# Neural Network solutions to Witsenhausen problem

# Jiaojiao Fan GTID:903565753

Email: jiaojiaofan@gatech.edu

Abstract—In this report, several neural networks with different structures are implemented to solve Witsenhausen problem. Other improving strategies include optimizers, initializations and forced function fixing. Finally, the result are compared with former people and a better result is obtained. Also, the shortcoming of the neural network also shows in this project. The neural network may be stuck into a near local minima.

### 1. Introduction

In this report, we proposed several solutions to the well-known and still unsolved Witsenhausen counterexample. [1] There have been some meaningful tries to detect the global minima of the min problem, such as Lee [2]and M. Barglietto [3] Some of their manipulations are also refered in this project. Other than that, thanks to the development of the neural networks, many other meaningful attempts are also taken such as input convex neural network (ICNN) structure [4] Different results would be listed to show the effect.

## 2. The Witsenhausen Counterexample

The Witsenhausen counterexample has been outstanding for more than 50 years. It is formulated by Hans Witsenhausen in 1968. [1] It is a counterexample to a natural conjecture that in a system with linear dynamics, Gaussian disturbance, and quadratic cost, affine control laws are optimal to minimize the cost. However, Witsenhausen courterexample, shown in figure below, has nonlinear control laws

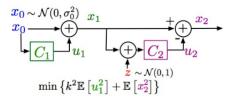


Figure 1. Witsenhausen couterexample

that outperform all linear laws.

$$f(x) = \gamma_1(x) + x \qquad g(x) = \gamma_2(x) \tag{1}$$

As a result,  $f(x_0) = x_1$  and  $g(x_1 + z) = u_2$ . Out goal is to minimize the quadratic cost: $k^2 \mathbb{E}[U_1^2] + \mathbb{E}[X_2^2]$ , which can also be written as

$$\min J(f,g) := k^2 \mathbb{E}[f(X_0) - X_0]^2 + \\ \mathbb{E}[(f(X_0) - g(f(X_0) + N))^2]$$
 (2)

In equation (2), there is a parameter  $k^2$ , which in fact determines the cost gap between the linear controller and nonlinear controller [3]. If  $k^2$  is smaller, the gap is bigger. For better comparison with the results got by previous researchers,  $k^2$  is set as 0.04 in this report.

In addition, it is already known that f(x) must have some strict property to be optimal [5]:

• Any optimal controller f is a strictly increasing unbounded piecewise real analytic function with a real analytic inverse

This means f has to be smooth enough. But interestingly, the neural network(NN) optimized result is exactly opposite from this property. The sharper the f becomes (opposite to smooth), the smaller the cost is.

## 3. Basic Neural Network Setup

The whold process could be generally separated into 4 parts:

- 1) Initialization setup for the f net and g net
- 2) Train the NNs using Gaussian distribution data. In this report, all data keeps the consistency:  $x_0 \sim \mathcal{N}(0, \sigma^2)$ , where  $\sigma = 5$ , and  $N \sim \mathcal{N}(0, 1)$ .
- 3) Fix f net and continue to train q net.

#### 3.1. Neural Network Architecture

Basically, two NNs are taken to represent f and g seperately using Pytorch structure. <sup>1</sup> For f net and g net, all layers are linear layers and activated by CELU [6] function, i.e.

$$CELU(x) = max(0, x) + min(0, \alpha * (exp(x/\alpha) - 1))$$
 (3)

is used since it makes the activation function continuously differentiable and improves the performance in initialization setup process. It is worth noting that CELU is convex and monotone increasing activation function. For f net structure,

1. The code of this project could be found at: https://github.com/sbyebss/Witsenhausen

we tried ICNN as f net's integral function: F net, i.e. f NN works as the derivative function of the function represented by F NN. We also have tried ResNet [7] since it performs much better than ordinary linear layer NNs.

## 3.2. Optimizer

For updating parameters of two NNs, we first use Stochastic Gradient Descent (SGD) and then used ADAM. In comparison, SGD performs much slower and becomes not stable while entering plateau of loss decreasing. ADAM increases the stability and speed a lot. So we focus on the better optimizer compared to ADAM later. There were a lot of variation of ADAM during the past several years. We mainly care whether the optimizer could lead the NN to the global optimizer. Some outstanding optimizers came out like AdaBound [8], RAdam [9] and Yogi [10]. We tried to use RAdam and Yogi but found Yogi didn't work well in this problem picture at all. So we finally proposed to use RAdam for main part of the project. But we will still list the result of RMSprop and Adam as the comparison.

## 3.3. Initialization

The f net and g net both have initialization. According to the previous researchers' work, we choose the 7-step-stair as basic two NN initialization. This is referring to the best result got from past researchers and the step parameters are from Yu-Chi Ho's group [9] as below:

$$f(x) = \begin{cases} 0 & 0 \le x < 3.25 \\ 6.5 & 3.25 \le x < 9.90 \\ 13.2 & 9.90 \le x < 16.65 \\ 19.9 & 16.65 \le x. \end{cases}$$
 (4)

According to symmetry, the  $x \ge 0$  part of f(x) could be obtained. We also tried segmented stair, which parameters are also obtained from Yu-Chi Ho's group [2]:

$$f(x) = \begin{cases} 0.00 & 0.00 \le x < 0.65 \\ 0.05 & 0.65 \le x < 1.95 \\ 0.10 & 1.95 \le x < 3.25 \\ 6.40 & 3.25 \le x < 4.58 \\ 6.45 & 4.58 \le x < 5.91 \\ 6.50 & 5.91 \le x < 7.24 \\ 6.55 & 7.24 \le x < 8.57 \\ 6.60 & 8.57 \le x < 9.90 \\ 13.10 & 9.90 \le x < 11.25 \\ 13.15 & 11.25 \le x < 12.60 \\ 13.20 & 12.60 \le x < 13.95 \\ 13.25 & 13.95 \le x < 15.30 \\ 13.20 & 15.30 \le x < 16.65 \\ 19.90 & 16.65 \le x. \end{cases}$$
 (5)

2. The pytorch optimizer source code could be found here: https://github.com/jettify/pytorch-optimizer

To distinguish, we call (4) as the 7-stair-init and (5) as the segmented-7-stair-init.

## 3.4. Test Data Set

In the process of training, we choose to get the test set loss for every 1000 epoches. Therefore, a good test data size need to be chosen. The variance of the empirical test loss is expected to be less than  $0.005^2$  to be compared successfully with previous results. If take the test data size as  $5 \times 10^5$ , the empirical variance of the test loss is 8.65791389375e-7. Finally, we chose the test data size as  $1 \times 10^5$ 

### Information for code

Different train process are edited in train\_specificName.py file.

The *modules* folder saves the network structure and ICNN weight clamping function.

The *runs* folder record the tensorboard loss file. Importing those files into tensorboard, we can compare the loss decreasing trend in one plot.

The *data* folder saves the main body training data as well as initialization training data.

The *test* folder saves some test files which verified some important rudimentary ideas.

The *modle* folder saves the different NN parameters for initialization or some trained solution NNs.

## Acknowledgment

Thanks to my advisor: Dr. Yongxin Chen and the other PhD studens: Rahul Singh and Qinsheng Zhang for their kind help. They give me a lot of ideas about how to improve the performance during the process of project.

#### References

- [1] H. S. Witsenhausen, "A counterexample in stochastic optimum control," *SIAM Journal on Control*, vol. 6, no. 1, pp. 131–147, 1968.
- [2] J. T. Lee, E. Lau, and Y.-C. Ho, "The witsenhausen counterexample: A hierarchical search approach for nonconvex optimization problems," *IEEE Transactions on Automatic Control*, vol. 46, no. 3, pp. 382–397, 2001.
- [3] M. Baglietto, T. Parisini, and R. Zoppoli, "Numerical solutions to the witsenhausen counterexample by approximating networks," *IEEE Transactions on Automatic Control*, vol. 46, no. 9, pp. 1471–1477, 2001.
- [4] B. Amos, L. Xu, and J. Z. Kolter, "Input convex neural networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 146–155.
- [5] Y. Wu and S. Verdú, "Witsenhausen's counterexample: A view from optimal transport theory," in 2011 50th IEEE Conference on Decision and Control and European Control Conference. IEEE, 2011, pp. 5732–5737.
- [6] J. T. Barron, "Continuously differentiable exponential linear units," arXiv preprint arXiv:1704.07483, 2017.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on* computer vision and pattern recognition, 2016, pp. 770–778.

- [8] L. Luo, Y. Xiong, Y. Liu, and X. Sun, "Adaptive gradient methods with dynamic bound of learning rate," arXiv preprint arXiv:1902.09843, 2019.
- [9] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, "On the variance of the adaptive learning rate and beyond," arXiv preprint arXiv:1908.03265, 2019.
- [10] M. Zaheer, S. Reddi, D. Sachan, S. Kale, and S. Kumar, "Adaptive methods for nonconvex optimization," in *Advances in neural infor*mation processing systems, 2018, pp. 9793–9803.