

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

List of Figures	iii
List of Tables	iv
1 Introduction	1
2 Background	1
2.1 Structural estimation	1
2.2 Neural networks	1
3 Adversarial estimation	2
3.1 Examples of discriminators	3
3.2 Generator objectives	4
3.2.1 Jensen-Shannon divergence	4
3.2.2 Wasserstein-p distance	5
4 Simulation	6
4.1 The Roy model	6
4.2 General simulation structure	7
4.3 Implementation details	7
4.3.1 Discriminators	7
4.3.2 Generators	8
4.4 Cross-sections of the loss landscape	8
4.5 Estimation	9
5 Conclusion	9
Appendix A Acknowledgement of system use	10
References	11

List of Figures

List of Tables

1 Introduction

Welcome to my thesis! It is based on the paper Kaji, Manresa, and Pouliot (2023).

2 Background

2.1 Structural estimation

Consider the problem of estimating the parameters of a structural economic model. For $k \in \{1, \dots, K\}$, let

$$Y_k = f_{\theta}(X_k, Z_k;), \quad (1)$$

where Y_k is a vector of outcome variables influenced by a vector of noise variables Z_k . The strength and functional form of the relationships between the variables is defined by a function f and its parameters θ .

A common approach to this problem is maximum likelihood estimation, that is, to find $\hat{\theta}$ such that

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \mathcal{L}_n(\theta; \mathbf{y}). \quad (2)$$

However, for some more sophisticated economic models, it is not easy or even possible to calculate the likelihood function.

This motivates further approaches, such as simulation methods, which attempts to infer θ based on a simulation of the true data. Most notable among these is perhaps the simulated method of moments.

The question naturally arises of how to judge whether the simulated distribution comes sufficiently close to the real distribution. This is one motivation for adversarial estimation. It is also intuitive that one approach to this involves classification. In machine learning, a popular tool for classification are neural networks, which I introduce next.

2.2 Neural networks

Definition

Training

3 Adversarial estimation

The basic idea of adversarial estimation is to structure the parameter estimation around two auxiliary models, called the *generator* and the *discriminator* (or *critic*). The generator $G(\hat{\theta}) : Z \rightarrow O$ creates simulated data based on a guess of the true parameter value $\hat{\theta}$. Given the real and the simulated data, the discriminator returns objective function for the generator, which I will call *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. In the classic GAN case involving neural networks, the loss function is the result of the discriminator solving a maximization problem:

$$\hat{\theta}_{adv} = \arg \min_{\theta \in \Theta} \max_{D \in \mathcal{D}} \text{loss}(D(X_i, G(\theta))). \quad (3)$$

Note that this game has a clear Nash-Equilibrium

This method is a variant of “Generative Adversarial Networks”, 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 and instead of estimating a parameter vector, noise is transformed into some output, such as an image. While GANs 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 nodes being fully connected and activation functions being used. 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 θ is what makes adversarial estimation meaningful. It wouldn’t be if θ were a long list of the weights and biases in a multi-layered neural network.

An implementation of [3](#) looks, generally, like [algorithm 1](#).

Algorithm 1 Adversarial estimation

```
Set necessary hyperparameters and initial values
Sample real observations
while Stopping criterion does not hold do
    Generate fake observations from the current generator
    Train the discriminator given the fake observations
    Calculate the loss
    Update  $\hat{\theta}$ 
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. The 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.

Algorithm 1 in Kaji, Manresa, and Pouliot (2023) illustrates one way to fill out the details of 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 algorithm (Diederik (2014)), which requires setting a range of hyperparameters. Their simulation code shows another way. There, they compare a range of estimators (including neural networks trained to completion) and update the generator using a gradient-free approach.

Now I discuss some of these terms in detail.

3.1 Examples of discriminators

All the following discriminators have in common that 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 unique Nash equilibrium from If the true densities p_0 and $p_\theta(x)$ are known, we get the *oracle discriminator*.

Definition 1 (Oracle discriminator). The **oracle discriminator** assigns

$$D_\theta(x) := \frac{p_0(x)}{p_0(x) + p_\theta(x)} \quad (4)$$

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

Kaji, Manresa, and Pouliot (2023) call this the *oracle discriminator*. Of course, p_0 and $p_\theta(x)$ are unknown in practice. Also, this discriminator is only optimal in the Nash equilibrium as off-equilibrium, it neglects to update from the prior probabilities p_0 and $p_\theta(x)$. Nevertheless, it is useful as a benchmark in simulations and has an interesting theoretical property: If the simulated sample size $m \rightarrow \infty$, θ_{oracle} approaches θ_{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 2 (Logistic discriminator). Let Λ be a sigmoid function with values in $(0, 1)$, and x^{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^{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

$$D(x) = \Lambda(\beta_0 + \sum_{k=1}^K \beta_k x_k^{mom}) \quad (5)$$

to every $x \in \cdot$.

Note that this classifier has to be calculated anew after each update of θ . 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 we different training procedures might result in different discriminators.

Definition 3 (Neural network discriminator). Define a classifier neural network $\mathcal{N} : \mathcal{X}^k \rightarrow [0, 1]$ by:

$$\mathcal{N}(x) = \sigma_L(W_L \sigma_{L-1}(W_{L-1} \cdots \sigma_1(W_1 x + b_1) \cdots + b_{L-1}) + b_L), \quad (6)$$

where $x \in \mathbb{R}^n$ is the input vector, L is the number of layers, W_i are weight matrices, b_i are bias vectors, σ_i are activation functions. Let θ_{train} be a collection of hyperparameters that specify the training of a neural network, such as a loss function, an initial guess, a number of training steps, a training algorithm, and a list of hyperparameters of the latter. Call $\mathcal{N}(\theta_{train})$ the network \mathcal{N} after it has been trained according to θ_{train} . Then the **neural network discriminator** assigns

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

to every $x \in \cdot$.

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 $D(x)$ into an objective for the generator is the following:

Definition 4 (Cross-entropy loss). The empirical cross-entropy loss (CE) is:

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

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

Theorem 5 (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, the 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. Luckily, there are ...

3.2.2 Wasserstein-p distance

Using the so-called Wasserstein-1 distance as an optimization target for the generator was first proposed by Arjovsky, Chintala, and Bottou (2017).

Definition 6 (Wasserstein-p distance). For two probability distribution \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}_r 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 7 (Kantorovich-Rubinstein duality).

$$W_1(\mathbb{P}_0, \mathbb{P}_G) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_0}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_G}[f(x)] \quad (8)$$

This dual representation of the Wasserstein-1 distance paves the way to approximating it using a neural network. The network estimates 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 than be approximated using a smoothed “soft-auction” algorithm. The divergence resulting from this algorithm is called the *Sinkhorn divergence*.

4 Simulation

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 chosing which sector to work in in each of two time periods. At the start of the game, nature determines the wages offered to each agent, by the following formulas: In the first time period, each agent i observes the (natural logarithms of the) wages $\log w_{i11}$ and $\log w_{i12}$ offered to them in the two sectors. Knowing their own discount factor β and the parameters γ_1 and γ_2 , they solve the dynamic programming problem and pick a sector for the first period. In the second period, they observe the wages $\log w_{i21}$ and $\log w_{i22}$ and pick a sector.

The researcher observes the realized wages $\log w_1$ and $\log w_2$ and as well as the corresponding sector choices $d_1, d_2 \in 1, 2$.

4.2 General simulation structure

First, I reproduce parts of the author’s simulation in the scientific Python stack, more precisely, using the packages `numpy`, `scipy`, and `scikit-learn` (Harris et al. (2020), Virtanen et al. (2020), and Pedregosa et al. (2011), respectively). My code is available

I program another simulation using PyTorch (Ansel et al. (2024)), a popular and highly developed neural network library for Python. Among the advantages of Pytorch is that it offers support for training neural networks on GPUs. Additionally, the library `GeomLoss` (Feydy et al. (2019)) builds on PyTorch. Its `SampleLoss` function provides the Sinkhorn approximations of the Wasserstein-1 and -2 distances. It is also optimized for running on GPUs by virtue of being built on `KeOps` (Charlier et al. (2021)). I generate plots using `Matplotlib` (Hunter (2007)).

The replication package can be downloaded from the journal website. It contains the author’s 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 (e.g. 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 created Figures 6 and 8. Thereafter, real and fake observations are generated and the estimation is performed on them. It is implemented as a constrained minimization of a loss function which in turn calculates the discriminators. 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.

4.3 Implementation details

4.3.1 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 author’s 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. Because the conjugate-gradient descent algorithm is not available to train `MLPClassifier`, I use the Adam algorithm (Diederik (2014)). It is popular for training neural networks and achieves comparable results in my case.

`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.

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

My simulations show that these modifications do not significantly alter the shape of the loss curves.

4.3.2 Generators

For the outer optimization loop that trains the generator, the authors use the third-party `fmin-searchcon` 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 algorithm with bounds and nonlinear constraints. I set an option to perform a version of the Nelder-Mead algorithm that's adapted to higher-dimensional problems, which shows improved convergence in my simulation.

I employ the `mp` module from Python's standard library to parallelize simulation runs on an HPC cluster (cf. A).

4.4 Cross-sections of the loss landscape

I reproduce the cross-section of the loss landscape from Figure

Looking at the wider intervals, it is striking that for some parameters, the loss is flat when moving to far away from the optimal value. Partly, this can be explained by the dynamics of the Roy model.

4.5 Estimation

5 Conclusion

This section concludes.

Appendix A 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>. [7]
- Arjovsky, Martin, Soumith Chintala, and Léon Bottou. 2017. *Wasserstein GAN*. arXiv: 1701.07875 [stat.ML]. [5, 6]
- 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. [6]
- 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>. [7]
- D’Errico, John. 2024. *fminsearchbnd, fminsearchcon*. <https://www.mathworks.com/matlabcentral/fileexchange/8277-fminsearchbnd-fminsearchcon>. MATLAB Central File Exchange. Accessed September 12, 2024. [8]
- Diederik, P Kingma. 2014. “Adam: A method for stochastic optimization.” (*No Title*). [3, 8]
- 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. [7]
- 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. [2]
- Goodfellow, Ian J., Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. *Generative Adversarial Networks*. eprint: [arXiv:1406.2661](https://arxiv.org/abs/1406.2661). [2, 5]
- Gulrajani, Ishaan, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. 2017. *Improved Training of Wasserstein GANs*. arXiv: 1704.00028 [cs.LG]. [6]
- 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>. [7]
- 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>. [7]
- Kaji, Tetsuya, Elena Manresa, and Guillaume Pouliot. 2023. “An adversarial approach to structural estimation.” *Econometrica* 91 (6): 2041–63. [1–3]
- 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. [8]

- 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. [7]
- 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>. [7]

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.

21. September 2024

Marvin Benedikt Riemer