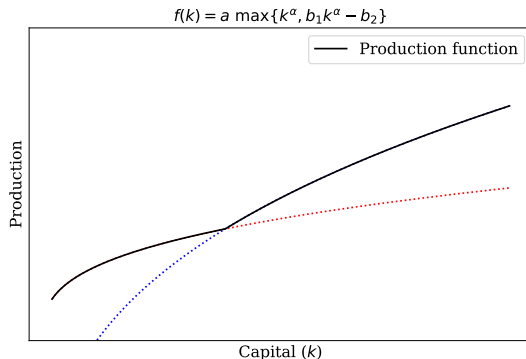


Neoclassical growth model: multiple steady-states and hysteresis

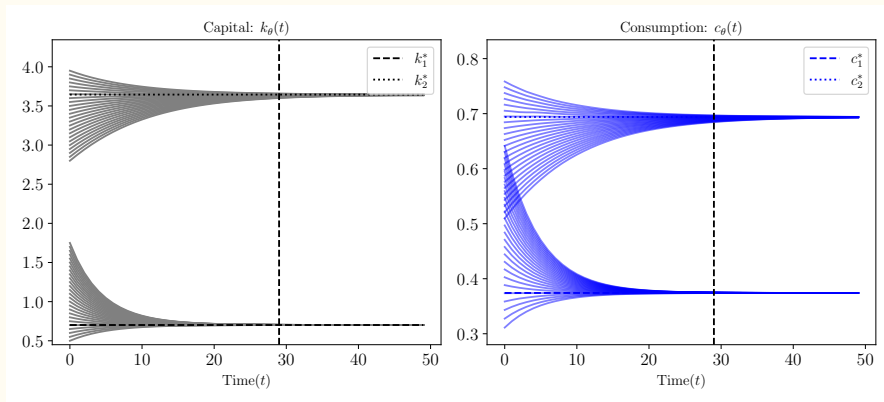
- When there are multiple steady states with saddle-path stability, each with its domain of attraction:
 - Can the inductive bias detect there are multiple basins of attraction?
 - How does the inductive bias move us toward the correct steady state for a given initial condition?
- Consider a non-concave production function:

$$f(k) \equiv a \max\{k^\alpha, b_1 k^\alpha - b_2\}$$

- Two steady-states k_1^* and k_2^* .
- The same numerical procedure.



Neoclassical growth model with non-concave production function: results



- Different initial conditions in $k_0 \in [0.5, 1.75] \cup [2.75, 4]$.
- In the vicinity of k_1^* and k_2^* the paths converge to the right steady-states.

**Deep learning is not the only
option**

Deep learning is not the only option: kernels

- Deep learning might be too “spooky”.
- We can use kernels methods, $K(\cdot, \cdot)$, instead of neural networks and control the RKHS norms.
- Focusing on continuous time equivalent of these problems.
- The same results, theoretical guarantees, very fast and robust.

How Inductive Bias in Kernel Methods Aligns with Optimality in Economic Dynamics

SUBMISSION 764

This paper examines the alignment of inductive biases in machine learning (ML), such as kernel machines, with structural models of economic dynamics. Unlike dynamical systems found in physical and life sciences, economic models are often specified by differential equations with a mixture of easy-to-enforce initial conditions and hard-to-enforce infinite horizon boundary conditions (e.g. transversality and no-ponzi-scheme conditions). We investigate algorithms using ridgeless kernel methods trained to fulfill the differential equations without explicitly fulfilling the boundary conditions. Our findings provide theoretical guarantees for cases where the inductive biases of these ML models are sufficient conditions to fulfill the infinite-horizon conditions. We then provide empirical evidence that ridgeless kernel methods are not only theoretically sound with respect to economic assumptions, but may even dominate classic algorithms in low to medium dimensions.

CONTENTS

Abstract	0
Contents	0
1 Introduction	1
2 Related Work	2
3 Setup	3
4 Method	5
5 Results	7
5.1 Asset Pricing	8
5.2 Neoclassical Growth Model	8
5.3 Neoclassical Growth Model with Multiple Steady-States	10
5.4 Other Examples	11
6 Conclusion	11

Optimal control framework

Consider the following problem arising in optimal control:

$$\dot{\mathbf{x}} = \mathbf{F}(\mathbf{x}(t), \mathbf{y}(t))$$

$$\dot{\boldsymbol{\mu}} = r\boldsymbol{\mu}(t) - \boldsymbol{\mu}(t) \odot \mathbf{G}(\mathbf{x}(t), \boldsymbol{\mu}(t), \mathbf{y}(t))$$

$$\mathbf{0} = \mathbf{H}(\mathbf{x}(t), \boldsymbol{\mu}(t), \mathbf{y}(t))$$

$$\mathbf{x}(0) = \mathbf{x}_0$$

- State variables $\mathbf{x}(t) \in \mathbb{R}^M$, initial condition \mathbf{x}_0 ; co-state variables $\boldsymbol{\mu}(t) \in \mathbb{R}^M$; jump variables $\mathbf{y}(t) \in \mathbb{R}^P$
- This problem is **ill-posed** and can have infinitely many solutions.

Transversality condition: an asymptotic boundary condition

$$\lim_{t \rightarrow \infty} e^{-rt} \mathbf{x}(t) \odot \boldsymbol{\mu}(t) = \mathbf{0}$$

- The transversality condition is an asymptotic boundary condition.
- We typically assume a finite time horizon T and shoot for the finite steady state \mathbf{x}^* , $\boldsymbol{\mu}^*$, and \mathbf{y}^* .
- This approach is straightforward in low dimensions but becomes significantly more challenging in high-dimensional settings.

Optimal control framework: an example, Ramsey–Cass–Koopmans model

- Classic Ramsey–Cass–Koopmans

$$\dot{k}(t) = f(k(t)) - c(t) - \delta k(t)$$

$$\dot{\mu}(t) = r\mu(t) - \mu(t)[f'(k(t)) - \delta]$$

$$0 = c(t)\mu(t) - 1$$

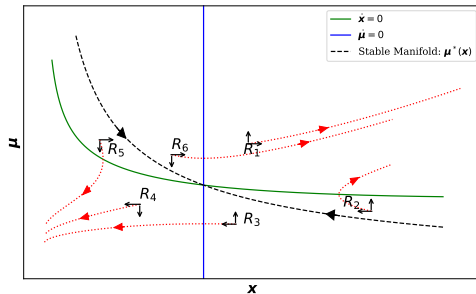
$$k(0) = k_0$$

$$0 = \lim_{t \rightarrow \infty} e^{-rt} k(t) \mu(t)$$

- $f(\cdot)$ is the production function, r discount rate, and δ is the depreciation.

What does the violation of the transversality condition look like?

- All paths solve the ordinary differential equations and the algebraic equation.
- The solutions that violate the transversality condition $\lim_{t \rightarrow \infty} \dot{\mu} = \infty$ and $\lim_{t \rightarrow \infty} \mu = \infty$
 - Diverges faster than e^{rt} .



Kernel approximation

Approximating the derivatives with a kernel:

$$\begin{aligned}\hat{\mathbf{x}}(t) &= \mathbf{x}_0 + \int_0^t \hat{\mathbf{x}}(\tau) d\tau, & \hat{\mu}(t) &= \hat{\mu}_0 + \int_0^t \hat{\mu}(\tau) d\tau, & \hat{\mathbf{y}}(t) &= \hat{\mathbf{y}}_0 + \int_0^t \hat{\mathbf{y}}(\tau) d\tau, \\ \hat{\mathbf{x}}(t) &= \sum_{j=1}^N \alpha_j^x K(t, t_j), & \hat{\mu}(t) &= \sum_{j=1}^N \alpha_j^\mu K(t, t_j), & \hat{\mathbf{y}}(t) &= \sum_{j=1}^N \alpha_j^y K(t, t_j)\end{aligned}$$

- \mathbf{x}_0 is given.
- $\hat{\mu}_0$, $\hat{\mathbf{y}}_0$, α^x , α^μ , and α^y are learnable parameters.
- $K(\cdot, \cdot)$ is the kernel.

Approximate solution: Algorithm

$$\begin{aligned} \min_{\substack{\hat{\mathbf{x}}(t) \in \mathcal{H}^M, \hat{\boldsymbol{\mu}}(t) \in \mathcal{H}^M, \\ \hat{\mathbf{y}}(t) \in \mathcal{H}^P}} \quad & \left(\sum_{m=1}^M \|\hat{\dot{\mathbf{x}}}^{(m)}\|_{\mathcal{H}}^2 + \sum_{m=1}^M \|\hat{\dot{\boldsymbol{\mu}}}^{(m)}\|_{\mathcal{H}}^2 \right) \\ \text{s.t.} \quad & \hat{\dot{\mathbf{x}}} = \mathbf{F}(\hat{\mathbf{x}}(t), \hat{\mathbf{y}}(t)) \\ & \hat{\dot{\boldsymbol{\mu}}} = r\hat{\boldsymbol{\mu}}(t) - \hat{\boldsymbol{\mu}}(t) \odot \mathbf{G}(\hat{\mathbf{x}}(t), \hat{\boldsymbol{\mu}}(t), \hat{\mathbf{y}}(t)) \\ & \mathbf{0} = \mathbf{H}(\hat{\mathbf{x}}(t), \hat{\boldsymbol{\mu}}(t), \hat{\mathbf{y}}(t)) \end{aligned}$$

- The objective function penalizes explosive paths.
- Constraints solve the "first order conditions".

Application: Growth with human and physical capital, a mid-size problem

$$\dot{k}(t) = i_k(t) - \delta_k k(t),$$

$$\dot{h}(t) = i_h(t) - \delta_h h(t),$$

$$\dot{\mu}_k(t) = r\mu_k(t) - \mu_k(t) [f_k(k(t), h(t)) - \delta_k],$$

$$\dot{\mu}_h(t) = r\mu_h(t) - \mu_h(t) [f_h(k(t), h(t)) - \delta_h],$$

$$0 = \mu_k(t)c(t) - 1$$

$$0 = \mu_k(t) - \mu_h(t)$$

$$0 = f(k(t), h(t)) - c(t) - i_k(t) - i_h(t),$$

for given initial conditions $k(0) = k_0$, $h(0) = h_0$, and two transversality conditions

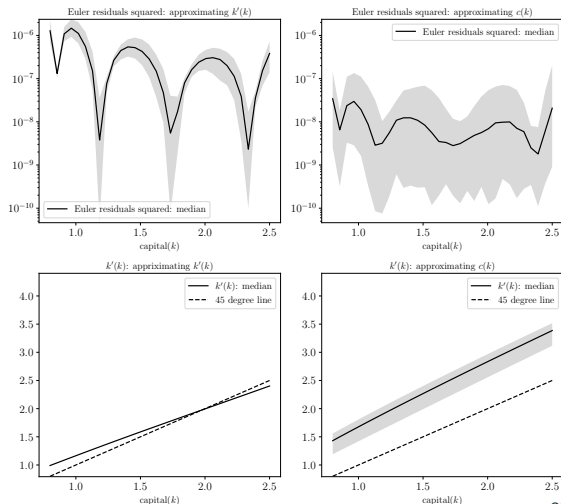
$$0 = \lim_{t \rightarrow \infty} e^{-rt} k(t) \mu_k(t),$$

$$0 = \lim_{t \rightarrow \infty} e^{-rt} h(t) \mu_h(t).$$

- Sequential models:
 - Shorter time-horizons.
 - Misspecification of growth.
- Recursive neoclassical growth model
 - Accurate short- and medium-run dynamics.
 - Accurate solutions even with TFP growth.
 - Deep learning solutions can go very wrong
 - We should use the information in the transversality condition to know what to approximate.

Deep learning solutions can be misleading: approximating capital vs. consumption

- Capital k is the state variable.
- Two options: approximating capital policy $k'_\theta(k)$ or $c_\theta(k)$
- left panels: results for $k'_\theta(k)$ approximation.
- Right panels: results for $c_\theta(k)$ approximation.
- Only the left panel results are correct. $k'_\theta(k)$ has a fixed point at the right steady state.
- However, the wrong solution has lower Euler residuals.



- Short- and medium-run accurate solutions can be obtained **without** strictly enforcing the long-run boundary conditions on the model's dynamics.
- Long-run (**global**) conditions can be replaced with appropriate regularization (**local**) to achieve optimal solutions, hence the title of the paper.
- Inductive bias provides a foundation for modeling forward-looking behavioral agents with self-consistent expectations.

Discussion: where to go from here?

- Can inductive bias/regularization be thought of as an equilibrium selection device?
 - In this paper it is used to select solutions.
- This method (mostly the kernel method) can be used for sampling high-dimensional state spaces when there is stochasticity.
 - Solve the deterministic in short-run and use the points as sample of the state-space.
 - Then solve the stochastic problem.

Appendix

Deep Learning: random initialization and non-convex optimization

