

# **Estimation with GANs**

Master's Thesis

Presented to the  
Department of Economics at the  
Rheinische Friedrich-Wilhelms-Universität Bonn

In Partial Fulfillment of the Requirements for the Degree of  
Master of Science (M.Sc.)

Supervisor: Prof. Dr. Joachim Freyberger

Submitted in September 2024 by

**Marvin Benedikt Riemer**

Matriculation Number: 2799234

# Contents

<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>iv</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Background</b>	<b>1</b>
2.1 Structural estimation	1
2.2 Neural networks	2
<b>3 Adversarial estimation</b>	<b>4</b>
3.1 Examples of classifier discriminators	5
3.2 Generator objectives	7
3.2.1 Jensen-Shannon divergence	7
3.2.2 Wasserstein-p distance	8
3.3 Theoretical properties	9
3.3.1 Theorems in KMP	9
3.3.2 Applicability to simulations	12
<b>4 Simulation study</b>	<b>12</b>
4.1 The Roy model	12
4.2 General simulation structure	13
4.3 Implementation details	15
4.3.1 Initial values	15
4.3.2 Discriminators	15
4.3.3 Generator	16
4.4 Cross-sections of the loss landscape	17
4.5 Estimation	20
4.5.1 Estimation with the Wasserstein-1 divergence	20
4.5.2 Uniform initialization in wide intervals	22
<b>5 Conclusion</b>	<b>24</b>
<b>Appendix A Wasserstein-2 estimation results with uniform initialization</b>	<b>27</b>
<b>Appendix B Acknowledgement of system use</b>	<b>28</b>
<b>References</b>	<b>29</b>

## List of Figures

1	Replication of Figure 6 in Kaji, Manresa, and Pouliot (2023)	17
2	Replication of Figure 8 in Kaji, Manresa, and Pouliot (2023)	17
3	Losses cross-sections plotted over wider intervals	18
4	Diagonal loss cross-sections. Rows: Jensen-Shannon divergence, Wasserstein-1. Columns: $\text{diag}_{1,2}$ , $\text{diag}_{1,-2}$ .	19
5	Loss cross-sections around $\theta_0$ : Jensen-Shannon divergence, Wasserstein-1, Wasserstein-2	20
6	Results of 200 bootstrap estimations with Wasserstein-1 loss	21
7	Results of 200 estimations with Wasserstein-1 loss and initial values that are uniform over broad intervals	23
8	Results of 200 estimations with Jensen-Shannon loss and initial values that are uniform over broad intervals	25
A.1	Results of 200 estimations with Wasserstein-2 loss and initial values that are uniform over broad intervals	27

## List of Tables

1	Parameter estimates	22
2	Parameter estimates and standard errors with uniform initialization in wide intervals	22

# 1 Introduction

A saying inspired by Box (1976) goes: “All models are wrong, but some are useful.” If a model is sufficiently useful, structural estimation promises to get some external validity out of it. However, there are models for which a likelihood function for the observations they predict is not available. This complicates the attempt to estimate them structurally.

Kaji, Manresa, and Pouliot (2023) (KMP) propose a remedy, which they call “adversarial estimation”. It is inspired by *Generative Adversarial Networks* (GANs), which were introduced to the field of machine learning by I. J. Goodfellow et al. (2014). So-called Wasserstein GANs, based on the family of Wasserstein distances, were introduced by Arjovsky, Chintala, and Bottou (2017) to mitigate some shortcomings of GANs. As I will demonstrate, replacing part of the classical GAN machinery borrowed from I. J. Goodfellow et al. (2014) with the Wasserstein distance can also improve the adversarial estimator.

My presentation and notation largely follow the previously mentioned papers. Another paper on using GANs in Econometrics is Athey et al. (2021).

This thesis is structured as follows. Section 2 provides some background on structural estimation and neural networks. Section 3 presents the adversarial estimation approach, discusses variants of it, lays out the theory described in Kaji, Manresa, and Pouliot (2023), and checks whether it can be applied to there in my simulation studies. Section 4 then describes and reports results from my simulation study, which replicates parts of Kaji, Manresa, and Pouliot (2023)’s study, and shows a case where the Wasserstein distance can improve the adversarial estimator. Section 5 concludes.

## 2 Background

### 2.1 Structural estimation

Consider the problem of estimating the parameters of a structural economic model. Let  $Z \sim P(Z)$  be a vector of random noise variables,  $\Theta$  a parameter space and  $\mathcal{X}$  a space of outcomes or observations. Then an economic model can be imagined as a function  $f$  that maps a draw from the noise  $Z$  and a parameter  $\theta \in \Theta$  to an outcome  $X \in \mathcal{X}$ :

$$X = f(Z, \theta). \tag{1}$$

$f$  might be quite complicated, in particular a system of equations. What makes this equation *structural* is that  $f$  is an implementation of a particular economic model together with the assumption that this model captures the true relationships between economic variables, rather

than just their statistical relationships. With  $f$  therefore given by the presumed model, and  $Z$  noise, the goal is to find a good estimate  $\hat{\theta}$  of the true parameter  $\theta$ , on the basis of samples  $X_1, \dots, X_n$ .

$Z \sim P(Z)$  induces a distribution  $p(X_1, \dots, X_n|\theta)$ . This suggests the standard approach of maximum likelihood estimation, that is, to find  $\hat{\theta}$  such that

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \log p(X_1, \dots, X_n|\theta). \quad (2)$$

The logarithm here is not strictly necessary; it is, however, usually taken to make the likelihood function more computationally tractable, since it does not affect the position of the maximum. However, for some more sophisticated economic models, it is not easy or even possible to give an analytical expression for the likelihood function.

This motivates further approaches, in particular simulation methods, which attempt to infer  $\theta$  based on a simulation of the true data. *Simulation* here refers to drawing random samples from the proposed data-generating process given  $\hat{\theta}$ , an estimate of  $\theta$ . Most notable among these is perhaps the simulated method of moments, which performs inference based on the moments of the simulated data.

The question naturally arises of how to judge whether the simulated distribution gets sufficiently close to the real distribution. This motivates the idea of introducing a second component that provides “feedback” on this question in the form of a criterion function to be minimized. Updating the estimate  $\hat{\theta}$  to minimize this criterion is the central idea behind adversarial estimation, which I will discuss in more detail in section 3.

One approach for defining the criterion is to base it on a classification, more precisely, on the probability, as estimated by a classifier, that a given data point was drawn from the real data rather than from a simulated distribution. In machine learning, a popular tool for classification are neural networks, which I introduce next.

## 2.2 Neural networks

A neural network, in the most general sense, is a directed graph that defines how a certain output should be calculated from a given input and is therefore also called *computation graph*. There are many different structures (also called *architectures*) of that graph. This thesis only considers so-called *feed-forward neural networks*, for which the graph consists of:

- an input layer with one node for each input variable.

- an ordered set of hidden layers. The number of nodes in each is chosen by the researcher. The nodes get incoming edges only from the previous and have outgoing edges only to the successive layer.
- an output layer with one node for each output variable.

Embedded in any architecture is the fundamental notion of neural networks, where nodes are imagined as neurons that calculate inputs and outputs according to their edges. A fully connected feed-forward neural network is called *Multi-layer Perceptron*.

**Definition 1** (Multilayer Perceptron). Let  $l = 1, \dots, L$  be an ordered set of layers, with 1 the input and  $L$  the output layer. Let  $I_l$  be the number of nodes in each layer with  $b_{l,1}, \dots, b_{l,I_l}$  and  $\sigma_{l,1}, \dots, \sigma_{l,I_l}$  their biases and differentiable activation functions, respectively. For the nodes in the input layer  $v_{1,i}$ , set their output equal the inputs. For the nodes  $v_{l,i}$ ,  $l \geq 2$ , set:

$$v_{l,i} = \sigma(b_{l,i} + \sum_{j=1}^{I_{l-1}} w_{l,j} \cdot v_{(l-1),j}) \quad (3)$$

and interpret the output of the last nodes  $v_{L,i}$  as the output of the network.

Then a neural network  $\mathcal{N} = (V, E, \psi)$  with nodes  $V$ , edges  $E$  and a parameter set  $\psi$ , which holds the weights, biases and activation functions is called a multilayer perceptron.

Note that if the last layer consists only of one node taking values in  $[0, 1]$ , the network can be interpreted as a classifier suitable to the real-or-fake-data problem introduced above.

In practice, neural networks mostly are used for machine learning applications and their parameters *trained* using a gradient descent or other optimization method. In particular, they are sometimes not trained until the optimization method converges, so it is reasonable to include a specification of the training process in the definition.

**Definition 2** (Trained feed-forward neural network). Let  $\theta_{\text{train}}$  be a set of hyperparameters specifying the training of a neural network  $\mathcal{N}$ , including at least:

- A set of initial parameters  $\psi_0$
- A training algorithm  $A$
- A stopping criterion

Then call  $\mathcal{N}(\theta_{\text{train}}) = (V, E, \psi(\theta_{\text{train}}))$  a feed-forward neural network trained according to  $\theta_{\text{train}}$ .

In practice,  $\psi_0$  is usually chosen randomly. Note that the training algorithm  $A$  might itself require the setting of further hyperparameters. So-called *stochastic gradient descent* algorithms are the most common training methods, of which Adam (Kingma and Ba (2017)) is the most popular choice. The stopping criterion might be either the convergence criterion of the training algorithm, or that a certain number of training steps having been performed.

### 3 Adversarial estimation

The basic idea of adversarial estimation is to employ two auxiliary models, called the *generator* and the *discriminator* (or *critic*). The generator  $G : \Theta \times \mathcal{Z} \rightarrow \mathcal{X}$  creates simulated data based on a guess  $\hat{\theta} \in \Theta$  of the true parameter value  $\theta_0$ . The discriminator  $D : \mathcal{X} \rightarrow \mathbb{D} \subset \mathbb{R}$  returns for each real and generated data point some value. This value is then used to construct the value of an objective function for the optimization of the generator, which I will call *criterion* or *loss*. This loss function can be any divergence or distance between the distributions of the real and simulated data, including functions that are directly analytically tractable and do not strictly require a separate discriminator to calculate. Assuming however the involvement of a discriminator, the estimation can be viewed as the result of the following minimax game:

$$\hat{\theta}_{adv} = \arg \min_{\theta \in \Theta} \max_{D \in \mathcal{D}} \text{loss}(D(X_i, G(\theta, Z))). \quad (4)$$

Note that this game has a unique Nash-Equilibrium, where  $\hat{\theta} = \theta_0$  and  $D(x) = 0.5$  for all  $x$ .

This method is a variant of *Generative Adversarial Networks* (GAN), first proposed by [I. J. Goodfellow et al. \(2014\)](#) (later published as [I. Goodfellow et al. \(2020\)](#)). There, two neural networks take the role of generator and discriminator. In particular, the discriminator is a classifier network that outputs the probability of a data point being a real rather than simulated observation. While GANs have achieved great success in image generation and related tasks, they are not directly suitable for structural estimation. One reason is that the functional form of the generator network is usually very complex, with feed-forward neural networks often being fully connected and activation functions introducing non-linearities. Relatedly, the exact architecture of a neural network is usually not chosen to be economically (or at all) interpretable, but rather as an imprecise “art” based on predictive performance. Therefore, one essential contribution of [Kaji, Manresa, and Pouliot \(2023\)](#) is to impose that the generator has the structure of an economic model. This model being fully specified by  $\theta$  is what makes adversarial estimation meaningful and the result interpretable. They wouldn’t be if  $\theta$  were a long list of the weights and biases in a multi-layered neural network.

An implementation of (4) looks, generally, like algorithm 1.

Note that while [Kaji, Manresa, and Pouliot \(2023\)](#) stress the importance of sampling noise only one time, [I. J. Goodfellow et al. \(2014\)](#) and [Athey et al. \(2021\)](#) draw it anew for every training step of the discriminator and generator.



---

**Algorithm 1:** Adversarial estimation

---

```
1: Set necessary hyperparameters and initial values
2: Sample real observations  $X \sim P_0$ 
3: Sample noise  $Z \sim P_Z$ 
4: while Stopping criterion does not hold do
5:   Generate fake observations from the current generator  $\tilde{X} \sim P_{\hat{\theta}}$ 
6:   if The discriminator requires training: then
7:     Train the discriminator given the fake observations
8:   end if
9:   Calculate the criterion:  $\text{loss}(D(X, G(\hat{\theta}, Z)))$ 
10:  Update  $\hat{\theta}$ 
11: end while
```

---

There are various ways to fill in the details of this algorithm. The stopping criterion might be a convergence criterion of the generator’s optimization problem, or simply a sufficiently high number of repetitions being reached. A classification discriminator might take various forms, which I discuss below. There are two canonical choices for the loss function, which I discuss afterwards. The updates of the generator can be done with a gradient descent algorithm if it is differentiable or at least smooth enough that calculating numerical gradients will not lead an optimizer astray. Otherwise, they should be performed with a gradient-free optimization procedure, for example, a simplex algorithm.

Algorithm 1 in Kaji, Manresa, and Pouliot (2023) illustrates one way to fill out the details of algorithm 1. They use convergence as a stopping criterion, a (not necessarily trained to completion) neural network discriminator, cross-entropy loss, and update the generator using a version of the popular Adam (Kingma and Ba (2017)).

Their simulation code demonstrates another way, which I discuss in detail in section 4.2. There, they compare a range of estimators (including neural networks trained to completion) and update the generator using a gradient-free approach.

Next I discuss some of these terms in detail.

### 3.1 Examples of classifier discriminators

All the following discriminators have in common that they are classifiers, so for a data point  $x$  they return a probability that it is from the real rather than the simulated data. How this probability is then turned into an objective for the generator will be discussed in the next subsection.

Recall the game-theoretic view of adversarial estimation. If the true densities  $p_0$  and  $p_{\theta}(x)$  are known, the discriminator has a best response depending only on these, rather than on the

samples (of real and simulated data). Kaji, Manresa, and Pouliot (2023) call the discriminator playing this best response the *oracle discriminator*.

**Definition 3** (Oracle discriminator). The **oracle discriminator** assigns

$$D_\theta(x) := \frac{p_0(x)}{p_0(x) + p_\theta(x)} \quad (5)$$

to every  $x \in \mathcal{X}$ .

Of course,  $p_0$  and  $p_\theta(x)$  are unknown in practice. Nevertheless, this discriminator is useful as a benchmark in simulations and has an interesting theoretical property: If the simulated sample size  $m \rightarrow \infty$ ,  $\theta_{\text{oracle}}$  approaches  $\theta_{MLE}$ .

A simple statistical method for classification is logistic regression. In line with the simulation study in Kaji, Manresa, and Pouliot (2023), I consider a version that regresses on some collection of features of the data points and moments of the data.

**Definition 4** (Logistic discriminator). Let  $\Lambda$  be a sigmoid function with values in  $(0, 1)$ , and  $x^{\text{mom}}$  an  $(i + j) \times k$ -matrix of features and moments of the data calculated for each data point. Let  $(\beta_0, \dots, \beta_k \in \mathbb{R}^{k+1})$  be coefficients of a logistic regression run with  $x^{\text{mom}}$  as a regressor and an output vector  $Y$  consisting of 0s and 1s for the simulated and true observations. Then the **logistic discriminator** assigns

$$D(x) = \Lambda\left(\beta_0 + \sum_{k=1}^K \beta_k x_k^{\text{mom}}\right) \quad (6)$$

to every  $x \in \mathcal{X}$ .

Note that this classifier has to be calculated anew after each update of  $\theta$ . While this calculation will usually be fast on modern computers, the same is not necessarily true of the potentially more powerful neural network discriminator. Therefore, neural networks are often not trained to completion in practice and different training procedures might result in different discriminators.

**Definition 5** (Neural network discriminator). Let  $\mathcal{N}(\theta_{\text{train}}) : \mathcal{X} \rightarrow [0, 1]$  be a classifier neural network trained according to  $\theta_{\text{train}}$ . Then the **neural network discriminator** assigns

$$D(x) = \mathcal{N}(\theta_{\text{train}})(x) \quad (7)$$

to every  $x \in \mathcal{X}$ .

To understand more deeply the loss landscape which a neural network discriminator builds for the generator, we must consider the loss function on which it is trained.

## 3.2 Generator objectives

### 3.2.1 Jensen-Shannon divergence

The classical way to turn the probabilities of a classifier discriminator  $D(x)$  into an objective for the generator is the following:

**Definition 6** (Cross-entropy loss). The empirical cross-entropy loss (CE) is:

$$\frac{1}{n} \sum_{i=1}^n \log D(X_i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(X_{i,\theta})).$$

For  $X_i \sim p_0$  and  $X_{i,\theta} \sim p_\theta$  The empirical cross-entropy is an approximation of the following:

**Definition 7** (Population cross-entropy loss). The empirical cross-entropy loss is:

$$\mathbb{E}_{X \sim p_0} \log D(X) + \mathbb{E}_{X_\theta \sim p_\theta} \log(1 - D(X_\theta)).$$

A discriminator that maximizes the cross-entropy loss thereby calculates the Jensen-Shannon divergence (plus a constant) for the generator to minimize.

**Theorem 8.** *If the discriminator is the oracle discriminator  $D = D_\theta$ , the population cross-entropy loss evaluates as*

$$\mathbb{E}_{X \sim p_0} \log D(X) + \mathbb{E}_{X_\theta \sim p_\theta} \log(1 - D(X_\theta)) = 2 \log 2 + 2 \text{JSD}(p_0 \mid p_\theta) \quad (8)$$

where JSD is the Jensen-Shannon divergence.

The proof for this is contained in the proof for theorem 1 in [I. J. Goodfellow et al. \(2014\)](#).

A neural network discriminator that is not trained to completion returns an approximation of the Jensen-Shannon divergence. In practice, only such an approximation is often used, for two reasons: First, training a neural network to completion multiple times for every gradient calculation of the generator’s optimizer can be very computationally costly, especially given that GANs in practice are often large neural networks and are applied to high-dimensional data sets. Second, an imprecise estimate of the gradient often still leads to convergence.

However, even the Jensen-Shannon divergence calculated by an optimal discriminator has a crucial disadvantage: The divergence is maximal if  $p_G$  and  $p_0$  have disjoint support. Therefore, there are regions of the loss landscape where even the optimal CE-discriminator provides a gradient of zero in every direction at every point to the generator. If the generator “ends up” in such a region or the initial guess is there, algorithm 1 is unlikely to converge.

### 3.2.2 Wasserstein-p distance

Using the so-called Wasserstein-1 distance as criterion for the generator was first proposed by [Arjovsky, Chintala, and Bottou \(2017\)](#).

**Definition 9** (Wasserstein-p distance). For two probability distributions  $\mathbb{P}_0$  and  $\mathbb{P}_G$ , let  $\Pi(\mathbb{P}_0, \mathbb{P}_G)$  be the set of all joint distributions  $\gamma(x, y)$  whose marginals are  $\mathbb{P}_0$  and  $\mathbb{P}_G$ . Then for  $p \geq 1$ ,

$$W_p(\mathbb{P}_0, \mathbb{P}_G) = \inf_{\gamma \in \Pi(\mathbb{P}_0, \mathbb{P}_G)} \left( \mathbb{E}_{(x,y) \sim \gamma} d(x, y)^p \right)^{1/p}$$

is the Wasserstein-p distance (also Wasserstein p-distance) between  $\mathbb{P}_0$  and  $\mathbb{P}_G$ .

A natural interpretation of this equation comes from the field of optimal transport. It quantifies how much probability mass has to be moved how far in order to transfer  $\mathbb{P}_0$  into  $\mathbb{P}_G$ , or vice versa, assuming that this transport is done optimally. Inspired by this image, the Wasserstein-1 distance is also called *Earth-Mover distance*.

The Wasserstein distances deliver a measure of the distance between two distributions that is strictly monotone even if they are non-overlapping. However, since they require a solution to the optimal transport problem, they can be demanding to calculate, especially in high-dimensional spaces. In the context of GANs, it is natural to consider approximating it using a neural network. To this end, the following fact is helpful:

**Theorem 10 (Kantorovich-Rubinstein duality).**

$$W_1(\mathbb{P}_0, \mathbb{P}_G) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_0}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_G}[f(x)] \quad (9)$$

This dual representation of the Wasserstein-1 distance paves the way to approximating it using a neural network. The network approximates the function  $f$  and is often called critic instead of discriminator since it does not return a probability anymore. Unfortunately, it is not trivial to regularize a neural network to obey the Lipschitz constraint. [Arjovsky, Chintala, and Bottou \(2017\)](#) clamp the weights of the neural network to lie in a compact space. They themselves describe this approach as “clearly terrible”, since there is no principled way to chose the clipping parameter and setting it too big or too small comes with difficult trade-offs.

[Gulrajani et al. \(2017\)](#) propose to penalize the norm of the gradient of the critic with respect to its input. This solution has been more widely accepted and is also used by [Athey et al. \(2021\)](#).

Of course, it is also possible to approximate the Wasserstein-1 distance without training a neural network. Since the Wasserstein distance is the solution to an optimal transport problem, it

can be derived using a “Pseudo-auction algorithm”. For differentiability, this algorithm can then be approximated using a smoothed “soft-auction” algorithm. The divergence resulting from this smoothed algorithm is called the *Sinkhorn divergence*.

### 3.3 Theoretical properties

#### 3.3.1 Theorems in KMP

Kaji, Manresa, and Pouliot (2023) contains three main theoretical results. In line with the focus of their paper, they are stated for the case with the (estimated) Jensen-Shannon distance as a criterion. They draw on the literature on M-estimation, particularly as presented in Chapter 3.2 of Vaart and Wellner (2023). I begin by laying out some of their notation.

Let  $\mathbb{M}_\theta(D) := \mathbb{P}_0 \log D + \mathbb{P}_\theta \log(1 - D)$  be the sample criterion function,  $\mathbb{M}_\theta(\hat{D}_\theta)$  the estimated, and  $\mathbb{M}_\theta(D_\theta)$  the oracle loss.

Let  $T_\theta(Z_i)$  be a map that represents the generator, and endow the parameter space  $\Theta$  with a distance  $h(\theta_1, \theta_2) := \sqrt{\int (\sqrt{p_{\theta_1}} - \sqrt{p_{\theta_2}})^2}$ .

They take the following symbols from Vaart and Wellner (2023):

The symbol “ $\lesssim$ ” means “is bounded from above up to a universal constant”. “ $\rightsquigarrow$ ” refers to weak convergence.

The first theorem states that for some reasonable conditions on the true and estimated criterion functions, which I discuss below, the adversarial estimator is indeed consistent.

**Theorem 11 (Theorem 1, KMP).** *Suppose*

- (1) *For every open  $G \in \Theta : \inf_{\theta \notin G} \mathbb{M}_\theta(D_\theta) > \mathbb{M}_{\theta_0}(D_{\theta_0})$ ,*
- (2)  *$\{\log D_\theta : \theta \in \Theta\}$  and  $\{\log D_\theta \circ T_\theta : \theta \in \Theta\}$  are  $P_0$ - and  $P_Z$ -Glivenko-Cantelli, respectively,*
- (3)  *$\sup_{\theta \in \Theta} \|\mathbb{M}_\theta(\hat{D}_\theta) - \mathbb{M}_\theta(D_\theta)\| \rightarrow 0$  in probability, and*
- (4)  *$\hat{\theta}$  satisfies  $\mathbb{M}_{\hat{\theta}}(\hat{D}_{\hat{\theta}}) \leq \inf_{\theta \in \Theta} \mathbb{M}_\theta(\hat{D}_{\hat{\theta}}) + o_P^*(1) \rightarrow 0$ .*

*Then  $h(\hat{\theta}, \theta_0) \rightarrow 0$ .*

Condition 1 is that the optimum  $\theta_0$  is identified by the true criterion. Condition 2 is that the generator and discriminator terms of the cross-entropy criterion converge to the true distributions  $P_0$  and  $P_Z$ , respectively. Condition 3 is that the estimated criterion converges uniformly to the true criterion in the entire parameter space. Condition 4 is that the estimated criterions of estimates  $\hat{\theta}$  converge to the estimands  $\theta$  at rate 1 in outer probability.

The second theorem gives a rate of convergence for the estimator and requires multiple assumptions. The first requires, roughly speaking, that the generator is parametric, that is, it

depends smoothly on the parameters, has non-exploding, non-vanishing Fisher information, therefore can be stably inverted, and that the inverted generator is also parametric in this sense.

**Assumption 1 (A1, KMP).** (1)  $\Theta$  is (a subset of) a Euclidean space.

(2)  $p_\theta$  is differentiable in  $\theta$  at every  $\theta \in \Theta$  for every  $x \in \mathcal{X}$  with the derivative continuous in both  $x$  and  $\theta$ .

(3) The maximum eigenvalue of the Fisher information  $I_\theta = P_\theta \dot{\ell}_\theta \dot{\ell}_\theta^\top$  is bounded uniformly in  $\theta \in \Theta$ .

(4) The minimum eigenvalue of  $I_\theta$  is bounded away from 0 uniformly in  $\theta \in \Theta$ .

The same is assumed for the “inverted” structural model  $\tilde{\mathcal{P}}_\theta = \{((p_0/p_\theta) \circ T_\theta) p_Z : \theta \in \Theta\}$ .

The second condition is on the growing synthetic sample size, and, barring computational constraints, can easily be guaranteed by the researcher.

**Assumption 2 (A2, KMP).**  $n/m \rightarrow 0$

The third assumption has two parts: The first is that the estimated criterion converges to its minimum for the true  $\theta$  at rate  $o_P^*(n^{-1})$ . The second part concerns itself with the estimation error of the excess oracle loss of  $\hat{\theta}$  over  $\theta$ . This estimation error is the difference between the excess loss under the estimated discriminator and the excess loss under the oracle discriminator, where the former is an estimate for the latter. Therefore, this condition imposes that the infimum over possible values of the true parameter vector  $\theta$  of the estimation error should be bounded by  $o_P^*(n^{-1})$

**Assumption 3 (A3, KMP).** There exists a sequence of open balls  $G_n := \{\theta \in \Theta : h(\theta, \theta_0) < \eta_n\}$  such that

(1)  $\eta_n \sqrt{n} \rightarrow \infty, \mathbb{M}_{\hat{\theta}}(\hat{D}_{\hat{\theta}}) \leq \inf_{\theta \in G_n} \mathbb{M}_\theta(\hat{D}_\theta) + o_P^*(n^{-1})$ , and

(2)  $\inf_{\theta \in G_n} [\mathbb{M}_{\hat{\theta}}(\hat{D}_{\hat{\theta}}) - \mathbb{M}_\theta(\hat{D}_\theta)] - [\mathbb{M}_{\hat{\theta}}(D_{\hat{\theta}}) - \mathbb{M}_\theta(D_\theta)] = o_P^*(n^{-1})$ .

As noted by KMP, satisfying the second part will involve the derivative of the estimated loss in  $\theta$  converging to that of the oracle. This can be empirically checked by plotting cross-sections of the loss around the original value, as the authors and I do in the simulation part.

The fourth assumption imposes two requirements to aid identification: That the true loss has approximately quadratic curvature near the optimum and that the true distribution  $P_0$  and the distribution generated by the true parameter  $P_{\theta_0}$  overlap around the optimum.

**Assumption 4 (A4, KMP).** (1) There exists an open set  $G \subset \Theta \subset \mathbb{R}^k$  containing  $\theta_0$  in which  $\mathbb{M}_\theta(D_\theta) - \mathbb{M}_{\theta_0}(D_{\theta_0}) \gtrsim h(\theta, \theta_0)^2$ .

$$(2) \quad h(\theta, \theta_0)^2 = O\left(\int D_{\theta_0}(\sqrt{p_{\theta_0}} - \sqrt{p_{\theta}})^2\right) \text{ as } \theta \rightarrow \theta_0.$$

Together, these yield:

**Theorem 12 (Theorem 2, KMP).** *Under Assumptions 1 to 4,  $h(\hat{\theta}, \theta_0) = O_P^*(n^{-1/2})$ .*

Regarding the efficiency of the estimator, consider this fifth assumption. It is a stronger version of assumption 1, but requires a higher degree of smoothness, namely the twice differentiability of the likelihood of the generator in the parameters and observations. The third condition also implies the quadratic shape of the criterion, as demanded by assumption 4 above.

**Assumption 5 (A5, KMP).** (1) *The parameter space  $\Theta$  is (a subset of) a Euclidean space  $\mathbb{R}^k$ .*

(2) *The structural model  $\{P_{\theta} : \theta \in \Theta\}$  has a likelihood that is twice differentiable in  $\theta$  at  $\theta_0$  for every  $x \in \mathcal{X}$  with the derivatives continuous in both  $x$  and  $\theta$ .*

(3) *The Fisher information matrix  $I_{\theta_0} := P_{\theta_0} \dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^{\top} = -P_{\theta_0} \ddot{\ell}_{\theta_0}$  is positive definite.*

(4) *The matrix  $\tilde{I}_{\theta_0} := 2P_{\theta_0} \left( D_{\theta_0} \dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^{\top} + (\ddot{\ell}_{\theta_0} + \dot{\ell}_{\theta_0} \dot{\theta}_{\theta_0}^{\top}) \log(1 - D_{\theta_0}) \right)$  is positive definite.*

(5)  *$T_{\theta}$  is continuously differentiable in  $\theta$  for every  $x \in \mathcal{X}$  and  $P_0$  has a likelihood that is continuously differentiable in  $x$ .*

Now Kaji, Manresa, and Pouliot (2023) arrive at their third theorem. It states the multinormal distribution towards which the standardized adversarial estimator weakly converges.

**Theorem 13 (Theorem 3, KMP).** *Under the conclusion of Theorem 2 and Assumptions 2, 3, and 5,*

$$\sqrt{n}(\hat{\theta} - \theta_0) = 2\tilde{I}_{\theta_0}^{-1} \sqrt{n} [\mathbb{P}_0(1 - D_{\theta_0}) \dot{\ell}_{\theta_0} - \mathbb{P}_{\theta_0} D_{\theta_0} \dot{\ell}_{\theta_0} - \tilde{\mathbb{P}}_0 \tau_n] + o_P^*(1) \rightsquigarrow N(0, \tilde{I}_{\theta_0}^{-1} V \tilde{I}_{\theta_0}^{-1}) \quad (10)$$

where  $V := \lim_{n \rightarrow \infty} 4P_{\theta_0} D_{\theta_0} (1 - D_{\theta_0}) \dot{\ell}_{\theta_0} \dot{\ell}_{\theta_0}^{\top}$ .

Efficiency requires correct specification of the generator model. This means that for the true parameter  $\theta_0$ , it simulates the true distribution  $P_0$ , for which the discriminator has the optimal response to assign probability 1/2.

**Assumption 6 (A6, KMP).** *The synthetic model  $\{P_{\theta} : \theta \in \Theta\}$  is correctly specified, that is,  $P_{\theta_0} = P_0$  and  $D_{\theta_0} \equiv 1/2$ .*

This allows a simplification of theorem 13 to the following efficiency result:

**Theorem 14 (Corollary 4, KMP).** *Under the conclusion of Theorem 3 and Assumption 6,*  
 $\sqrt{n}(\hat{\theta} - \theta_0) = I_{\theta_0}^{-1} \sqrt{n} (\mathbb{P}_0 - \mathbb{P}_{\theta_0}) \dot{\ell}_{\theta_0} + o_P^*(1) \rightsquigarrow N(0, I_{\theta_0}^{-1}).$

### 3.3.2 Applicability to simulations

Kaji, Manresa, and Pouliot (2023) partially discuss the applicability of the theoretical results to their simulation study. It seems plausible that the requirements of theorem 11 are fulfilled by at least one their discriminators, and indeed their results show convergence and amicable loss landscapes. The discrete choice nature of the Roy model (which will be introduced below) seems problematic for the results on convergence and efficiency, which depend on the generator. First, it is not obvious that it can be inverted, as required by assumption 1. Second, because  $d_1$  and  $d_2$  are (crucial) parts of the observation, the generator is not differentiable and therefore not twice differentiable along  $x$ , as required by Assumptions 1 and 5.

Regarding my own simulations, of course the unchanged Roy model generator will preserve these difficulties. Focusing on the Sinkhorn approximation of the Wasserstein discriminator, it is beyond the scope of this thesis to check whether the results above can be transferred to this case. However, consider the simple example of learning only a location parameter. The population Wasserstein-1 distance will not show quadratic curvature near the optimum, violating assumptions 4 and 5. Therefore, the Wasserstein-2 loss seems more promising to uphold the theoretical results.

## 4 Simulation study

The authors simulate the estimation of the Roy model, a discrete choice model which has intractable likelihood for certain parameter values.

### 4.1 The Roy model

The Roy model models a set of agents choosing which sector to work in in each of two time periods. At the start of the game, nature determines the (natural logarithms of the) wages offered to each agent, by the following formulas:

$$\log w_{i1s} = \mu_s + \varepsilon_{i1s} \quad (11)$$

$$\log w_{i2s} = \mu_s + \gamma \mathbf{1}_{d_{i1}=s} + \varepsilon_{i2s} \quad (12)$$

where the noise of the offered wages is distributed as follows:

$$\begin{bmatrix} \varepsilon_{i11} \\ \varepsilon_{i12} \\ \varepsilon_{i21} \\ \varepsilon_{i22} \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho_s \sigma_1 \sigma_2 & \rho_t \sigma_1^2 & \rho_s \rho_t \sigma_1 \sigma_2 \\ \rho_s \sigma_1 \sigma_2 & \sigma_2^2 & \rho_s \rho_t \sigma_1 \sigma_2 & \rho_t \sigma_2^2 \\ \rho_t \sigma_1^2 & \rho_s \rho_t \sigma_1 \sigma_2 & \sigma_1^2 & \rho_s \sigma_1 \sigma_2 \\ \rho_s \rho_t \sigma_1 \sigma_2 & \rho_t \sigma_2^2 & \rho_s \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix} \right). \quad (13)$$



In the first time period, each agent  $i$  observes the wages  $\log w_{i11}$  and  $\log w_{i12}$  offered to them in the two sectors. Knowing their own discount factor  $\beta$  and the parameters  $\gamma_1$  and  $\gamma_2$ , they solve the dynamic programming problem and pick a sector for the first period that maximizes their expected discounted sum of wages. In the second period, the wages  $\log w_{i21}$  and  $\log w_{i12}$  are revealed to them and they pick a sector.

The researcher observes the realized wages  $\log w_1$  and  $\log w_2$  as well as the corresponding sector choices  $d_1, d_2 \in 1, 2$ . Broadly speaking, the parameters  $\mu_1, \mu_2$ , and to a lesser degree  $\gamma_1$  and  $\gamma_2$ , are location parameters of the distributions of offered wages, while  $\sigma_1, \sigma_2, \rho_s$  and  $\rho_t$  determine the shape and correlations (shapes of the joint distributions).

Note the selection effects resulting from only realized wages being observed: in particular, if  $\mu_i$  is sufficiently below  $\mu_j$ , sector  $i$  will never be picked and  $\mu_i$  will not be identified.

If  $\rho_t = 0$ , a likelihood function for observations from the Roy model is available.

## 4.2 General simulation structure

First, I reproduce parts of the authors' simulation in Python (Python Software Foundation (2023)). I mostly build on the scientific Python stack, more precisely, the packages `numpy`, `scipy`, and `scikit-learn` (Harris et al. (2020), Virtanen et al. (2020), and Pedregosa et al. (2011), respectively).

I implement the Sinkhorn divergence as an approximation of the Wasserstein distance using the `SampleLoss` function from the `GeomLoss` (Feydy et al. (2019)) package. It is also optimized for running on GPUs by virtue of being built on `KeOps` (Charlier et al. (2021)). `GeomLoss` is built on `PyTorch` (Ansel et al. (2024)), which I therefore integrate into the scientific Python stack for my simulation. I also use `PyTorch` for a GPU-accelerated calculation of the Jensen-Shannon divergence.

For my scientific Python code, I use the `mp` module from Python's standard library to parallelize simulation runs on an HPC cluster (cf. appendix B). I generate plots using `Matplotlib` (Hunter (2007)). My code is available upon request.

The replication package of Kaji, Manresa, and Pouliot (2023) can be downloaded from the journal website. It contains the authors' simulation code, written in Matlab. As the authors state in the readme file, the simulations for the Roy model are contained in the files `main_roy.m` (Figures 6, 7) and `main_case.m` (Figures 8, 9, and Table I). They draw on functions in other files to simulate data and calculate losses.

Both main files share a general structure: after setting parameters of the simulation itself (for example sample sizes, number of simulation runs) and the Roy model, the values of loss

functions are calculated along a linear grid and then rendered to create Figures 6 and 8. The grid describes seven- or eight-dimensional “cross-hairs”, where one parameter is varied while the others remain fixed at the true value. Thereafter, real and fake observations are generated and an initial guess is generated. Then, the Roy model is estimated using multiple methods, which are implemented as constrained minimizations. The constraints are bounds on the parameters of the Roy model, on which the authors do not further elaborate, but which are likely added for computational efficiency. Where necessary, an additional nonlinear constraint enforces that the guesses of the minimizer stay within the support of the Roy model.

---

**Algorithm 2:** Initial guess in `main_roy.m`

---

```

1: Input: True parameter vector  $\theta$ , lower and upper bounds  $L$ ,  $U$ , the support of the Roy model  $\mathcal{S}$ 
2: while  $\theta_{0,p} \notin \mathcal{S}$ , for each parameter  $p$  of  $\theta$  to be estimated: do
3:   Sample noise  $u_p \sim \mathcal{N}(0, 0.2)$ 
4:   Set  $\theta_{0,p} := \theta_p + u_p$ 
5:   Clip  $\theta_{0,p} := \min(\max(\theta_{0,p}, L_p), U_p)$ 
6: end while

```

---

In `main_roy.m`, the authors plot cross-sections of the loss landscapes generated by the oracle and neural network discriminator as well as the loss implied by MLE. Then, they generate the initial guess  $\theta_0$  using algorithm 2. Following simple maximum likelihood estimation with  $\theta_0$  as an initial value, they perform adversarial estimation with:

- the oracle discriminator, using the result of MLE as initial value
- the Logistic discriminator, using the result of the previous step as initial value
- the neural network discriminator, again using the result of the oracle step as initial value

They also estimate the Roy model using Indirect Inference and optimally-weighted SMM, but don’t use the results.

`Main_case.m` simulates the untractable likelihood case. Therefore, an initial guess  $\theta_0$  is drawn similarly to algorithm 2, but the first draw is used. Logistic regression is then performed using  $\theta_0$  as an initial value, and the result used to initialize adversarial estimation. Again, loss-curves for the neural network and logistic discriminator are plotted, including for  $\rho_t$ . To study the properties of the adversarial estimators, Kaji, Manresa, and Pouliot (2023) then perform the bootstrap, sampling with replacement from the noise and the true observations independently. They perform estimation with the logistic discriminator using the previous logistic estimate as the initial value. The result of this serves as the initial value for estimation with the neural network discriminator, as well as for the Indirect Inference and optimally-weighted SMM estimators (but the authors don’t report the Indirect Inference results in the paper).

## 4.3 Implementation details

### 4.3.1 Initial values

As mentioned above, Kaji, Manresa, and Pouliot (2023) use results of estimation procedures as initial guesses for other estimation procedures. This is unproblematic in the sense that an estimator will hopefully converge to the optimum independent of its starting value and therefore the final distribution of estimates should be largely independent of the initializations. However, about their Figure 7, they write: “The resulting estimators are comparable with MLE”, referring to the oracle and the neural network discriminators. This is unsurprising, perhaps even disappointing, given that the theoretically optimal oracle estimator is initialized with the result of MLE and the neural network estimator with the result of the oracle estimation.

This issue is less pronounced for their Figure 9, because the setting of initial values does not involve the impossible-in-practice oracle discriminator. However, it is still not surprising that the neural network estimators results are “comparable” to the logistic estimator with which it is initialized.

In the README.pdf of the replication package, Kaji, Manresa, and Pouliot (2023) propose to start adversarial the adversarial estimation by pre-estimation with a logistic discriminator. Their Figures 7, 9, and table I should perhaps rather be seen as indicating neural network adversarial estimation does not bring large improvements after this first step in the Roy model case.

To achieve a proper comparison, I reproduce their choice of initial values in my simulation of the Wasserstein discriminator. I will leave out the pre-estimation in my simulation of initial values drawn from wider intervals, because then I’ll be interested in the more pure properties of the different criterion functions.

### 4.3.2 Discriminators

I reproduce the MLE and oracle discriminator by translating the `logroypdf.m` function that Kaji, Manresa, and Pouliot (2023) have provided for the case with  $\rho_t$  known to be zero. Kaji, Manresa, and Pouliot (2023) use two logistic regression discriminators. The first in `main_case.m`, is:

$$D_{\text{loss1}}(x) : \Lambda(\beta^\top x^{\text{mom}}) \quad (14)$$

with  $\beta = (\beta_0, \dots, \beta_7)$  and  $x^{\text{mom}} = (1, \log w_1, d_1, \log w_2, d_2, (\log w_1)^2, (\log w_2)^2)$ . As illustrated in Figure 8 of the working paper version (compare also section 4.4),  $\rho_t$  is not identified using this discriminator. Therefore, they add the cross-moment between  $(\log w_1)$  and  $(\log w_2)$  in `main_case.m`:

$$D_{\text{loss2}}(x) : \Lambda(\beta^\top x^{\text{mom}}) \quad (15)$$

with  $\beta = (\beta_0, \dots, \beta_8)$  and  $\mathbf{x}^{mom} = (1, \log w_1, d_1, \log w_2, d_2, (\log w_1)^2, (\log w_2)^2), \log w_1 \cdot \log w_2$ . I employ `sklearn.linear_model.LogisticRegression` for both logistic regression discriminators.

The authors' code for the neural network discriminator is in `NND.m`. It uses Matlab's `patternnet` and `train`. The scientific Python stack comes with limited support for neural networks, but I can sufficiently approximate the authors' discriminator using `sklearn.neural_network.MLPClassifier`.

Following the authors, I create a net with 1 hidden layer containing 10 nodes, followed by the tanh activation function. Inspecting sklearn's source code reveals that a logistic output activation function is automatically set. The authors train their network with a conjugate gradient descent algorithm. Because this is not available to train `MLPClassifier`, I use the Adam algorithm (Kingma and Ba (2017)).

`MLPClassifier`'s default convergence criteria cause my code to raise warnings about non-convergence of the discriminator nets. This is not completely mitigated even by setting `max_iter` (the maximum number of iterations of the optimizer) to 2000 (10 times the default value), at the cost of a longer runtime. Nevertheless, the networks converge well enough under the default settings. Leaving `max_iter` at 200, but increasing `tol`, the tolerance of the convergence criterium, five- or tenfold mitigates the warnings but results in flatter and less smooth loss functions. Therefore, I leave the default settings and accept the warnings.

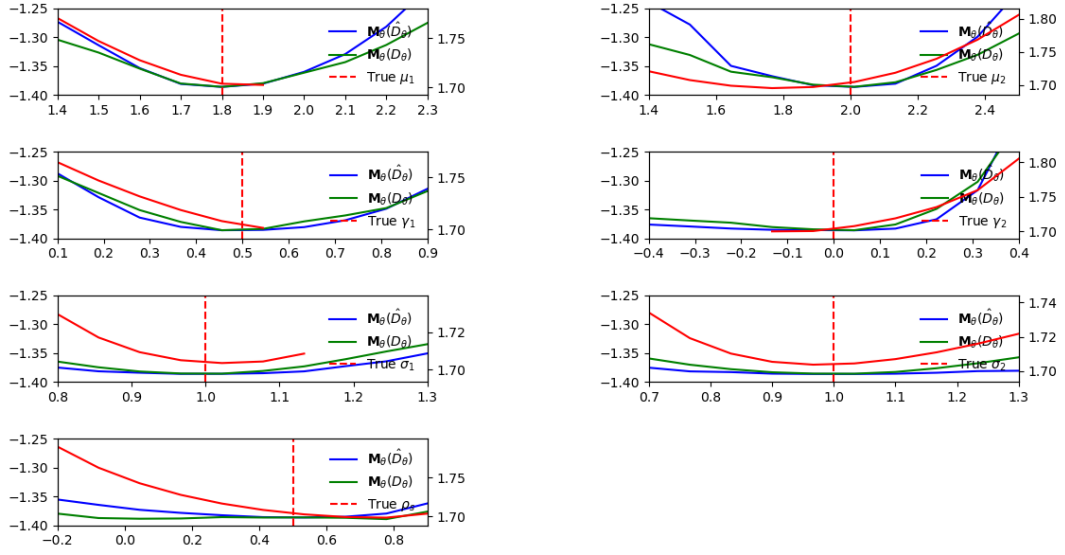
The authors also set the normalization and regularization parameters of `patternnet`. Since these are handled differently in `MLPClassifier`, I do not translate this adaption.

Figures 1 and 2 in section 4.4 suggest that these modifications do not perceptively alter the loss landscape.

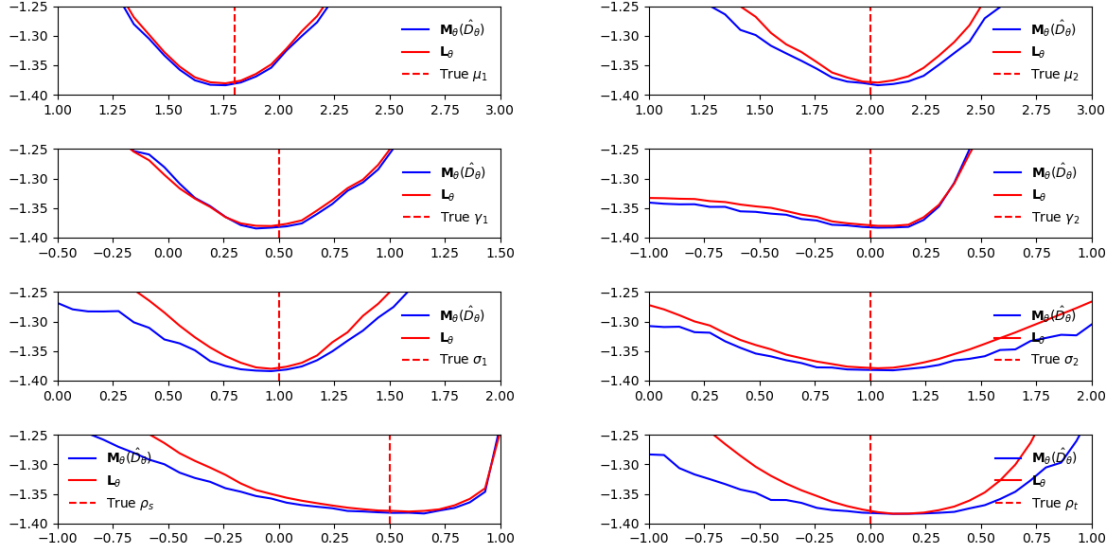
### 4.3.3 Generator

Both true and simulated observations from the roy model are generated by drawing uniform noise, transforming it into the multinormally distributed shocks  $(\varepsilon_{i11}, \varepsilon_{i12}, \varepsilon_{i21}, \varepsilon_{i22})$  and then, based on some given  $\theta$ , calculating the decisions of the agents. The authors provide an option to smooth the observations, but do not use it, since the loss crosssections (cf. section 4.4) look sufficiently smooth and the estimations work nevertheless.

For the outer optimization loop that trains the generator, the authors use the third-party `fminsearchcon` function (D'Errico (2024)). This is a wrapper function that adds support for bounds and nonlinear constraints to Matlab's built-in `fminsearch`, which employs the Nelder-Mead simplex algorithm (Lagarias et al. (1998)) to minimize a function without computing gradients. I employ `scipy.optimize.minimize`, which natively supports the Nelder-Mead



**Figure 1.** Replication of Figure 6 in Kaji, Manresa, and Pouliot (2023)

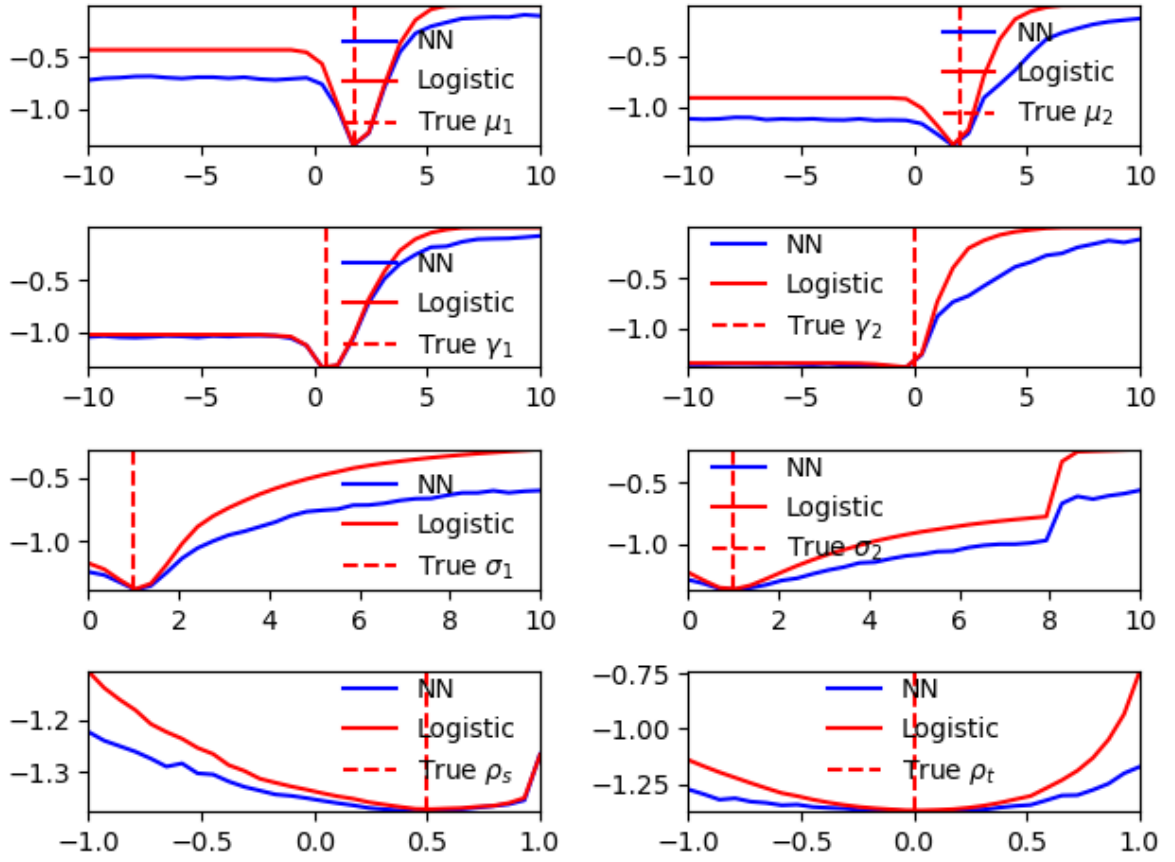


**Figure 2.** Replication of Figure 8 in Kaji, Manresa, and Pouliot (2023)

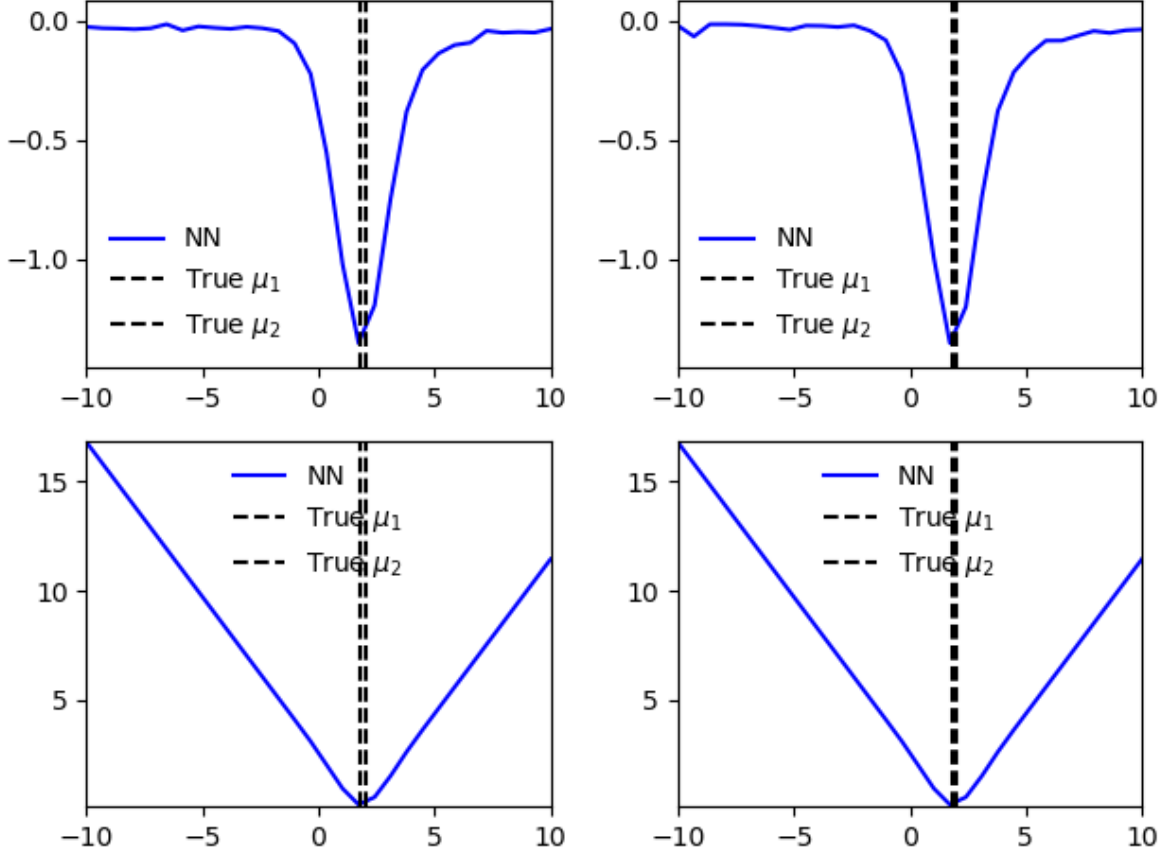
algorithm with bounds and nonlinear constraints. I set an option to perform a version of the Nelder-Mead algorithm that's adapted to the dimensionality of the problem (Gao and Han (2012)), which shows improved convergence in my simulation.

#### 4.4 Cross-sections of the loss landscape

Figures 1 and 2 match Figures 6 and 8 in Kaji, Manresa, and Pouliot (2023), confirming that my replication in the scientific Python stack, including of the neural network discriminator, is sufficiently close to the original.



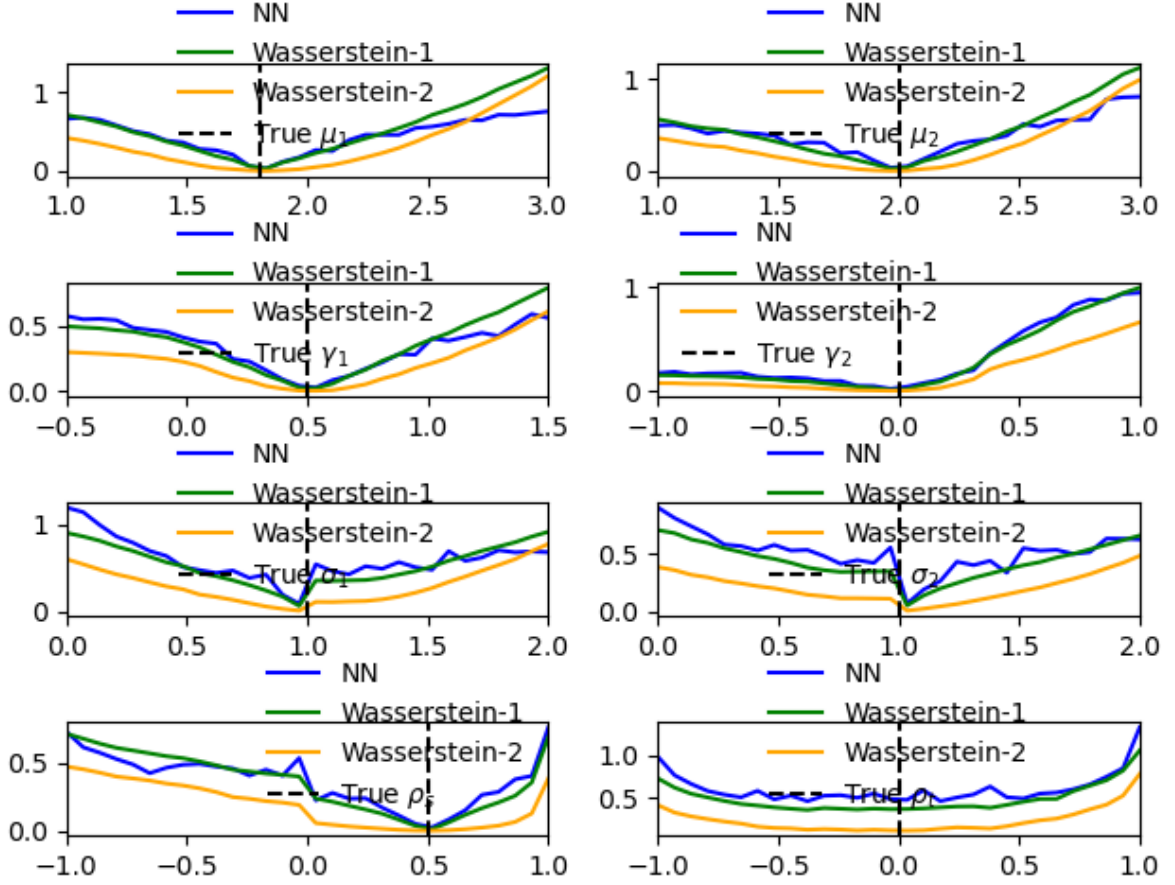
**Figure 3.** Losses cross-sections plotted over wider intervals



**Figure 4.** Diagonal loss cross-sections. Rows: Jensen-Shannon divergence, Wasserstein-1. Columns:  $\text{diag}_{1,2}$ ,  $\text{diag}_{1,-2}$ .

Figure 3 shows a variant of figure 2 plotted over wider intervals. It is striking that for some parameters, the loss is flat when moving too far away from the optimal value. Partly, this can be explained by the discrete choice nature of the Roy model. Consider for example  $\mu_1$ . If it becomes too small relative to  $\mu_2$  while  $\gamma_1$  and  $\gamma_2$  are held constant, the agents stop choosing sector 1. Since only their chosen sectors and wages are observed, in such cases there are no observations that help to narrow down the value of  $\mu_1$  (except to bound it from above). Similar arguments explain the flatness towards the tail of the cross-hairs for all four location parameters  $\mu_1$ ,  $\mu_2$ ,  $\gamma_1$ , and  $\gamma_2$ .

Recall from section 3.2.1 that there is another reason for the loss-function to become flat, at least for the neural network estimator with cross-entropy loss. Namely, the constant Jensen-Shannon divergence for disjoint distributions. To isolate the effect, I rotate part of the cross-hairs to look at the loss along the diagonal  $\text{diag}_{1,2} = \{(\mu_1, \mu_2) = (m, m); m \in \mathbb{R}\}$ . This way, the fake distribution can become disjoint from the real distribution without affecting the agents' choices. Note, however that there is a small distortion because  $\mu_1 = 1.8 \neq 2.0 = \mu_2$ .



**Figure 5.** Loss cross-sections around  $\theta_0$ : Jensen-Shannon divergence, Wasserstein-1, Wasserstein-2

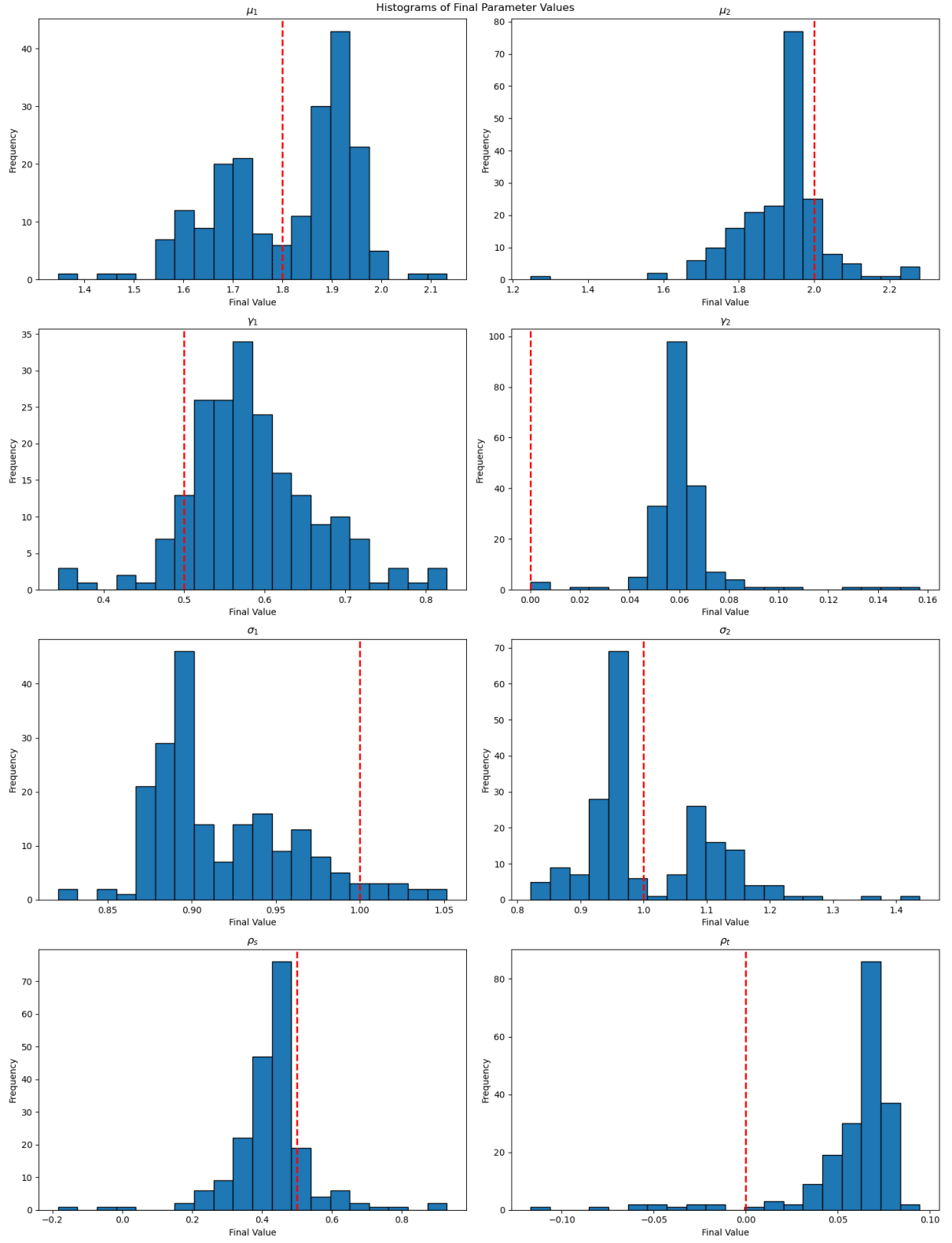
Figure 4 shows the results: In the top row, a neural network discriminator approximates the JS divergence. In the bottom row, the Wasserstein-1 distance as approximated by the Sinkhorn divergence is plotted. The left column shows clearly that the Jensen-Shannon divergence becomes constant when  $m$  is small or big enough that the distributions become disjoint, while the Wasserstein-1 distance provides constant gradients. The right column plots the diagonal  $\text{diag}_{1,-2} = \{(\mu_1, -\mu_2) = (m, m); m \in \mathbb{R}\}$ . Disjointness is also realized in this case for large absolute values of  $m$ .

## 4.5 Estimation

### 4.5.1 Estimation with the Wasserstein-1 divergence

I reproduce the simulations in section 3.2.2 of Kaji, Manresa, and Pouliot (2023) using the approximate Wasserstein-1 estimator implemented in `geomloss.SamplesLoss` and 200 bootstrap samples. The first row of table 1 shows the results. The point estimates are always at least as close to the true parameter values as those reported in Table I of Kaji, Manresa, and Pouliot (2023). The standard errors are also comparable or tighter except for  $\mu_1$  and  $\mu_2$ .





**Figure 6.** Results of 200 bootstrap estimations with Wasserstein-1 loss

**Table 1.** Parameter estimates

	$\mu_1$	$\mu_2$	$\gamma_1$	$\gamma_2$	$\sigma_1$	$\sigma_2$	$\rho_s$	$\rho_t$
Wasserstein-1-discriminator	1.81 (0.14)	1.91 (0.12)	0.58 (0.08)	0.06 (0.02)	0.92 (0.04)	1.01 (0.10)	0.43 (0.12)	0.06 (0.03)
True values	1.80	2.00	0.50	0.00	1.00	1.00	0.00	0.50

**Table 2.** Parameter estimates and standard errors with uniform initialization in wide intervals

	$\mu_1$	$\mu_2$	$\gamma_1$	$\gamma_2$	$\sigma_1$	$\sigma_2$	$\rho_s$	$\rho_t$
W1 uniform init	-0.97 (4.48)	-0.74 (4.19)	-1.31 (4.72)	-2.36 (4.95)	1.78 (1.88)	2.05 (1.99)	-0.10 (0.67)	-0.05 (0.53)
W2 uniform init	-0.99 (4.54)	-0.57 (4.27)	-1.93 (5.15)	-1.92 (4.77)	1.83 (1.79)	1.72 (1.70)	-0.00 (0.71)	-0.04 (0.54)
JS uniform unit	-0.46 (6.03)	-0.17 (5.96)	-0.11 (5.81)	0.64 (5.88)	4.92 (2.92)	5.14 (2.89)	-0.08 (0.59)	0.05 (0.57)
True values	1.80	2.00	0.50	0.00	1.00	1.00	0.00	0.50

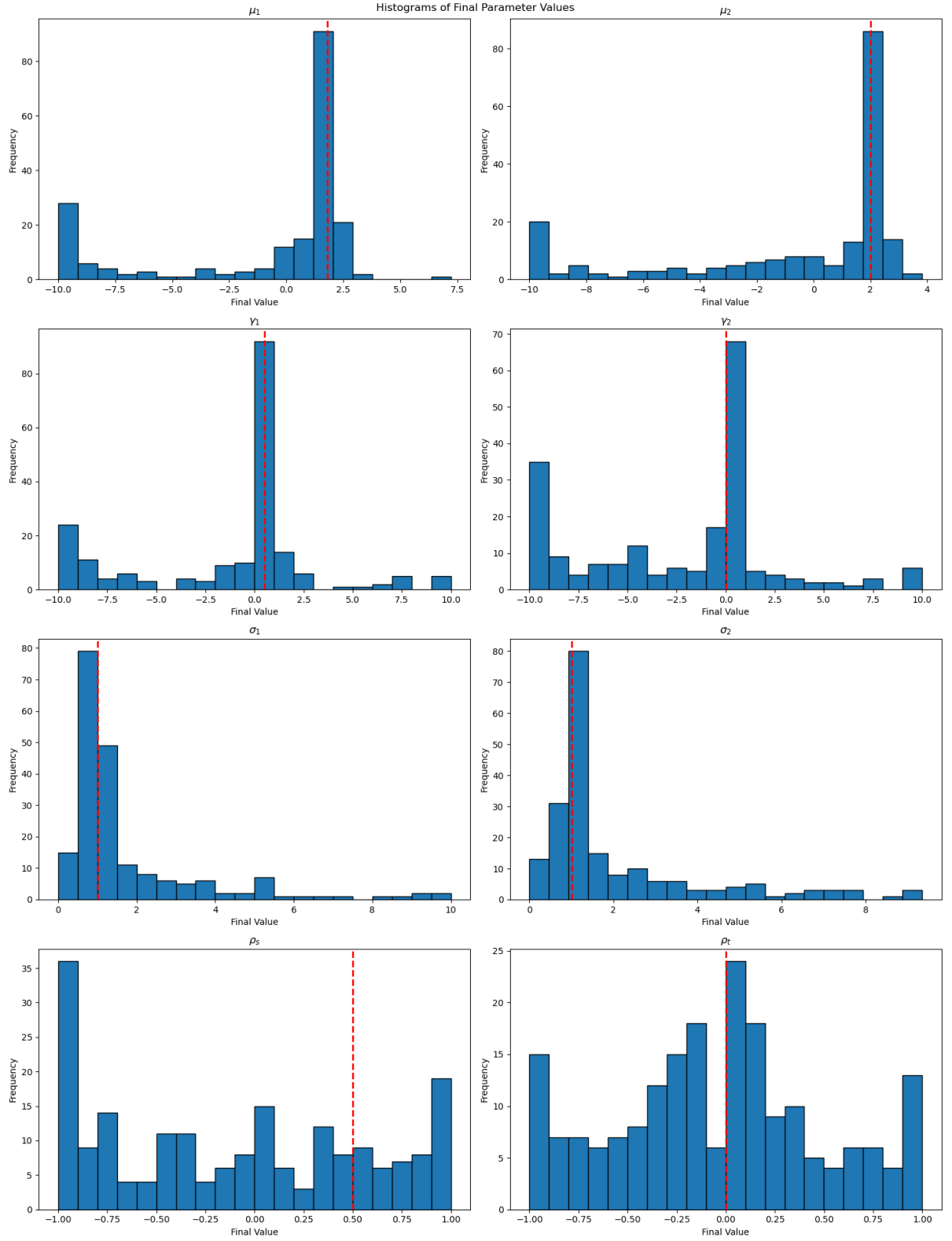
Figure 6 visualizes the distribution of estimates.

#### 4.5.2 Uniform initialization in wide intervals

The Wasserstein estimators should have an advantage over those using Jensen-Shannong divergence when starting far away from the true value. To verify this, I run a simulation in which initial values are drawn over a uniform distribution from a set of wide parameter bounds:

- $\mu_{1,0} \in [-10, 10]$
- $\mu_{2,0} \in [-10, 10]$
- $\gamma_{1,0} \in [-10, 10]$
- $\gamma_{2,0} \in [-10, 10]$
- $\sigma_{1,0} \in [0, 10]$
- $\sigma_{2,0} \in [0, 10]$
- $\rho_{s,0} \in [-1, 1]$
- $\rho_{t,0} \in [-1, 1]$

Note that the wages in the Roy model are in log space. Under a literal interpretation it is not realistic that, for example,  $\exp(\mu_{1,0})$ , the initial guess of the average wage in sector one, would be even  $e^5 \approx 148$  times too low. The goal of this simulation is to generally investigate



**Figure 7.** Results of 200 estimations with Wasserstein-1 loss and initial values that are uniform over broad intervals

the robustness of the estimators to starting values that are far from the true value. It can also be understood as an imperfect metaphor for higher-dimensional models, where the problem of non-overlapping supports will become even more severe.

For the Wasserstein-1 estimator, the results are summarized in the second row of table 2 and plotted in figure 7. The numerical comparison shows that they are much less precise than the estimates considered in the previous section. However, checking the histograms reveals that the mode of the estimate is clearly distinct and close to the true value for the location parameters and the variance. There is also a mode, albeit less clear, near the true  $\rho_t$ , but  $\rho_s$  does not get estimated at all. In all cases, the distribution of estimates is quite wide, and in some cases, a significant amount is at the bounds of the considered intervals.

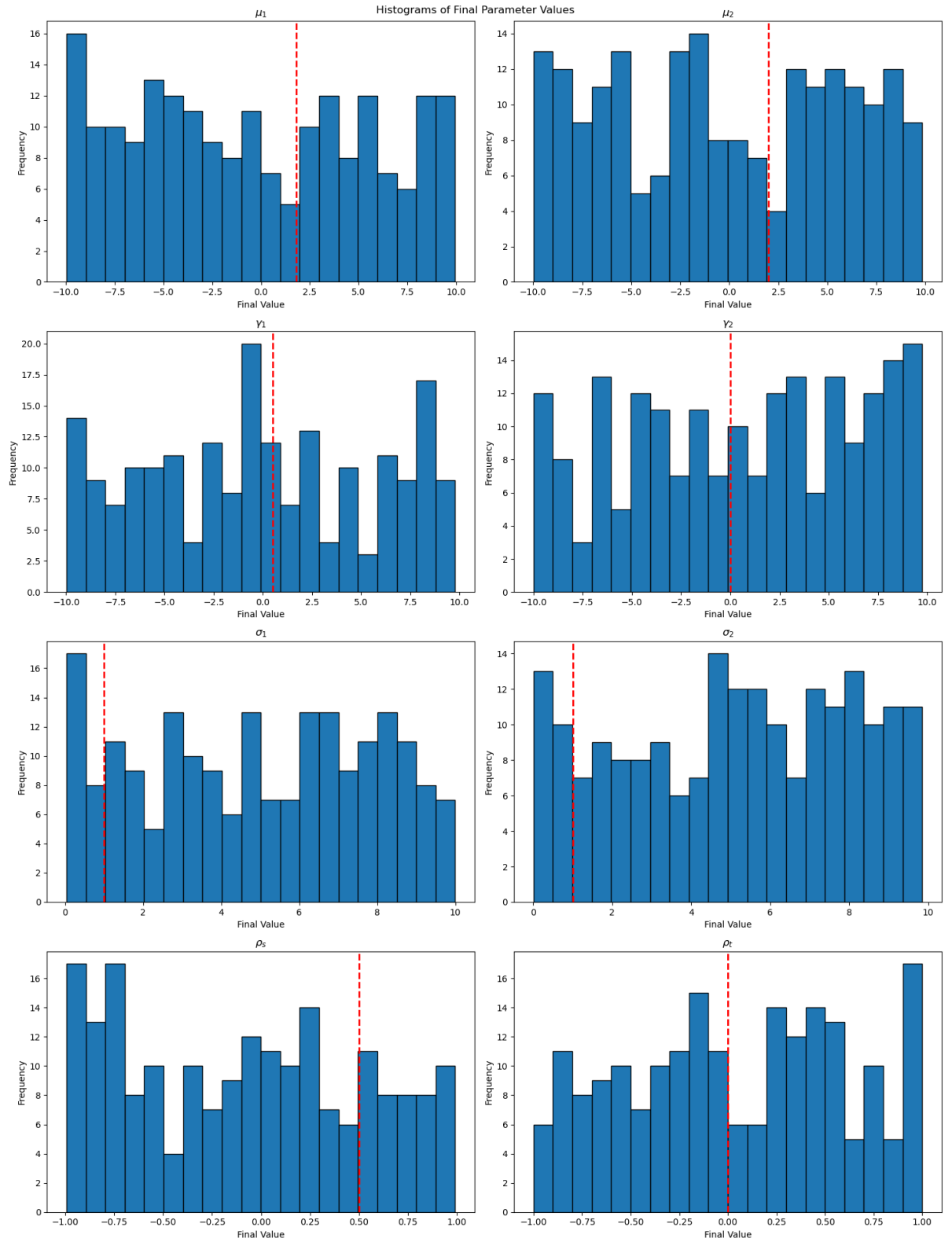
I repeat the same experiment with the Wasserstein-2 estimator. The results are reported in table 2 and figure A.1 in appendix A. They show a similar picture, so the smoother loss surface of the Wasserstein-2 discriminator does not bring any meaningful advantage in this simulation study.

The comparison with the Jensen-Shannon estimator is striking. The values reported in table 2 as well as the histograms plotted in figure 8 show that the estimates have barely improved from their uniform initialization. Indeed, inspecting the simulation outcomes reveals that the estimator has converged in the default time in 0 out of 200 repetitions.

## 5 Conclusion

I have argued and demonstrated in the simulation study that the Wasserstein distances (and their approximation by the Sinkhorn divergences) can improve the adversarial estimator. The central reason for this is that the Wasserstein distances provide non-flat gradients to the optimizer of the generating model in situations where the Jensen-Shannon divergence fails. In practice, the Sinkhorn approximation implemented by Feydy et al. (2019) can also be calculated quickly if a GPU is available. In my simulation study, it was much faster than training a neural network to get the Jensen-Shannon divergence. While the central idea of Kaji, Manresa, and Pouliot (2023) was to take I. J. Goodfellow et al. (2014) and replace the generator network with a structural economic model, Feydy et al. (2019) have allowed me to also replace the discriminator network with a faster approximation.

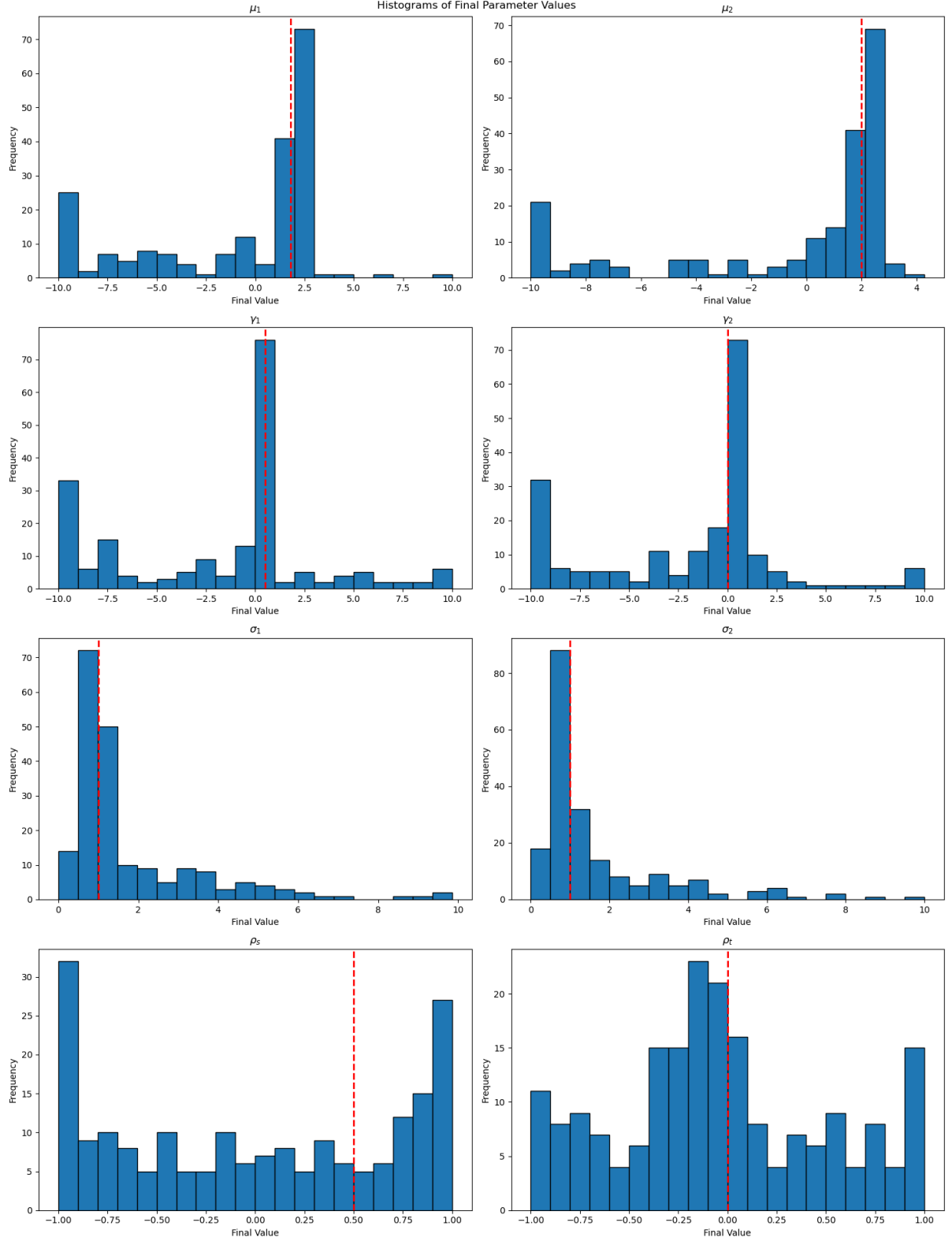
One clear direction for future work is to investigate the theoretical properties of the Wasserstein/Sinkhorn adversarial estimators. A more open ended direction is however to keep an eye on the continuing developments in machine learning and artificial intelligence. Perhaps they will



**Figure 8.** Results of 200 estimations with Jensen-Shannon loss and initial values that are uniform over broad intervals

continue to be, like the GAN, which was published 10 years ago, “useful models” that inspire econometric research, even if there are no (black-box) neural networks left at the end.

# Appendix A Wasserstein-2 estimation results with uniform initialization



**Figure A.1.** Results of 200 estimations with Wasserstein-2 loss and initial values that are uniform over broad intervals

## **Appendix B   Acknowledgement of system use**

The author gratefully acknowledges the granted access to the Marvin cluster hosted by the University of Bonn.



## References

- Ansel, Jason, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, et al. 2024. “PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation.” In *29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24)*. ACM. <https://doi.org/10.1145/3620665.3640366>. [13]
- Arjovsky, Martin, Soumith Chintala, and Léon Bottou. 2017. *Wasserstein GAN*. arXiv: 1701.07875 [stat.ML]. [1, 8]
- Athey, Susan, Guido W Imbens, Jonas Metzger, and Evan Munro. 2021. “Using Wasserstein Generative Adversarial Networks for the design of Monte Carlo simulations.” *Journal of Econometrics*, 105076. [1, 4, 8]
- Box, George E. P. 1976. “Science and Statistics.” *Journal of the American Statistical Association* 71 (356): 791–99. <https://doi.org/10.1080/01621459.1976.10480949>. [1]
- Charlier, Benjamin, Jean Feydy, Joan Alexis Glaunès, François-David Collin, and Ghislain Durif. 2021. “Kernel Operations on the GPU, with Autodiff, without Memory Overflows.” *Journal of Machine Learning Research* 22 (74): 1–6. <http://jmlr.org/papers/v22/20-275.html>. [13]
- D’Errico, John. 2024. *fminsearchbnd*, *fminsearchcon*. <https://www.mathworks.com/matlabcentral/fileexchange/8277-fminsearchbnd-fminsearchcon>. MATLAB Central File Exchange. Accessed September 12, 2024. [16]
- Feydy, Jean, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and Gabriel Peyré. 2019. “Interpolating between Optimal Transport and MMD using Sinkhorn Divergences.” In *The 22nd International Conference on Artificial Intelligence and Statistics*, 2681–90. [13, 24]
- Gao, Fuchang, and Lixing Han. 2012. “Implementing the Nelder-Mead simplex algorithm with adaptive parameters.” *Computational Optimization and Applications* 51 (1): 259–77. [17]
- Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. “Generative adversarial networks.” *Communications of the ACM* 63 (11): 139–44. [4]
- Goodfellow, Ian J., Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. *Generative Adversarial Networks*. arXiv: 1406.2661 [stat.ML]. [1, 4, 7, 24]
- Gulrajani, Ishaan, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. 2017. *Improved Training of Wasserstein GANs*. arXiv: 1704.00028 [cs.LG]. [8]
- Harris, Charles R., K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, et al. 2020. “Array programming with NumPy.” *Nature* 585 (7825): 357–62. <https://doi.org/10.1038/s41586-020-2649-2>. [13]
- Hunter, J. D. 2007. “Matplotlib: A 2D graphics environment.” *Computing in Science & Engineering* 9 (3): 90–95. <https://doi.org/10.1109/MCSE.2007.55>. [13]

- Kaji, Tetsuya, Elena Manresa, and Guillaume Pouliot.** 2023. “An adversarial approach to structural estimation.” *Econometrica* 91 (6): 2041–63. [1, 4–6, 9, 11–15, 17, 20, 24]
- Kingma, Diederik P., and Jimmy Ba.** 2017. *Adam: A Method for Stochastic Optimization*. arXiv: 1412.6980 [cs.LG]. [3, 5, 16]
- Lagarias, Jeffrey C, James A Reeds, Margaret H Wright, and Paul E Wright.** 1998. “Convergence properties of the Nelder–Mead simplex method in low dimensions.” *SIAM Journal on optimization* 9 (1): 112–47. [16]
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, et al.** 2011. “Scikit-learn: Machine Learning in Python.” *Journal of Machine Learning Research* 12: 2825–30. [13]
- Python Software Foundation.** 2023. *Python Language Reference, version 3.12*. Python Software Foundation. <https://www.python.org/>. [13]
- Vaart, A. W. van der, and Jon A. Wellner.** 2023. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer International Publishing. <https://doi.org/10.1007/978-3-031-29040-4>. [9]
- Virtanen, Pauli, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, et al.** 2020. “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python.” *Nature Methods* 17: 261–72. <https://doi.org/10.1038/s41592-019-0686-2>. [13]

# **Selbstständigkeitserklärung**

Ich versichere hiermit, dass ich die vorstehende Masterarbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe, dass die vorgelegte Arbeit noch an keiner anderen Hochschule zur Prüfung vorgelegt wurde und dass sie weder ganz noch in Teilen bereits veröffentlicht wurde. Wörtliche Zitate und Stellen, die anderen Werken dem Sinn nach entnommen sind, habe ich in jedem einzelnen Fall kenntlich gemacht.

30. September 2024

Marvin Benedikt Riemer