

Estimating Nonlinear Heterogeneous Agent Models with Neural Networks

Hanno Kase, Leonardo Melosi, Matthias Rottner

September 4, 2025

Introduction

- Over last three decades, researchers have focused on integrating heterogeneity in macroeconomic models.
- **DGE Models: Enables study of distributional issues, economic fluctuations, and stabilization policies in a micro-founded framework.**
- Accounting for these non-linear dynamics is essential to understand macroeconomic events, including recurring and prolonged periods at Zero Lower Bound (ZLB), deep recessions and recent rise in inflation.
- This paper introduces Neural-Network (NN) based method for solving and estimating models with heterogeneous agents, fully incorporating non-linear dynamics.

A neural-network framework for estimation

- Begin with a primer on neural networks.
- How NN can be used efficiently to characterize policy functions of a model treating parameters as psedo-state variables.
- How NN can be trained to characterize the deterministic steady state of the model.
- Training the NN particle filter and estimating likelihood functions.
- Solving Models:
 - Model #1: small scale linearized DSGE model
 - Model #2: Nonlinear RANK model
 - Model #3: Nonlinear HANK model

What are neural networks?

Consider a function:

$$Y = \psi(X)$$

where

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}, \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

are column vectors.

This function can then be approximated by a neural network as:

$$Y \approx \psi_{\text{NN}}(X \mid W) \tag{1}$$

with W denoting the weights that define the NN.

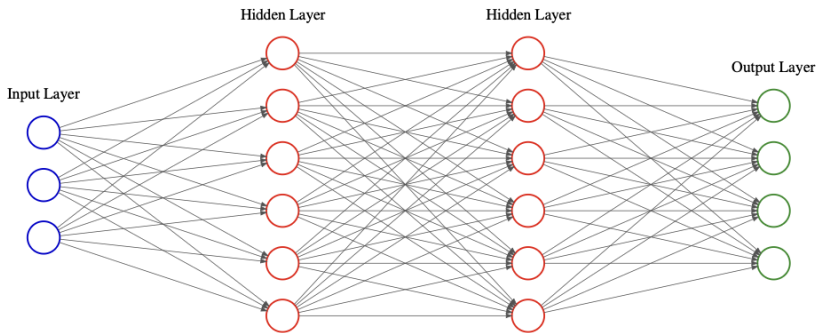


Figure 1: Illustration of a neural network. An illustrative example of a NN with three inputs (blue circles), two hidden layers (the two rows of red circles) with six neurons (the red circles) each, and four outputs (the green circles).

Neural networks

A single neuron is an affine function with:

- inputs $x_1, x_2 \dots x_M$
- weights $w_1, w_2 \dots w_M$
- bias w_0
- nonlinear activation function $g(\cdot)$

From all of this, we get a single output \tilde{y}

$$\tilde{y} = g \left(w_0 + \sum_{i=1}^M w_i x_i \right) \quad (2)$$

Neural networks

In matrix notation, each layer of the NN can be written as

$$y = g(\mathcal{A}^h(x)) \quad \text{where} \quad \mathcal{A}^h(x) \equiv W^h x \quad (3)$$

x , y and W are all vectors now.

Optimization of the weights depend on the back-propagation algorithm, which computes all partial derivatives of the loss with respect to the weights of the NN.

Accuracy of the NN can be achieved by adjusting the learning rate during the training.

Neural-network solution method

Consider a large class of DGE models by the following transition equations:

$$\mathbf{S}_t = f(\mathbf{S}_{t-1}, \mathbf{v}_t | \Theta) \quad (4)$$

Solving this model gives O policy functions that map state variables to a set of control variables ψ_t for given values of the model parameters

$$\Theta : \psi_t = \psi(\mathbf{S}_t | \Theta)$$

These policy functions satisfy set of K equilibrium conditions derived from the model equations

$$F(\psi(\mathbf{S}_t | \Theta)) = 0 \quad (5)$$

Extended policy functions:

The parameters to be estimated are considered as pseudo-state variables when approximating policy functions

$$\psi_t = \psi(\mathbb{S}_t, \tilde{\Theta} | \bar{\Theta}) \quad (6)$$

Θ is split into:

- $\tilde{\Theta}$ - set of parameters to estimate
- $\bar{\Theta}$ - set of parameters that are fixed in estimation

The transition equation then becomes:

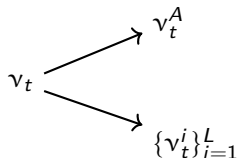
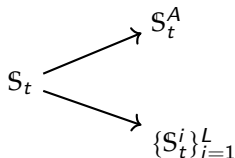
$$\mathbb{S}_t = f(\mathbb{S}_{t-1}, \mathbf{v}_t, \tilde{\Theta} | \bar{\Theta}) \quad (7)$$

Neural-network solution method

Agents' heterogeneity:

Discretize the model population by keeping track of a finite number L of agents.

Transition equation can be written in the same form as in Equation (4), but now deal with both individuals and aggregate levels.



Neural-network solution method

Two sets of policy functions to approximate:

- Aggregate policy functions: $\psi^A(\cdot)$
- Individual policy function: $\psi^I(\cdot)$

The extended forms would look like:

$$\psi_t^A = \psi^A(S_t, \tilde{\Theta} | \bar{\Theta}) \quad \text{and} \quad \psi_t^i = \psi^I(S_t^i, S_t, \tilde{\Theta} | \bar{\Theta}) \quad (8)$$

Then, the number of variables we have:

- $S = S^i \times L + S^A$
- $U = U^i \times L + U^A$
- $O = O^i \times L + O^A$

Neural-network solution method

Neural networks and training:

Set up two NNs $\psi_{NN}^I(\cdot)$ and $\psi_{NN}^A(\cdot)$ to approximate individual and aggregate policy functions $\psi^I(\cdot)$ and $\psi^A(\cdot)$.

Train the two NNs by choosing their weights W to minimize weighted squared residual error for the equilibrium conditions:

$$\Phi^L = \sum_{k=1}^K \alpha_k [F_k(\psi_{NN}(S_t, \tilde{\Theta}|\bar{\Theta}))]^2 \quad (9)$$

In this framework, the number of optimality conditions is as much as the number of policy functions. So,

$$K \geq O^I \times L + O^A$$

Neural-network solution method

For faster computations in GPUs, the authors have also proposed the following methods (parallelization):

Two NNs are trained on a batch with size B

- B pseudo-randomly drawn vectors, $\mathbf{S}_{t,b}$
- parameters to estimate: $\tilde{\Theta}_b$

In each step, minimize the loss functions Φ^L averaged across batch size B .

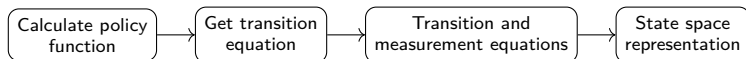
$$\bar{\Phi}^L = \frac{1}{B} \sum_{b=1}^B \sum_{k=1}^K \alpha_k [F_k(\psi_{NN}(\mathbf{S}_{t,b}, \tilde{\Theta}_b | \bar{\Theta}))]^2 \quad (10)$$

Neural Network Particle Filter

In DSGE models, likelihood functions are calculated in two steps:

1. Solution step:

The model's policy functions are characterized here.



2. Evaluation step:

Having a filter on state space representation can give us the likelihood function.

Work of a filter: models one step ahead predictive densities over available sample period

Neural Network Particle Filter

The filter looks like:

$$p_{\bar{\Theta}}(y_t | y_{t-1}, \dots, y_1, \tilde{\Theta})$$

The product of model's predictive densities in every period $y \in 1, \dots, T$ is the likelihood function:

$$\mathcal{L}_{\bar{\Theta}} = \prod_{t=1}^T p_{\bar{\Theta}}(y_t | y_{t-1}, \dots, y_1, \tilde{\Theta}) \quad (11)$$

Traditionally, we require MC filters (Fernandez-Villaverde and Rubio-Ramirez, 2007), and even in this case, this is used as a benchmark

Training the NN particle filter

NN particle filter can be trained in three steps:

- **Step 1:** Obtain several 1000s of quasi-random draws from the parameter space.
- **Step 2:** Execute the standard particle filter to derive likelihoods at each of these draws
- **Step 3:** NN particle filter is trained to minimize the MSE between predicted and actual likelihood values.

Solving Model #1

Solving a linearized small scale DSGE model:

There's a NK-DSGE model with a TFP shock. Model is log-linearized around the unique steady-state equilibrium, which results in following equations:

$$\hat{Y}_t = E_t \hat{Y}_{t+1} - \sigma^{-1} \left(\hat{R}_t - E_t \hat{\Pi}_{t+1} \right), \quad (12)$$

$$\hat{\Pi}_t = \kappa \left(\hat{Y}_t - \hat{Y}_t^* \right) + \beta E_t \hat{\Pi}_{t+1}, \quad (13)$$

$$\hat{R}_t = \phi_{\Pi} \hat{\Pi}_t + \phi_Y \hat{X}_t, \quad (14)$$

$$\hat{Y}_t^* = \omega \hat{A}_t, \quad (15)$$

$$\hat{A}_t = \rho_A \hat{A}_{t-1} + \sigma_A \epsilon_t^A, \quad (16)$$

Solving Model #1

Here are some definitions:

- Output: $\hat{Y}_t = \frac{Y_t - Y}{Y}$
- Inflation: $\hat{\Pi}_t = \Pi_t - \Pi$
- Nominal interest rate: $\hat{R}_t = R_t - R$
- Output in the flex-price economy: $\hat{Y}_t^* = \frac{Y_t^* - Y^*}{Y^*}$
- TFP: $\hat{A}_t = \frac{A_t - A}{A}$
- Output gap: $\hat{X}_t = \hat{Y}_t - \hat{Y}_t^*$
- Shock: $\epsilon_t^A \sim \mathcal{N}(0, 1)$

Solving Model #1

After some manipulations in the equations, they become:

$$\hat{X}_t = E_t \hat{X}_{t+1} - \sigma^{-1} \left(\phi_{\Pi} \hat{\Pi}_t + \phi_Y \hat{X}_t - E_t \hat{\Pi}_{t+1} - \hat{R}_t^* \right), \quad (17)$$

$$\hat{\Pi}_t = \kappa \hat{X}_t + \beta E_t \hat{\Pi}_{t+1}, \quad (18)$$

$$\hat{R}_t^* = \rho_A \hat{R}_{t-1}^* + \sigma(\rho_A - 1) \omega \sigma_A \epsilon_t^A, \quad (19)$$

where \hat{R}_t^* is the natural rate of interest which follows an exogenous process derived from the TFP process.

Solving Model #1

The analytical solution for the output gap and the inflation gap is as follows:

$$\hat{X}_t = \frac{(1 - \beta\rho_A)}{(\sigma(1 - \rho_A) + \theta_Y)(1 - \beta\rho_A) + \kappa(\theta_\Pi - \rho_A)} \hat{R}_t^*, \quad (20)$$

$$\hat{\Pi}_t = \frac{\kappa}{(\sigma(1 - \rho_A) + \theta_Y)(1 - \beta\rho_A) + \kappa(\theta_\Pi - \rho_A)} \hat{R}_t^*. \quad (21)$$

The analytical solution for the output gap and the inflation gap depends both on the state variable \hat{R}_t^* and the parameters.

Mapping of Model #1 to the general framework

The state variable, structural shock and the control variables are as follows:

$$\mathbf{S}_t = \{\hat{R}_t^\star\}, \quad \text{and} \quad \mathbf{v}_t = \{\epsilon_t^A\}, \quad \text{and} \quad \boldsymbol{\psi}_t = \{\hat{X}_t, \hat{\Pi}\} \quad (22)$$

In this model, there are no fixed parameters. So $\bar{\Theta}$ is an empty set.

The estimated parameters are as follows:

$$\tilde{\Theta} = \{\beta, \sigma, \eta, \phi, \theta_\Pi, \theta_Y, \rho_A, \sigma_A\} \quad (23)$$

The NN is then trained to determine the output gap and inflation:

$$\begin{pmatrix} \hat{X}_t \\ \hat{\Pi}_t \end{pmatrix} = \psi_{NN}(\mathbf{S}_t, \tilde{\Theta} \mid \bar{\Theta}). \quad (24)$$

Mapping of Model #1 to the general framework

- Use 500,000 iterations to train the NN by minimizing the residual error for the Euler equation and Phillips curve.
- NN with five hidden layers with 256 neurons each.
- Continuously differentiable exponential linear units (CELU)
- AdamW optimizer
- Learning rate: Cosine annealing LR

For parallel computing, $B = 1000$.

Mapping of Model #1 to the general framework

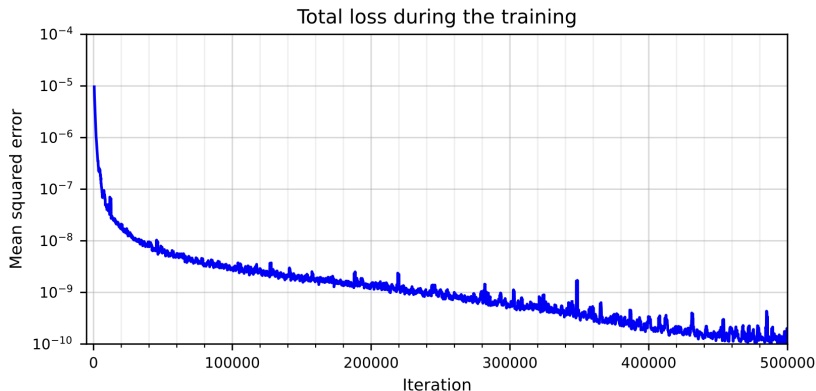


Figure 2: Convergence of the neural network solution. The figure shows the dynamics of the mean squared error. The loss is calculated as a moving average over 1000 iterations. The vertical axis has a logarithmic scale.

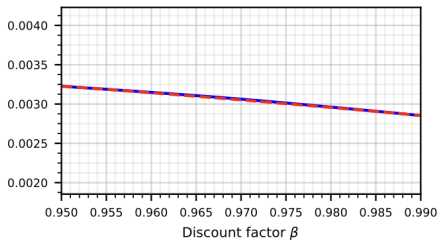
Mapping of Model #1 to the general framework

With the analytical solutions of the linearized DSGE model already in place, authors have considered the following values for the parameters:

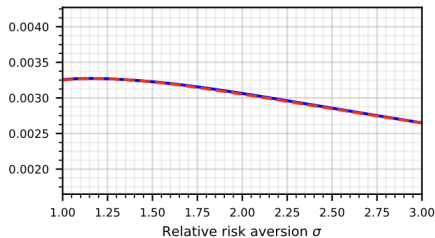
Parameters		LB	UB	Parameters		LB	UB
β	Discount factor	0.95	0.99	θ_{Π}	Mon.pol. inflation response	1.25	2.5
σ	Relative risk aversion	1	3	θ_Y	Mon.pol. output response	0.0	0.5
η	Inverse Frisch elasticity	1	4	ρ_A	Persistence TFP shock	0.8	0.95
φ	Price duration	0.5	0.9	σ_A	Std. dev. TFP shock	0.02	0.1

Table 1: Model parameters values – first proof of concept. The table shows the parameters of the three-equation New Keynesian model. The lower bound (LB) and upper bound (UB) show how we truncate the parameter space when training the NN to approximate the model solution.

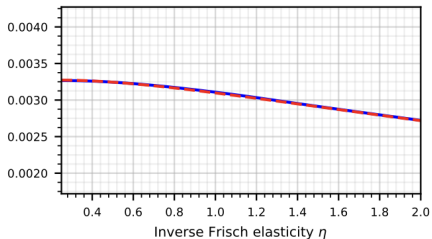
PF $\hat{\Pi}_t$ conditioned on β



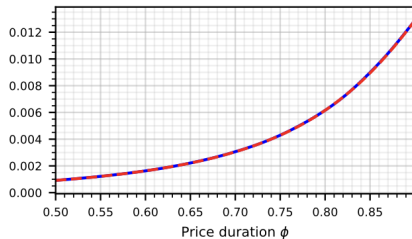
PF $\hat{\Pi}_t$ conditioned on σ



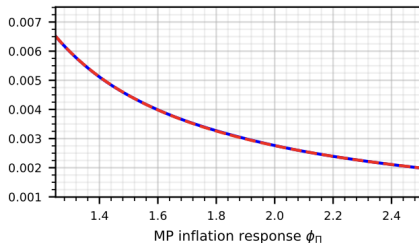
PF $\hat{\Pi}_t$ conditioned on η



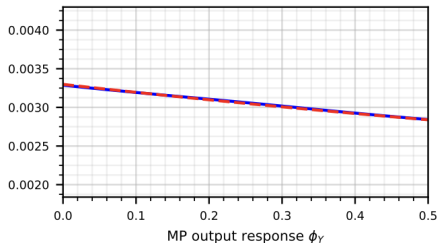
PF $\hat{\Pi}_t$ conditioned on ϕ



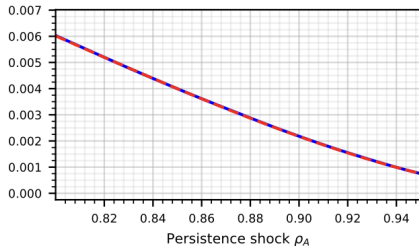
PF $\hat{\Pi}_t$ conditioned on ϕ_π



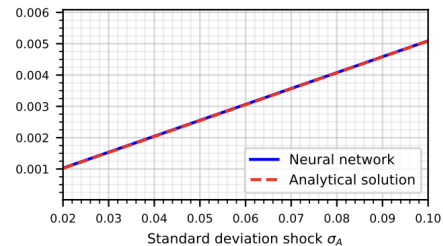
PF $\hat{\Pi}_t$ conditioned on ϕ_Y



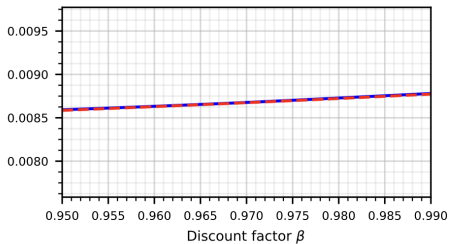
PF $\hat{\Pi}_t$ conditioned on ρ_A



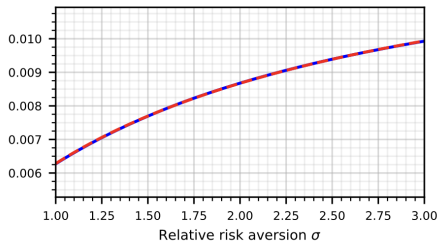
PF $\hat{\Pi}_t$ conditioned on σ_A



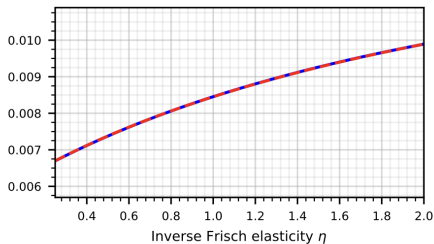
PF \hat{X}_t conditioned on β



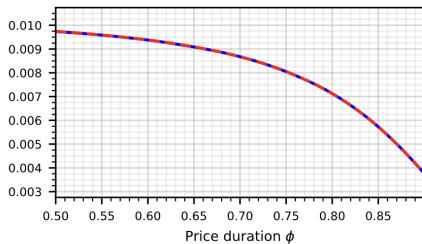
PF \hat{X}_t conditioned on σ



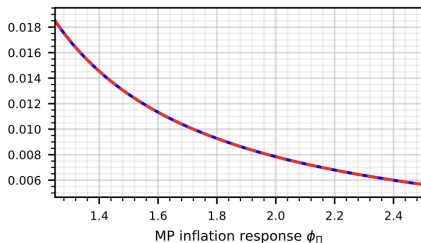
PF \hat{X}_t conditioned on η



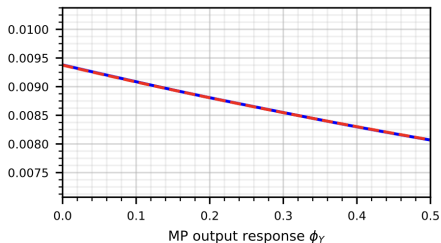
PF \hat{X}_t conditioned on ϕ



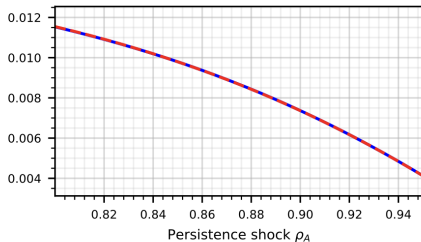
PF \hat{X}_t conditioned on ϕ_π



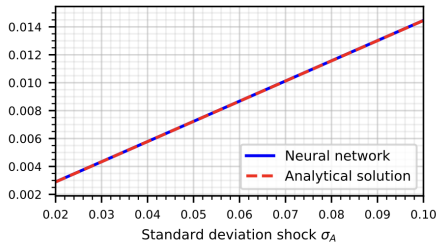
PF \hat{X}_t conditioned on ϕ_Y



PF \hat{X}_t conditioned on ρ_A



PF \hat{X}_t conditioned on σ_A



Training NN particle filter for Model #1

- There are 10,000 quasi-random parameter draws
- Might be potential issues with overfitting
- 80-20 split (training - validation splits)
- From the 8000 parameter draws (training), choose a batch size of 100 particles. Hence, optimization step is performed 80 times in each epoch.
- There are 20,000 epochs. Therefore, there are 1.6 million steps.

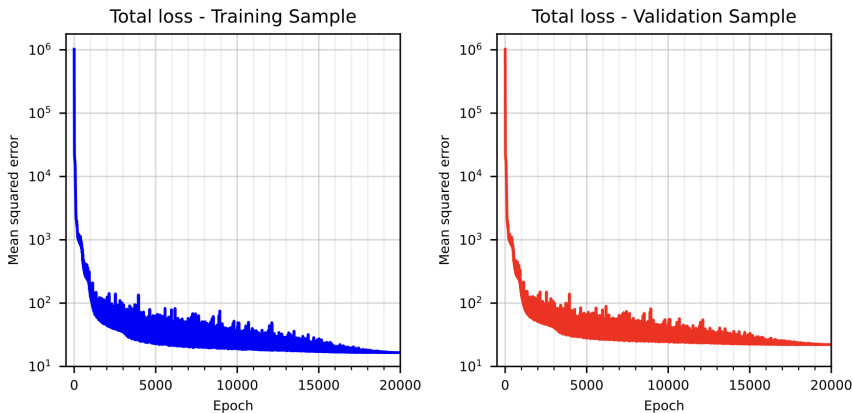


Figure 4: Training and validation convergence. The figure shows the total loss over the training sample (left) and over the validation sample (right). An epoch is completed when all the training or validation sample points are utilized. The vertical axis is expressed in a logarithmic scale.

Dealing with further overfitting issues in Model #1

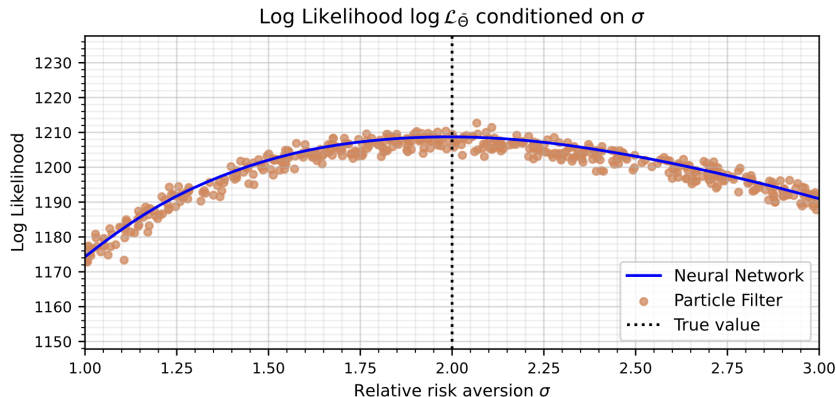
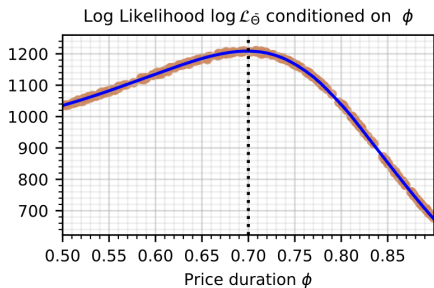
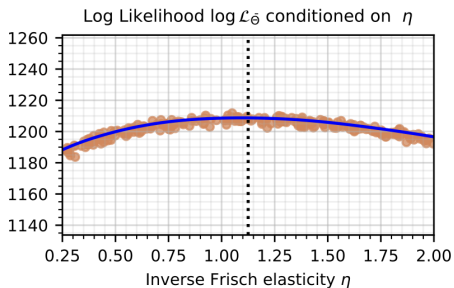
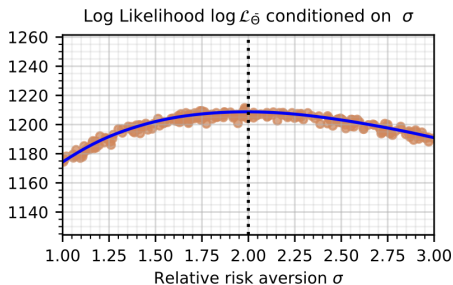
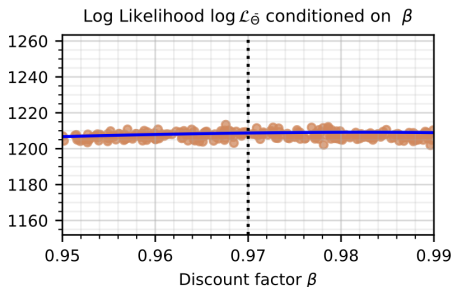
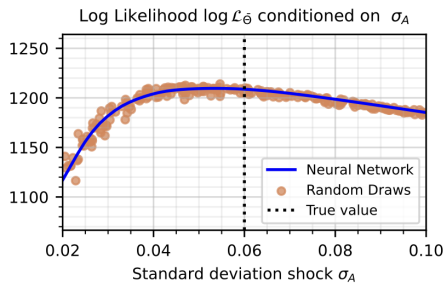
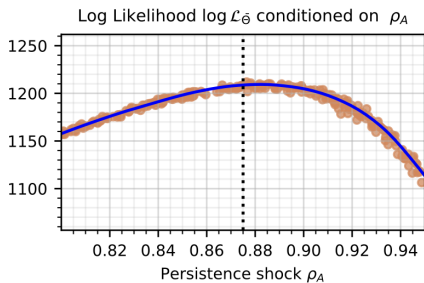
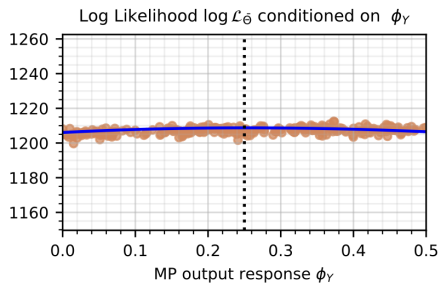
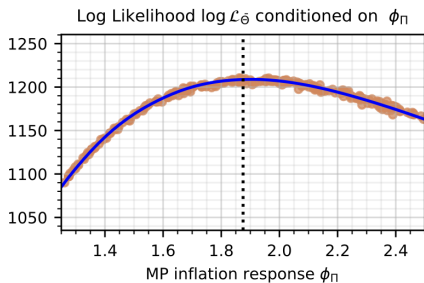


Figure 5: Accuracy in likelihood evaluations: NN particle filter vs. standard particle filter. The logarithm of the likelihood of the model as a function of the risk aversion parameter σ . The value of the fixed parameters is set to the middle of their bounds.





What's next?

- Solving a nonlinear RANK model using NN
- Solving a nonlinear HANK model using NN