

E-Companion to *Using Neural Networks To Guide Data-driven Operational Decisions*

Saman Lagzi, Ningyuan Chen, Joseph Milner

Table EC.1 Table of variables and parameters.

Acronym	Meaning
f_*	True unknown objective function
Y^i	A noisy observation of f_* in epoch i
\mathbf{x}	Observed covariates
\mathbf{z}	Decision taken
d_X	Dimension of the observed covariates
d_Z	Dimension of the decision taken
f_θ	A generic DNN function
$f_{\hat{\theta}}$	A realized fitted DNN function
$\sigma(\cdot)$	The activation function
$\sigma'_{h,l}(\mathbf{x}, \mathbf{z})$	The point-wise first derivative of the activation function at the h -th node and l -th layer
$\sigma''_{h,l}(\mathbf{x}, \mathbf{z})$	The point-wise second derivative of the activation function at the h -th node and l -th layer
$\mathbf{z}^*(\mathbf{x}_0)$	The true optimal solution given \mathbf{x}_0
$\hat{\mathbf{z}}(\mathbf{x}_0)$	The DNN based optimal solution given \mathbf{x}_0
∂	The partial derivative
K_*	The Lipschitz constant of f_*
$K_{\hat{\theta}}$	The Lipschitz constant of $f_{\hat{\theta}}$
K	The Lipschitz constant of $\ f_* - f_{\hat{\theta}}\ $
$\mu_{X,Z}$	The marginal distribution of (\mathbf{X}, \mathbf{Z})
δ_N	The generalization bound with N data points
c_1	The lower bound on the density function of $\mu_{X,Z}$
c_3	Min. radius of a ball away from the optimal solution to have probability measure c_4
c_4	Probability measure of a ball of radius c_3 away from the optimal solution
c_2	Lower bound on the increase of f_* for solutions far away from optimal solution
c_5	Quadratic increase rate of f_* for solutions near the optimal solution
$L_{d_X+d_Z}$	Constant for volume of a $d_X + d_Z$ dimensional ball

EC.1. Using Neural Networks to Fit the Objective Function

As mentioned in Section 3, we propose to find the best fit objective function $\hat{f}(\mathbf{x}, \mathbf{z})$ from the data $\{(\mathbf{X}^i, \mathbf{Z}^i, Y^i)\}_{i=1}^N$ within the class of DNN functions \mathcal{F}_{NN} . That is, assuming a quadratic loss function, and a clear specification of \mathcal{F}_{NN} , we seek:

$$\hat{f} = \arg \min_{f \in \mathcal{F}_{\text{NN}}} \sum_{i=1}^N (f(\mathbf{X}^i, \mathbf{Z}^i) - Y^i)^2. \quad (\text{EC.1})$$

The class \mathcal{F}_{NN} is the set of functions that can be represented through a particular DNN. The function class \mathcal{F}_{NN} is determined by the architecture of the network and the form of interactions between the nodes. We consider *feedforward* DNNs consisting of:

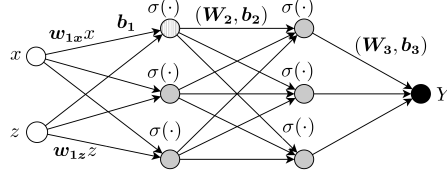


Figure EC.1 The illustration of a fully connected feedforward neural network with $L = 2$ layers and $H = 3$ nodes in each layer. The activation function is applied to the input at each of the nodes.

- L layers (the depth of the network),
- H nodes per layer (the width)¹,
- fully connected layers,
- an activation function $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$.

An example of a network architecture is illustrated in Figure EC.1. The neural network takes as input two scalars (x, z) and outputs a single scalar (y). It has $L = 2$ layers and each layer has $H = 3$ width (note that we do not count the first and the last layers).

The commonly use activation functions are:

1. **sigmoid:** $\sigma(x) = 1/(1 + e^{-x})$
2. **Swish:** $\sigma(x) = x/(1 + e^{-x})$
3. **ReLU:** $\sigma(x) = \max\{0, x\}$.

Given a network architecture $\mathcal{F}_{\text{NN}}(L, H, \sigma)$, a function withing the class is specified by a parameterization $\theta = (\mathbf{W}_{1x}, \mathbf{W}_{1z}, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_L, \mathbf{b}_L, \mathbf{w}_{L+1}, b_{L+1})$, where $\mathbf{W}_{1x} \in \mathbb{R}^{H \times d_x}$, $\mathbf{W}_{1z} \in \mathbb{R}^{H \times d_z}$, $\mathbf{b}_k \in \mathbb{R}^H$ for $k = 1, \dots, L$, $\mathbf{W}_k \in \mathbb{R}^{H \times H}$ for $k = 2, \dots, L$, $\mathbf{w}_{L+1} \in \mathbb{R}^H$ and $b_{L+1} \in \mathbb{R}$.

Given (\mathbf{x}, \mathbf{z}) and a parameter θ ,

$$f_{\theta}(\mathbf{x}, \mathbf{z}) = \mathbf{w}_{L+1}^{\top} \sigma \left(\dots \sigma \left(\mathbf{W}_3 \sigma \left(\mathbf{W}_2 \sigma \left(\mathbf{W}_{1x} \mathbf{x} + \mathbf{W}_{1z} \mathbf{z} + \mathbf{b}_1 \right) + \mathbf{b}_2 \right) + \mathbf{b}_3 \right) + \dots \right) + b_{L+1}.$$

Here, $\sigma(\cdot) : \mathbb{R}^H \rightarrow \mathbb{R}^H$ represents the vectorized activation function that applies $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ elementwise.

By this formulation, \mathcal{F}_{NN} is the set of all functions $f_{\theta}(\mathbf{x}, \mathbf{z})$. One may use Figure EC.1 to better understand $f_{\theta}(\mathbf{x}, \mathbf{z})$. Consider the left-most layer (the input layer), the edges leading to the second layer (also referred to as the first hidden layer) and the term $\mathbf{w}_{1x} \mathbf{x} + \mathbf{w}_{1z} \mathbf{z} + \mathbf{b}_1$. The upper edges coming from node x represent $\mathbf{w}_{1x} \mathbf{x}$ and the lower edges coming from node z represent $\mathbf{w}_{1z} \mathbf{z}$. The vector \mathbf{b}_1 is referred to as the constant term (or bias). Observe that $\mathbf{w}_{1x} \mathbf{x} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^H$. In this case the scalar x is transferred by the weights \mathbf{w}_{1x} into H inputs to the

¹This is without loss of generality since H can be set equal to the number of nodes in the widest layer, and for other layers, one can consider dummy nodes with weight set to zero and the constant parameter set to be equal to the root point of the activation function (0 for ReLU and Swish and a large negative number such as -100 for sigmoid).

second layer, and similarly for z . At each node in the second layer the function $\sigma(\cdot) : \mathbb{R}^1 \rightarrow \mathbb{R}^1$ is applied to the H -dimensional term $\mathbf{w}_{1x}x + \mathbf{w}_{1z}z + \mathbf{b}_1$ element wise to produce H outputs. Then, recursively, in a feedforward fashion, the results are fed to the edges between the second and third layer, represented by \mathbf{W}_2 and \mathbf{b}_2 , and so on. In the output layer, an inner product using weight \mathbf{w}_{L+1} is applied to the output of the penultimate layer (the last hidden layer). Adding the scalar constant term b_{L+1} gives the scalar output.

Given observations $\{(\mathbf{X}^i, \mathbf{Z}^i, Y^i)\}_{i=1}^N$, we seek f_θ that solves (EC.1). That is, we train the DNN by selecting \mathbf{W} 's and \mathbf{b} 's that minimize the quadratic loss. A common approach is to use Stochastic Gradient Descent (SGD) which relies on the observation that the gradient of f_θ with respect to θ can be given in closed form, when σ is differentiable (Rumelhart et al. 1986). SGD has shown good empirical performance (Zhang et al. 2021). See Goodfellow et al. (2016) for an overview of SGD applied to fitting DNNs.

EC.2. Optimize the Fitted Neural Network

In this section we focus on the case where \mathcal{P} is a box and discuss how to choose \mathbf{z} given a covariate \mathbf{X}_0 so that the output of the fitted DNN is optimized. Let $\hat{\theta}$ denote the fitted DNN (the optimized weights \mathbf{W} and constant terms \mathbf{b}). Then, following (5), we would want to find

$$\hat{\mathbf{z}}(\mathbf{X}_0) \in \arg \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z}). \quad (\text{EC.2})$$

We assume

$$\mathcal{P} = \{\mathbf{z} \in R^{d_Z} \mid \underline{\mathbf{Z}}_i \leq \mathbf{z}_i \leq \bar{\mathbf{Z}}_i, \forall i \leq d_Z\},$$

where $\underline{\mathbf{Z}}_i$ and $\bar{\mathbf{Z}}_i$ indicate the lowest and highest values of the i^{th} coordinate of \mathbf{Z} in the data.

Consider Figure EC.1. We are suggesting that we have \mathbf{X}_0 and $f_{\hat{\theta}}$ and want to choose \mathbf{z} to minimize Y . Solving (EC.2) may be NP-hard, depending on the choice of the activation function (Katz et al. 2017, Anderson et al. 2020).

Define $\hat{\mathbf{z}}(\mathbf{x})$ to be a local minimum if there exists $\delta > 0$, such that

$$\forall \mathbf{z} \in \mathcal{P} : \|\mathbf{z} - \hat{\mathbf{z}}(\mathbf{x})\|^2 \leq \delta, f_{\hat{\theta}}(\mathbf{x}, \hat{\mathbf{z}}(\mathbf{x})) \leq f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}). \quad (\text{EC.3})$$

We suggest using a gradient-based approach to finding a local minimum to (EC.2). We can use backpropagation to derive the gradient of $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ with respect to \mathbf{z} in closed form. Hence, we may hope to efficiently find all the local minima to (EC.2). Moreover, we can use the tools developed for SGD in the training of DNNs to find the $\hat{\mathbf{z}}(\mathbf{x})$ that minimizes $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$. Similar to the the backpropagation algorithm (Rumelhart et al. 1986), we use the layer by layer structure of the DNN to calculate the derivative of $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ with respect to \mathbf{z} .

We can find a local minimum of $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ using a gradient descent algorithm. Let $\hat{g} = \frac{\partial f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}$ and let \hat{K} be the Lipschitz constant of $\frac{\partial f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}$ and let $[\mathbf{z}]^+$ be the projection operator onto the feasible set \mathcal{P} , defined as

$$[\mathbf{z}]_i^+ = \begin{cases} \bar{Z}_i & \text{if } z_i \geq \bar{Z}_i, \\ Z_i & \text{if } z_i \leq Z_i, \\ z_i & \text{otherwise.} \end{cases}$$

Algorithm 2 Projected gradient update

Require: \hat{K} , $0 < s < \frac{2}{\hat{K}}$

Require: Starting point decision value $\mathbf{z}^{(0)} \in \mathcal{P}$

$k \leftarrow 1$

while Stopping criterion not met **do**

 Compute gradient: $\hat{g} \leftarrow \frac{\partial f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}^{(k)})}{\partial \mathbf{z}}$

 Apply update: $\hat{\mathbf{z}}^{(k)} \leftarrow \mathbf{z}^{(k)} - s\hat{g}$

 Project back to feasible set: $\mathbf{z}^{(k+1)} \leftarrow [\hat{\mathbf{z}}^{(k)}]^+$

$k \leftarrow k + 1$

end while

There is a large body of literature on how to choose the step size (s) in Algorithm 2 to ensure faster convergence (Bertsekas 1999). However, if $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ is continuously differentiable, \mathcal{P} is a compact and nonempty convex set, and $\frac{\partial f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}$ is Lipschitz continuous with a known Lipschitz constant, Algorithm 2 as stated will converge in the limit to a stationary point of $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ (Bertsekas (1999), Proposition 2.3.2). If $\frac{\partial f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}$ is only locally Lipschitz continuous or its

Lipschitz constant is not known, one can alternatively use the celebrated Armijo rule to determine the step size schedule that will guarantee convergence in the limit (Bertsekas (1999), Proposition 2.3.3).

REMARK EC.1. The differentiability of σ is imperative for Proposition 1. If σ is a non-differentiable function such as ReLU, as pointed out by several recent works such as Bolte and Pauwels (2021) and Bianchi et al. (2022), backpropagation for ReLU networks does not necessarily compute any kind of known derivative or even subdifferential of $f_{\hat{\theta}}$ and may not even return the correct derivative of $f_{\hat{\theta}}$ at points where it is differentiable. Moreover, we do not know of any research that shows in such a case, similar to Proposition 1, convergence to a stationary point of $f_{\hat{\theta}}(\mathbf{x}, \cdot)$ would result. In particular, the results of Bianchi et al. (2022) which are developed for SGD, can be extended to the deterministic case of Algorithm 2 for ReLU activated DNNs, with backpropagation used to calculate \hat{g} . However their results only guarantee convergence to the set

of zeros of the conservative field of $f_{\hat{\theta}}(\mathbf{x}, \cdot)$ and this set is not necessarily identical to the set of its stationary points. (See [Bolte and Pauwels \(2021\)](#) for a comprehensive review). \square

EC.3. The Choice of Activation Functions

The choice of activation functions has a substantial impact on the performance of DNNs, as documented in a number of empirical studies (for example, see [Ramachandran et al. 2017](#)).

Although the field of deep learning is evolving rapidly and new activation functions are constantly proposed, ReLU remains one of the most popular choices ([Ramachandran et al. 2017](#), [Agarap 2018](#), [Nwankpa et al. 2018](#)). In addition to the empirical success, DNNs with ReLU activation functions possess benign theoretical properties, achieving the minimax convergence rate ([Schmidt-Hieber 2020](#)).

In our framework, since the quality of the data-driven decision is closely related to how well f_* can be learned from the data and ReLU has had tremendous empirical success in deep learning, it seems a natural choice to use ReLU. However, we point out a surprising observation that is specific to our framework using DNNs for data-driven decision-making: using ReLU may lead to *worse* data-driven decisions compared to smooth activation functions such as sigmoid and Swish.

We illustrate this observation using the following toy example.

EXAMPLE EC.1 (USING RELU LEADS TO WORSE OUTCOMES). Consider a scenario without covariates and a one-dimensional objective function $f_*(z) = z - z^2$, $z \in [0, 1]$. The function is shown by the dashed lines in [Figure EC.2](#). We compare the ability of two DNNs to fit this function. The DNNs have one hidden layer and 29 nodes. They differ only in their activation function, one with ReLU and one with Swish. We then sample a large number of observations of $f_*(z)$ and fit the DNNs. The ReLU network has a smaller out of sample prediction error (MSE 7.24×10^{-6}) compared to that of the Swish network (MSE 1.86×10^{-5}). After fitting the DNNs, we solve for the optimal \hat{z} for the DNNs. In [Figure EC.2](#) we observe the optimal \hat{z} for the Swish network is much closer to z^* than that of the ReLU network.

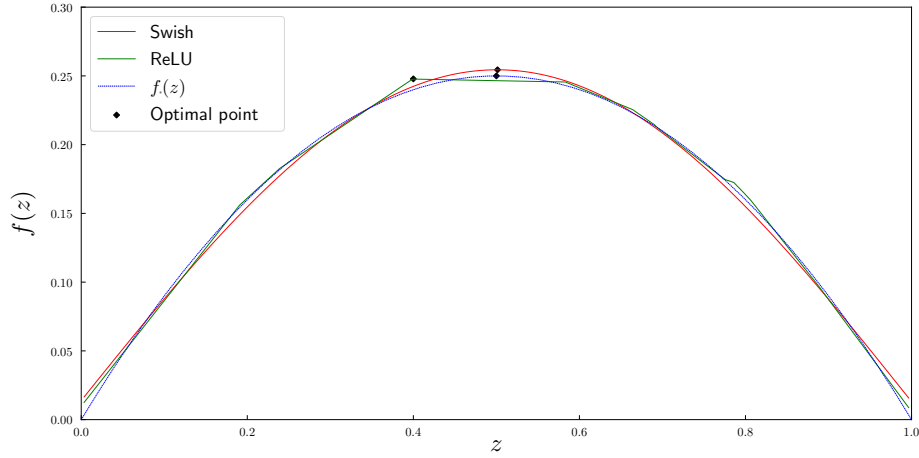


Figure EC.2 The illustration of the prediction value for DNNs using ReLU and Swish activation functions.

EC.4. Proofs

Proof of Proposition 1 As σ is differentiable everywhere $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ is differentiable. Thus, we can apply the chain rule layer by layer, inductively.

Then, we observe

$$\frac{\partial \sigma_{h,1}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \sigma'(\mathbf{w}_{h,1\mathbf{x}}^{\top} \mathbf{x} + \mathbf{w}_{h,1\mathbf{z}}^{\top} \mathbf{z} + b_{h,1}) \cdot \mathbf{w}_{h,1\mathbf{z}}.$$

So

$$\frac{\partial \boldsymbol{\sigma}_1(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \text{diag}(\sigma'_{1,1}(\mathbf{x}, \mathbf{z}), \dots, \sigma'_{H,1}(\mathbf{x}, \mathbf{z})) \times \mathbf{W}_{1\mathbf{z}}.$$

Furthermore,

$$\frac{\partial \boldsymbol{\sigma}_2(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \frac{\partial \boldsymbol{\sigma}_2(\mathbf{x}, \mathbf{z})}{\partial \boldsymbol{\sigma}_1(\mathbf{x}, \mathbf{z})} \times \frac{\partial \boldsymbol{\sigma}_1(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \text{diag}(\sigma'_{1,2}(\mathbf{x}, \mathbf{z}), \dots, \sigma'_{H,2}(\mathbf{x}, \mathbf{z})) \times \mathbf{W}_2 \times \frac{\partial \boldsymbol{\sigma}_1(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}.$$

Assume for some $l \geq 2$, that

$$\frac{\partial \boldsymbol{\sigma}_l(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \text{diag}(\sigma'_{1,l}(\mathbf{x}, \mathbf{z}), \dots, \sigma'_{H,l}(\mathbf{x}, \mathbf{z})) \times \mathbf{W}_l \times \frac{\partial \boldsymbol{\sigma}_{l-1}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}.$$

Then, we have

$$\frac{\partial \boldsymbol{\sigma}_{l+1}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \frac{\partial \boldsymbol{\sigma}_{l+1}(\mathbf{x}, \mathbf{z})}{\partial \boldsymbol{\sigma}_l(\mathbf{x}, \mathbf{z})} \times \frac{\partial \boldsymbol{\sigma}_l(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \text{diag}(\sigma'_{1,l+1}(\mathbf{x}, \mathbf{z}), \dots, \sigma'_{H,l+1}(\mathbf{x}, \mathbf{z})) \times \mathbf{W}_{l+1} \times \frac{\partial \boldsymbol{\sigma}_l(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}}.$$

Then, finally,

$$\frac{\partial f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \frac{\partial f_{\hat{\theta}}}{\partial \boldsymbol{\sigma}_L(\mathbf{x}, \mathbf{z})} \times \frac{\partial \boldsymbol{\sigma}_L(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}} = \mathbf{w}_{L+1}^{\top} \times \frac{\partial \boldsymbol{\sigma}_L(\mathbf{x}, \mathbf{z})}{\partial \mathbf{z}},$$

which completes the proof. \square

Proof of Proposition 2 We note that by assumption $\sigma(\cdot)$ is twice continuously differentiable and hence $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$ is twice continuously differentiable. Moreover, we can apply the chain rule to the nodes in each layer, inductively. Thus, the proof follows from applying the chain rule to $\sigma(\cdot)$ and its derivative. To elaborate, it is trivial to show that

$$\begin{aligned}\frac{\partial \sigma_{h,1}(\mathbf{x}, \mathbf{z})}{\partial z_j} &= w_{h1}^j \sigma'_{h,1}(\mathbf{x}, \mathbf{z}), \\ \frac{\partial^2 \sigma_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} &= w_{h1}^j w_{h1}^k \sigma''_{h,l}(\mathbf{x}, \mathbf{z}).\end{aligned}$$

Now, assume $l \geq 2$, from the chain rule we have

$$\frac{\partial \sigma_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_j} = \frac{\partial \sigma_{h,l}(\mathbf{x}, \mathbf{z})}{\partial (\sum_{h'=1}^H w_{h',l}^{h'} \sigma_{h',l-1}(\mathbf{x}, \mathbf{z}) + b_{h,l})} \frac{\partial (\sum_{h'=1}^H w_{h',l}^{h'} \sigma_{h',l-1}(\mathbf{x}, \mathbf{z}) + b_{h,l})}{\partial z_j} = \sigma'_{h,l}(\mathbf{x}, \mathbf{z}) \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j} \right),$$

while

$$\frac{\partial^2 \sigma_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} = \frac{\partial \sigma'_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_k} \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j} \right) + \frac{\partial (\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j})}{\partial z_k} \sigma'_{h,l}(\mathbf{x}, \mathbf{z})$$

by applying the chain rule to the term $\frac{\partial \sigma'_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_k}$, we get

$$\begin{aligned}\frac{\partial^2 \sigma_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} &= \frac{\partial \sigma'_{h,l}(\mathbf{x}, \mathbf{z})}{\partial (\sum_{h'=1}^H w_{h',l}^{h'} \sigma_{h',l-1}(\mathbf{x}, \mathbf{z}) + b_{h,l})} \frac{\partial (\sum_{h'=1}^H w_{h',l}^{h'} \sigma_{h',l-1}(\mathbf{x}, \mathbf{z}) + b_{h,l})}{\partial z_k} \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j} \right) \\ &\quad + \sigma'_{h,l}(\mathbf{x}, \mathbf{z}) \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial^2 \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} \right),\end{aligned}$$

which leads to

$$\frac{\partial^2 \sigma_{h,l}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} = \sigma''_{h,l}(\mathbf{x}, \mathbf{z}) \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_k} \right) \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j} \right) + \sigma'_{h,l}(\mathbf{x}, \mathbf{z}) \left(\sum_{h'=1}^H w_{h',l}^{h'} \frac{\partial^2 \sigma_{h',l-1}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} \right).$$

The last part of the Proposition follows trivially from noticing that $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}) = \sum_{h=1}^H w_{L+1}^{h'} \sigma_{h,L}(\mathbf{x}, \mathbf{z})$, and hence:

$$\frac{\partial^2 f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k} = \sum_{h=1}^H w_{L+1}^h \frac{\partial^2 \sigma_{h,L}(\mathbf{x}, \mathbf{z})}{\partial z_j \partial z_k}.$$

□

Proof of Lemma 1. From Algorithm 1 we have:

$$\hat{\mathbf{z}}(\mathbf{X}_0) \in \arg \min_{\mathbf{z} \in \mathcal{P}} \|\mathbf{z} - \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)\|.$$

Thus we must have:

$$\|\hat{\mathbf{z}}(\mathbf{X}_0) - \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)\| \leq \max_{\mathbf{z} \in \bar{\mathcal{P}}} \min_{\mathbf{z}' \in \mathcal{P}} \|\mathbf{z} - \mathbf{z}'\|.$$

Moreover, $f_{\hat{\theta}}$ is Lipschitz continuous, as it is a composition of Lipschitz continuous functions. Hence, if $K_{\hat{\theta}}$ is its Lipschitz constant, then we have:

$$|f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}(\mathbf{X}_0)) - f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0))| \leq K_{\hat{\theta}} (\max_{\mathbf{z} \in \bar{\mathcal{P}}} \min_{\mathbf{z}' \in \mathcal{P}} \|\mathbf{z} - \mathbf{z}'\|).$$

Now, if $f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) \leq \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})$ then by definition we have

$$f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) \leq \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z}) \leq f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}(\mathbf{X}_0))$$

which implies

$$|f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}(\mathbf{X}_0)) - \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})| \leq K_{\hat{\theta}}(\max_{\mathbf{z} \in \bar{\mathcal{P}}} \min_{\mathbf{z}' \in \mathcal{P}} \|\mathbf{z} - \mathbf{z}'\|).$$

Hence, assume $f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) \geq \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})$. We have

$$\begin{aligned} & |f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}(\mathbf{X}_0)) - \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})| = \\ & |f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}(\mathbf{X}_0)) - f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) + f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) - \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})| \leq \\ & |f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}(\mathbf{X}_0)) - f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0))| + |f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) - \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})| \leq \\ & K_{\hat{\theta}}(\max_{\mathbf{z} \in \bar{\mathcal{P}}} \min_{\mathbf{z}' \in \mathcal{P}} \|\mathbf{z} - \mathbf{z}'\|) + |f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) - \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})| \leq \\ & K_{\hat{\theta}}(\max_{\mathbf{z} \in \bar{\mathcal{P}}} \min_{\mathbf{z}' \in \mathcal{P}} \|\mathbf{z} - \mathbf{z}'\|) + |f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)) - \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z})| \end{aligned}$$

where the last inequality follows from the fact that by definition $\bar{\mathcal{P}}$ is a superset of \mathcal{P} and hence

$$\min_{\mathbf{z} \in \bar{\mathcal{P}}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z}) \leq \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{X}_0, \mathbf{z}) \leq f_{\hat{\theta}}(\mathbf{X}_0, \hat{\mathbf{z}}_{\bar{\mathcal{P}}}(\mathbf{X}_0)),$$

which completes the proof. \square

Proof of Proposition 3. We first express $f_*(\mathbf{x}, z)$ in a closed form and then its weak derivatives with respect to its input. Note that the following expressions are all in the almost-everywhere sense, which is how weak derivatives are defined. Because

$$f_*(\mathbf{x}, z) = \int_{-\infty}^{z-g(\mathbf{x})} h(z-g(\mathbf{x})-u)f_{\mu}(u)du + \int_{z-g(\mathbf{x})}^{\infty} b(g(\mathbf{x})+u-z)f_{\mu}(u)du,$$

we have

$$\frac{\partial f_*(\mathbf{x}, z)}{\partial z} = \int_{-\infty}^{z-g(\mathbf{x})} h f_{\mu}(u)du - \int_{z-g(\mathbf{x})}^{\infty} b f_{\mu}(u)du \leq \max\{h, b\}M_{\epsilon}\|\bar{\epsilon} - \underline{\epsilon}\|$$

This suggests that

$$\frac{\partial^2 f_*}{\partial z^2} = (h-b)f_{\mu}(z-g(\mathbf{x}))$$

and according to the assumptions on g and f_{μ} , we have

$$\frac{\partial^2 f_*}{\partial z^2} \in \mathcal{W}^{\min\{m,p\},\infty}([-1, 1]^{d_X+1}, (h-b)M_{\epsilon}).$$

Similarly, we have

$$\frac{\partial^2 f_*}{\partial z \partial x_i} = (b-h)\frac{\partial g}{\partial x_i}f_{\mu}(z-g(\mathbf{x})).$$

By the assumption on g and f_μ ,

$$\frac{\partial^2 f_*}{\partial z \partial x_i} \in \mathcal{W}^{\min\{m, p-1\}, \infty}([-1, 1]^{d_X+1}, (b-h)M_\epsilon M_g^2).$$

Furthermore, we have

$$\frac{\partial f_*}{\partial x_i} = \int_{z-g(\mathbf{x})}^{\infty} b \frac{\partial g}{\partial x_i} f_\mu(u) du - \int_{-\infty}^{z-g(\mathbf{x})} h \frac{\partial g}{\partial x_i} f_\mu(u) du \leq \max\{h, b\} M_\epsilon M_g \|\bar{\epsilon} - \underline{\epsilon}\|.$$

Finally,

$$\frac{\partial^2 f_*}{\partial x_i \partial x_j} = \frac{\partial g}{\partial x_j} (h-b) \left(\frac{\partial g}{\partial x_i} \right) f_\mu(z-g(\mathbf{x})) + \int_{z-g(\mathbf{x})}^{\infty} b \frac{\partial^2 g}{\partial x_i \partial x_j} f_\mu(u) du - \int_{-\infty}^{z-g(\mathbf{x})} h \frac{\partial^2 g}{\partial x_i \partial x_j} f_\mu(u) du.$$

By the assumption that $f_\mu(u) = 0, \forall u \notin [\underline{\epsilon}, \bar{\epsilon}]$ and given g is bounded and in $\mathcal{W}^{p, \infty}([-1, 1]^{d_X}, M_g)$, there must exist $\hat{M} < \infty$ such that

$$\frac{\partial^2 f_*}{\partial x_i \partial x_j} \in \mathcal{W}^{\min\{m, (p-2)^+\}, \infty}([-1, 1]^{d_X}, \hat{M}).$$

Given the fact that $\underline{\epsilon} \leq \epsilon \leq \bar{\epsilon}$, while g is bounded and in $\mathcal{W}^{p, \infty}([-1, 1]^{d_X}, M_g)$, we must have $f_*(\mathbf{x}, z) \leq \bar{M}$ for some $\bar{M} < \infty$. Moreover, remember we assumed the marginal probability measure of (\mathbf{X}, Z) is absolutely continuous with respect to the Lebesgue measure, while $[-1, 1]^{d_X+1}$, which is a compact set, is the support of (\mathbf{x}, z) . As such, for $\beta = \min\{m+1, p\}$, there exists $M < \infty$ such that

$$f_*(\mathbf{x}, z) \in \mathcal{W}^{\beta, \infty}([-1, 1]^{d_X+1}, M).$$

□

Proof of Proposition 4. By Assumption 2, we have that $\mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2] \leq \delta_N$. Recall that $\underline{u} > 0$ is the lower bound for the conditional PMF $\mu_{\cdot|\mathbf{X}}$ for all \mathbf{x} and \mathbf{z} . It is easy to show that the marginal PMF $\mu_Z \geq \underline{u}$ as well:

$$\mu_Z(\mathbf{z}) = \int_{\mathbf{x} \in [-1, 1]^{d_X}} \mu_{\cdot|\mathbf{X}}(\mathbf{z}) d\mu_{\mathbf{X}} \geq \int_{\mathbf{x} \in [-1, 1]^{d_X}} \underline{u} d\mu_{\mathbf{X}} = \underline{u}.$$

Then we have for all $\mathbf{z}' \in \mathcal{P}$:

$$\mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2 | \mathbf{Z} = \mathbf{z}'] \times \underline{u} \leq \mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2 | \mathbf{Z} = \mathbf{z}'] \times \Pr(\mathbf{Z} = \mathbf{z}').$$

Moreover, by the law of total probability, we have

$$\begin{aligned} \mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2 | \mathbf{Z} = \mathbf{z}'] \times \Pr(\mathbf{Z} = \mathbf{z}') &\leq \sum_{\mathbf{z} \in \mathcal{P}} \mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2 | \mathbf{Z} = \mathbf{z}] \times \Pr(\mathbf{Z} = \mathbf{z}) \\ &= \mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2]. \end{aligned}$$

As a result, we have

$$\mathbb{E}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2 | \mathbf{Z} = \mathbf{z}'] \leq \frac{\delta_N}{\underline{u}}.$$

For any $(\mathbf{x}_0, \mathbf{z}')$, where \mathbf{x}_0 is not an extreme point, consider a ball $B \subset \mathbb{R}^{d_X}$ of radius b that contains \mathbf{x}_0 , and the ball itself is fully contained in $[-1, 1]^{d_X}$ (later we will show that b will shrink in δ_N and hence if δ_N is small enough, this is always possible). Let $\mathbf{x}^{min} \in \arg \min_{\mathbf{x} \in B} |f_*(\mathbf{x}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}')|$ and $\mathbf{x}^{max} \in \arg \max_{\mathbf{x} \in B} |f_*(\mathbf{x}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}')|$. Because $d\mu_{\mathbf{X}} \geq \underline{u}dx$ and $\mu_{\mathbf{z}} \leq 1$:

$$d\mu_{\mathbf{X}|\mathbf{z}} = \frac{d\mu_{\mathbf{X}} \times \mu_{\cdot|\mathbf{X}}(\mathbf{z})}{\mu_{\mathbf{z}}} \geq \underline{u}^2 dx$$

while we can recall that a ball of radius b in the d_X -dimensional space has volume $L_{d_X} b^{d_X}$ (see Definition 1), we have

$$\begin{aligned} \underline{u}^2 L_{d_X} b^{d_X} (f_*(\mathbf{x}^{min}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}'))^2 &\leq \int_{\mathbf{x} \in B} (f_*(\mathbf{x}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}'))^2 d\mu_{\mathbf{X}|\mathbf{z}'} \\ &\leq \int_{\mathbf{x} \in B} (f_*(\mathbf{x}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}'))^2 d\mu_{\mathbf{X}|\mathbf{z}'} \leq \frac{\delta_N}{\underline{u}} \end{aligned}$$

Hence, we can upper bound $|f_*(\mathbf{x}^{min}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}')|$ as:

$$|f_*(\mathbf{x}^{min}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}')| \leq \sqrt{\frac{\delta_N}{\underline{u}^3 L_{d_X} b^{d_X}}}.$$

By Assumption 3, the function $|f_*(\mathbf{x}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z}')|$ is Lipschitz continuous with a constant K . Moreover, because $\|\mathbf{x}^{max} - \mathbf{x}^{min}\|_2 \leq 2b$, we have

$$|f_*(\mathbf{x}_0, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}')| \leq |f_*(\mathbf{x}^{max}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}')| \leq \sqrt{\frac{\delta_N}{\underline{u}^3 L_{d_X} b^{d_X}}} + 2bK := \epsilon. \quad (\text{EC.4})$$

We note that this inequality holds for all $\mathbf{z}' \in \mathcal{P}$.

Now consider $\hat{\mathbf{z}}(\mathbf{x}_0) \in \arg \min_{\mathbf{z} \in \mathcal{P}} f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z})$, the optimal DNN solution. Furthermore, with a slight abuse of notation let $\mathbf{z}^* \in \arg \max_{\mathbf{z} \in \mathcal{P}} f_*(\mathbf{x}_0, \mathbf{z})$, an optimal solution of f_* in \mathcal{P} .

Applying the inequality (EC.4) first for $\mathbf{z}' = \hat{\mathbf{z}}(\mathbf{x}_0)$ and then for $\mathbf{z}' = \mathbf{z}^*(\mathbf{x}_0)$, we have

$$f_*(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0)) \leq f_{\hat{\theta}}(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0)) + \epsilon \leq f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + \epsilon \leq f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + 2\epsilon.$$

To complete the proof, we note that by choosing a suitable b for diminishing δ_N , we can shrink ϵ .

In particular, let

$$b = \left(\frac{(d_X)^2 \delta_N}{16K^2 \underline{u}^3 L_{d_X}} \right)^{\frac{1}{d_X+2}}.$$

At this point note that unless \mathbf{x}_0 is an extreme point, if δ_N is small enough, then b will also be small for the ball B to be fully contained in $[-1, 1]^{d_X}$. By plugging the obtained b into

$$\epsilon = \sqrt{\frac{\delta_N}{\underline{u}^3 L_{d_X} b^{d_X}}} + 2bK$$

we obtain the bound

$$2\left(\frac{\delta_N d_X}{4K\underline{u}^3 L_{d_X}} + 2\delta_N^{\frac{1}{d_X+2}} \left(\frac{K^{d_X} d_X^2}{16\underline{u}^3 L_{d_X}}\right)^{\frac{1}{d_X+2}}\right).$$

If \mathbf{x}_0 is an extreme point, then let the ball B be centered around \mathbf{x}_0 . We have

$$\frac{1}{2^{d_X}} \underline{u}^2 L_{d_X} b^{d_X} (f_*(\mathbf{x}^{min}, \mathbf{z}') - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}'))^2 \leq \frac{\delta_N}{\underline{u}}$$

giving raise to the following alternative bound:

$$2\left(\frac{\delta_N d_X}{4K\underline{u}^3 2^{d_X} L_{d_X}} + 2\delta_N^{\frac{1}{d_X+2}} \left(\frac{K^{d_X} d_X^2}{16\underline{u}^3 2^{d_X} L_{d_X}}\right)^{\frac{1}{d_X+2}}\right),$$

which completes the proof. \square

Proof of Lemma 2 By Assumption 2, we have that $\mathbb{E}\|f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z})\|^2 \leq \delta_N$. Now let $(\mathbf{x}^{min}, \mathbf{z}^{min}) \in \arg \min_{(\mathbf{x}, \mathbf{z}) \in B \cap [-1, 1]^{d_X} \times \mathcal{P}} \|f_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})\|$ and $(\mathbf{x}^{max}, \mathbf{z}^{max}) \in \arg \max_{(\mathbf{x}, \mathbf{z}) \in B \cap [-1, 1]^{d_X} \times \mathcal{P}} \|f_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})\|$.

We have:

$$\int_{(\mathbf{x}, \mathbf{z}) \in B \cap [-1, 1]^{d_X} \times \mathcal{P}} \|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})\|^2 d\mu_{X, Z} \leq \int_{(\mathbf{x}, \mathbf{z}) \in B \cap [-1, 1]^{d_X} \times \mathcal{P}} \|f_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})\|^2 d\mu_{X, Z},$$

while we know :

$$\int_{(\mathbf{x}, \mathbf{z}) \in B \cap [-1, 1]^{d_X} \times \mathcal{P}} \|f_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})\|^2 d\mu_{X, Z} \leq \int_{(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P}} \|f_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})\|^2 d\mu_{X, Z} \leq \delta_N.$$

$$\implies S \|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})\|^2 \leq \delta_N$$

Thus we can see that:

$$|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})| \leq \sqrt{\frac{\delta_N}{S}}. \quad (\text{EC.5})$$

At this point, we need to relate $|f_*(\mathbf{x}^{max}, \mathbf{z}^{max}) - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}^{max})|$ to $|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})|$. We have

$$\begin{aligned} & |f_*(\mathbf{x}^{max}, \mathbf{z}^{max}) - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}^{max})| - |f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})| \leq \\ & \|f_*(\mathbf{x}^{max}, \mathbf{z}^{max}) - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}^{max})\| - \|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})\|. \end{aligned}$$

Moreover, since we know $|f_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})|$ is K -Lipschitz continuous, and by definition $\|(\mathbf{x}^{max}, \mathbf{z}^{max}) - (\mathbf{x}^{min}, \mathbf{z}^{min})\| \leq 2b$ (both points lie in a ball with radius b), we have

$$\|f_*(\mathbf{x}^{max}, \mathbf{z}^{max}) - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}^{max})\| - \|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})\| \leq 2bK.$$

Therefore, it follows that

$$\|f_*(\mathbf{x}^{max}, \mathbf{z}^{max}) - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}^{max})\| - \|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})\| \leq 2bK.$$

Hence, we have

$$\|f_*(\mathbf{x}^{max}, \mathbf{z}^{max}) - f_{\hat{\theta}}(\mathbf{x}^{max}, \mathbf{z}^{max})\| \leq \|f_*(\mathbf{x}^{min}, \mathbf{z}^{min}) - f_{\hat{\theta}}(\mathbf{x}^{min}, \mathbf{z}^{min})\| + 2bK \leq \sqrt{\frac{\delta_N}{S}} + 2bK,$$

where the last inequality follows from (EC.5). \square

Proof of Lemma 3. We proceed with this proof by constructing an alternative target function \bar{f}_* that agrees with f_* everywhere except within an area of zero measure. To that avail let Assumptions 1 (continuous version) to 4 hold. Consider \mathbf{x}_0 for which condition 1 in the lemma holds. In other words, for all $(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P}$ and $\|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\| \leq r_1$, the PDF of $\mu_{X,Z}$ is 0. Consider an alternative \bar{f}_* that has the exact same value as f_* everywhere, except for $(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P}$ where $\|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\| \leq r_1$. For M defined in Assumption 1, let

$$\bar{f}_*(\mathbf{x}, \mathbf{z}) = \begin{cases} \frac{(4+\epsilon_M)(r_1 - \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\|)}{r_1} M + f_*(\mathbf{x}, \mathbf{z}), & \text{if } (\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P} \text{ \& } \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\| \leq r_1 \\ f_*(\mathbf{x}, \mathbf{z}), & \text{otherwise.} \end{cases}$$

\bar{f}_* is continuous and hence, by the universal approximation result for neural networks (e.g. [Leshno et al. \(1993\)](#)), for any $\epsilon > 0$ there exists $f_{\hat{\theta}}$ that for all $(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P}$, $|\bar{f}_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})| \leq \epsilon$. Meanwhile, it is trivial to check that for any δ_N , if ϵ small enough, we have

$$\mathbb{E}_{\mu_{X,Z}}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2] \leq \delta_N.$$

Let $\hat{\mathbf{z}}(\mathbf{X}_0)$ be the optimizer of $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$. Then, it is easy to see that given

$$\bar{f}_*(\mathbf{x}, \mathbf{z}) \geq M + M\epsilon_M/2, \quad \forall (\mathbf{x}, \mathbf{z}) : \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\| \leq r_1/2,$$

it must be that for some $\epsilon < M\epsilon_M/4$,

$$\|\hat{\mathbf{z}}(\mathbf{x}_0) - \mathbf{z}^*(\mathbf{x}_0)\| \geq r_1/2,$$

which completes the proof for condition 1.

The proof for the case under condition 2 is similar. Assume for some $\mathbf{z}' \in \mathcal{P}$, there exists $r_2 > 0$ such that $\|(\mathbf{x}_0, \mathbf{z}') - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\| \geq r_1/2 + r_2$, and we have $\int_{(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P} : \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}')\| \leq r_2} d\mu_{X,Z} = 0$.

We construct \bar{f}_* such that it agrees with f_* everywhere except for $(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P} : \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}')\| \leq r_2$ as follows:

$$\bar{f}_*(\mathbf{x}, \mathbf{z}) = \begin{cases} -\frac{(4+\epsilon_M)(r_2 - \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}')\|)}{r_2} M + f_*(\mathbf{x}, \mathbf{z}), & \text{if } (\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P} \text{ \& } \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}')\| \leq r_2 \\ f_*(\mathbf{x}, \mathbf{z}), & \text{otherwise.} \end{cases}$$

\bar{f}_* is continuous and hence, by the universal approximation result for neural networks, for any $\epsilon > 0$ there exists $f_{\hat{\theta}}$ that for all $(\mathbf{x}, \mathbf{z}) \in [-1, 1]^{d_X} \times \mathcal{P}$, $|\bar{f}_*(\mathbf{x}, \mathbf{z}) - f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})| \leq \epsilon$. Meanwhile, it is trivial to check that for any δ_N , if ϵ small enough, we have

$$\mathbb{E}_{\mu_{X,Z}}[(f_*(\mathbf{X}, \mathbf{Z}) - f_{\hat{\theta}}(\mathbf{X}, \mathbf{Z}))^2] \leq \delta_N.$$

Let $\hat{\mathbf{z}}(\mathbf{X}_0)$ be the optimizer of $f_{\hat{\theta}}(\mathbf{x}, \mathbf{z})$. Then, it is easy to see that given

$$\bar{f}_*(\mathbf{x}, \mathbf{z}) \leq -M - M\epsilon_M/2, \quad \forall (\mathbf{x}, \mathbf{z}) : \|(\mathbf{x}, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}')\| \leq r_2/2,$$

it must be that for some $\epsilon < M\epsilon_M/4$,

$$\|\hat{\mathbf{z}}(\mathbf{x}_0) - \mathbf{z}'\| \leq r_2,$$

which suggests we have

$$\|\hat{\mathbf{z}}(\mathbf{x}_0) - \mathbf{z}^*(\mathbf{x}_0)\| \geq r_1/2.$$

This completes the proof for condition 2. □

Proof of Proposition 5. We first will show that $\|\hat{\mathbf{z}}(\mathbf{x}_0) - \mathbf{z}^*(\mathbf{x}_0)\| \leq r$ where r is defined in Assumption 5. By Assumption 6 we have $c_2 - 2c_3K > 0$. Thus $c_2 - 2(c_3 + r')K > 0$, for some $0 < r' < r$. Let δ_N be small enough such that $c_2 - 2(c_3 + r')K > \sqrt{\frac{\delta_N}{c_4}} + \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} c_1 r'^{d_X+d_Z}}}$. Define $S_{r'} = \int_{(\mathbf{x}, \mathbf{z}) \in B_{r'}} d\mu_{X,Z}$, where $B_{r'}$ is a ball of radius r' centered at $(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))$. Then we have

$$c_2 - 2(c_3 + r')K > \sqrt{\frac{\delta_N}{c_4}} + \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} c_1 r'^{d_X+d_Z}}} \geq \sqrt{\frac{\delta_N}{c_4}} + \sqrt{\frac{\delta_N}{S_{r'}}}. \quad (\text{EC.6})$$

The last inequality follows from Assumption 5 and the fact that

$$\int_{(\mathbf{x}, \mathbf{z}) \in B_{r'}} d\mu_{X,Z} \geq c_1 L_{d_X+d_Z} r'^{d_X+d_Z}.$$

Now, consider some $\mathbf{z} \in \mathcal{P}$ such that $\|\mathbf{z} - \mathbf{z}^*(\mathbf{x}_0)\| > r$. Then $\|(\mathbf{x}_0, \mathbf{z}) - (\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))\| > r$. It directly follows from the definition of $\hat{\mathbf{z}}(\mathbf{x}_0)$ and Lemma 2 by considering a ball of radius r' surrounding $(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))$ that:

$$f_{\hat{\theta}}(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0)) \leq f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) \leq f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + \sqrt{\frac{\delta_N}{S_{r'}}} + 2r'K. \quad (\text{EC.7})$$

Moreover, from Assumption 6, we have that $f_*(\mathbf{x}_0, \mathbf{z}) - f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) \geq c_2$. Thus

$$f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + \sqrt{\frac{\delta_N}{S_{r'}}} + 2r'K \leq f_*(\mathbf{x}_0, \mathbf{z}) - c_2 + \sqrt{\frac{\delta_N}{S_{r'}}} + 2r'K. \quad (\text{EC.8})$$

Furthermore, given (EC.6), from Inequalities (EC.7) and (EC.8) we have:

$$f_{\hat{\theta}}(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0)) \leq f_*(\mathbf{x}_0, \mathbf{z}) - c_2 + \sqrt{\frac{\delta_N}{S_{r'}}} + 2r'K < f_*(\mathbf{x}_0, \mathbf{z}) - \sqrt{\frac{\delta_N}{c_4}} - 2c_3K. \quad (\text{EC.9})$$

Now, we need to bound the error of $f_{\hat{\theta}}$ at $(\mathbf{x}_0, \mathbf{z})$. By Assumption 6 we know that a $d_X + d_Z$ dimensional ball of radius c_3 around $(\mathbf{x}_0, \mathbf{z})$ has a probability measure of at least c_4 . Hence we can use Lemma 2, suggesting:

$$|f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}) - f_*(\mathbf{x}_0, \mathbf{z})| \leq \sqrt{\frac{\delta_N}{c_4}} + 2c_3K.$$

This implies:

$$f_*(\mathbf{x}_0, \mathbf{z}) - \sqrt{\frac{\delta_N}{c_4}} - 2c_3K \leq f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z})$$

From Inequality (EC.9) we finally have:

$$f_{\hat{\theta}}(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0)) \leq f_*(\mathbf{x}_0, \mathbf{z}) - c_2 + \sqrt{\frac{\delta_N}{S_{r'}}} + 2r'K < f_*(\mathbf{x}_0, \mathbf{z}) - \sqrt{\frac{\delta_N}{c_4}} - 2c_3K \leq f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}).$$

Thus, we have proven that for any \mathbf{z} such that $\|\mathbf{z} - \mathbf{z}^*(\mathbf{x}_0)\| > r$, we have:

$$f_{\hat{\theta}}(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0)) < f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}),$$

which proves:

$$\|\hat{\mathbf{z}}(\mathbf{x}_0) - \mathbf{z}^*(\mathbf{x}_0)\| \leq r.$$

In other words, the optimal decision output by the fitted NN is within a neighborhood of the actual optimal decision. Next we will find the smallest such neighborhood and show that they are getting smaller as δ_N decreases.

For some $r_1 < r$, we want to show that for any \mathbf{z} such that $\|\mathbf{z} - \mathbf{z}^*(\mathbf{x}_0)\| > r_1$, we have

$$f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) < f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}).$$

Consider $\mathbf{z}' \in \mathcal{P}$ such that $r_1 < \|\mathbf{z}' - \mathbf{z}^*(\mathbf{x}_0)\| \leq r$. We know from Assumption 7 that:

$$f_*(\mathbf{x}_0, \mathbf{z}') > f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + c_5r_1^2.$$

At this point, we would like to lower bound $f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}')$ using $f_*(\mathbf{x}_0, \mathbf{z}')$. To that avail, consider a $(d_X + d_Z)$ dimensional ball B_{r_2} of radius $r_2 \leq r_1/2$ centered at $(\mathbf{x}_0, \mathbf{z}'')$. The value \mathbf{z}'' is chosen such that \mathbf{z}'' is a convex combination of \mathbf{z}' and $\mathbf{z}^*(\mathbf{x}_0)$, and such that $\|\mathbf{z}' - \mathbf{z}''\| = r_2$. Because \mathcal{P} is

convex according to Assumption 1, we can always find such a ball with $(\mathbf{x}_0, \mathbf{z}'') \in \mathcal{P}$. Note that by construction, $(\mathbf{x}_0, \mathbf{z}') \in B_{r_2}$. From Lemma 2 we have:

$$f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}') \geq f_*(\mathbf{x}_0, \mathbf{z}') - \sqrt{\frac{\delta_N}{S}} - 2r_2K \geq f_*(\mathbf{x}_0, \mathbf{z}') - \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} - 2r_2K, \quad (\text{EC.10})$$

where $S = \int_{(\mathbf{x}, \mathbf{z}) \in B_{r_2} \cap [-1, 1]^{d_X} \times \mathcal{P}} d\mu_{X,Z}$. The last inequality follows from Assumption 5 as we have

$$\int_{(\mathbf{x}, \mathbf{z}) \in B_{r_2} \cap [-1, 1]^{d_X} \times \mathcal{P}} d\mu_{X,Z} \geq c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}.$$

Next, we would like also to develop an upper bound on $f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))$. Thus, let us consider another Euclidean ball around $(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0))$ with radius r_2 . Then, by Lemma 2 and Assumption 5, we know that

$$f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) \leq f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} + 2r_2K. \quad (\text{EC.11})$$

Recall that by construction we have $\|\mathbf{z}' - \mathbf{z}^*(\mathbf{x}_0)\| > r_1$. From Assumption 7, (EC.11) implies that:

$$f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) \leq f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) + \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} + 2r_2K < f_*(\mathbf{x}_0, \mathbf{z}') - c_5 r_1^2 + \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} + 2r_2K. \quad (\text{EC.12})$$

Therefore, if the following inequality holds

$$\sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} + 2r_2K \leq \frac{c_5 r_1^2}{2}, \quad (\text{EC.13})$$

then for the right-hand side of Inequality (EC.11) and from Inequality (EC.10) we will have:

$$f_*(\mathbf{x}_0, \mathbf{z}') - c_5 r_1^2 + \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} + 2r_2K \leq f_*(\mathbf{x}_0, \mathbf{z}') - \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} - 2r_2K \leq f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}').$$

This, along with Inequality (EC.12), proves that

$$f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) < f_{\hat{\theta}}(\mathbf{x}_0, \mathbf{z}').$$

To find the smallest r_1 for which Inequality (EC.13) holds, we note that r_1 will achieve its lowest value (that satisfies Inequality (EC.13)) when r_2 is set to be the minimizer of the left-hand side of Inequality (EC.13), which happens to be convex with respect to r_2 . To prove the convexity, we notice that:

$$\frac{d\left(\sqrt{\frac{\delta_N}{L_{d_X+d_Z} c_1 r_2^{d_X+d_Z}}} + 2r_2K\right)}{dr_2} = -\frac{d_X+d_Z}{2} \sqrt{\frac{\delta_N}{L_{d_X+d_Z} c_1}} r_2^{-\frac{d_X+d_Z+2}{2}} + 2K,$$

while:

$$\frac{-\frac{d_X+d_Z}{2} \sqrt{\frac{\delta_N}{L_{d_X+d_Z} c_1}} r_2^{-\frac{d_X+d_Z+2}{2}} + 2K}{dr_2} = \frac{(d_X+d_Z)(d_X+d_Z+2)}{4} \sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z}}} r_2^{-\frac{d_X+d_Z+4}{2}} \geq 0.$$

Hence, the best r_2 would satisfy the first-order condition of the left hand side of Inequality (EC.13) with respect to r_2 . Therefore, we have:

$$r_2 = \left(\frac{(d_X+d_Z)^2 \delta_N}{16K^2 c_1 L_{d_X+d_Z}} \right)^{\frac{1}{d_X+d_Z+2}}.$$

Then by plugging r_2 into the left hand side of Inequality (EC.13), we will have:

$$\sqrt{\frac{\delta_N}{c_1 L_{d_X+d_Z} r_2^{d_X+d_Z}}} + 2r_2 K = \frac{\delta_N (d_X+d_Z)}{4K c_1 L_{d_X+d_Z}} + 2K \left(\frac{(d_X+d_Z)^2 \delta_N}{16K^2 c_1 L_{d_X+d_Z}} \right)^{\frac{1}{d_X+d_Z+2}}.$$

Hence, setting

$$r_1 = \sqrt{2 \frac{\frac{\delta_N (d_X+d_Z)}{4K c_1 L_{d_X+d_Z}} + 2K \left(\frac{(d_X+d_Z)^2 \delta_N}{16K^2 c_1 L_{d_X+d_Z}} \right)^{\frac{1}{d_X+d_Z+2}}}{c_5}},$$

completes the proof. We note that as long as δ_N is small enough we have:

$$r_2 = \left(\frac{(d_X+d_Z)^2 \delta_N}{16K^2 c_1 L_{d_X+d_Z}} \right)^{\frac{1}{d_X+d_Z+2}} \leq \sqrt{\frac{\frac{\delta_N (d_X+d_Z)}{4K c_1 L_{d_X+d_Z}} + 2K \left(\frac{(d_X+d_Z)^2 \delta_N}{16K^2 c_1 L_{d_X+d_Z}} \right)^{\frac{1}{d_X+d_Z+2}}}{2c_5}} = \frac{r_1}{2} < \frac{r}{2},$$

where the inequality holds since its left hand side diminishes much quicker in δ_N .

To finalize the proof, we note that since $\|\hat{\mathbf{z}}(\mathbf{x}_0) - \mathbf{z}^*(\mathbf{x}_0)\| \leq r_1$, it follows directly from Assumption 3 that:

$$|f_*(\mathbf{x}_0, \mathbf{z}^*(\mathbf{x}_0)) - f_*(\mathbf{x}_0, \hat{\mathbf{z}}(\mathbf{x}_0))| \leq K_* \sqrt{2 \frac{\frac{\delta_N (d_X+d_Z)}{4K c_1 L_{d_X+d_Z}} + 2K \left(\frac{(d_X+d_Z)^2 \delta_N}{16K^2 c_1 L_{d_X+d_Z}} \right)^{\frac{1}{d_X+d_Z+2}}}{c_5}}.$$

□

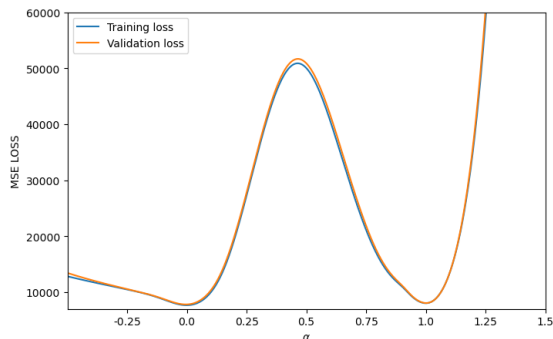


Figure EC.3 MSE loss of the linear interpolation of small and large batch DNNs.

EC.5. Additional Investigations and Figures

In some applications of neural networks such as image classification, researchers have pointed out that neural networks trained with SGD with large batch sizes may not generalize well (Keskar et al. 2016, Li et al. 2018). The existing research suggests this could be because large batch methods tend to converge to sharp minimizers of the training and testing loss functions (Keskar et al. 2016). Given the importance of generalization bounds in our paper, we look into such potential sharpness in Equation (6). Similar to Li et al. (2018), to visualize such sharpness, we consider two batch sizes of 50 (small batch) and 1000 (large batch) in the newsvendor experiment. On a data set comprised of 20,000 training data points and 5,000 test data points, we look at the converged DNNs of the small batch SGD and that of the large batch one. Let $\hat{\theta}^S$ denote the parameters of the fitted DNN with small batch, and $\hat{\theta}^L$ that of the large batch. Following Li et al. (2018), let $f_{\hat{\theta}^\alpha}$ be the neural network with parameters $\alpha \times \hat{\theta}^L + (1 - \alpha) \times \hat{\theta}^S$. Figure EC.3 portrays the training and validation MSE loss of $f_{\hat{\theta}^\alpha}$ as a function of α varying for $-.5$ to 1.5 . The result suggests that in the application, such sharpness of loss function for neural networks trained with large batch size SGD may not be present.

Figure EC.4 depicts the training and validation error as a function of training epochs for one of the DNNs in Example 1 when training with 20,000 data points and validating with 5,000. It portrays a robust effect across the experiments. The training error does not shrink to zero as observed in other applications of DNNs (most probably because we use relatively small DNNs), and further, the validation error tends to be slightly larger than the training error. In other words, we do not observe benign overfitting documented in, e.g., Bartlett et al. (2020). This is because we focus on the classic statistical learning regime in which the data size is not huge and the number of parameters is significantly less than the number of data points. Moreover, the figure provides evidence that in the newsvendor experiment, the SGD method used to train the DNNs does not suffer from unstable convergence, as reported elsewhere in the literature (Ahn et al. 2022).

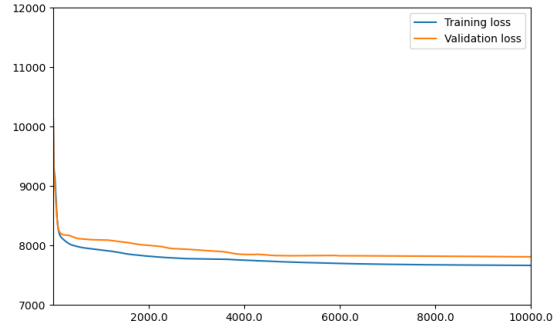


Figure EC.4 Training versus Validation Loss.

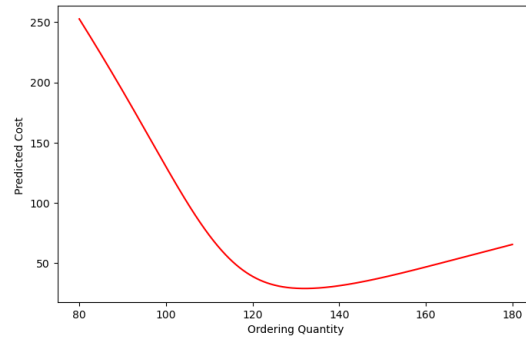


Figure EC.5 The illustration of the effect of order quantity on predicted cost.

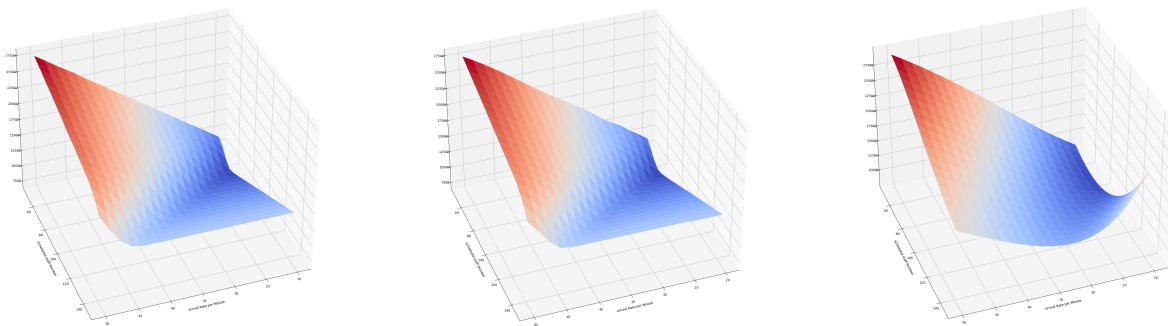


Figure EC.6 $G/GI/z + G$ actual and predicted costs. Left is the actual cost, middle is the DNN prediction and right is the PR prediction.

Figure EC.5 depicts the predicted newsvendor cost as a function of order quantity when all other covariates are fixed (after training with 250 days of data for one of the DNNs). The figure suggests a smooth and well-behaved objective function that strongly resembles the actual newsvendor objective function in Equation 13. Hence, if the true objective function does not suffer from local optima, we expect the fitted DNNs prediction value to exhibit a similar trend. In Example 1 we optimize the DNNs from multiple starting points and all of the results are nearly identical suggesting a lack of local optima issues.

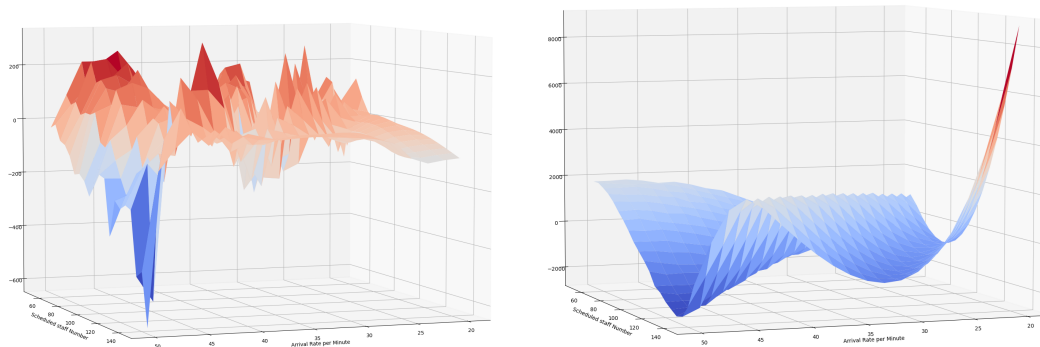


Figure EC.7 Prediction error of DNN and Polynomial Regression in the $G/GI/z+G$ case.

Figure EC.6 presents the actual cost (left figure) and predicted cost (DNN center figure, PR right figure) of the $G/GI/z+G$ queue for various arrival rates and staffing levels. Figure EC.7 compliments Figure EC.6. The left-hand side figure suggests that when predicting the daily cost of the call center staffing problem (Example 3) for the $G/GI/z+G$ case, the error by the DNN is quite small (generally less than 1%) and centered around zero throughout the tested ranges of staffing levels and arrival rates. The error is never larger than 5% of the actual objective function (portrayed in the left figure in Figure EC.6). The right figure in Figure EC.7 suggests the same cannot be said about Polynomial Regression method that we tested and it may suffer from very large errors in certain regions.

Figures EC.8a, EC.8b and EC.8c pertain to the predicted optimal cost using each methodology, in the linear staffing case, increasing convex staffing and $G/GI/z+G$ case, respectively. We emphasize that these results are not simulated and are the values obtained by plugging each method's optimal staffing level into its approximation of Equation (22).

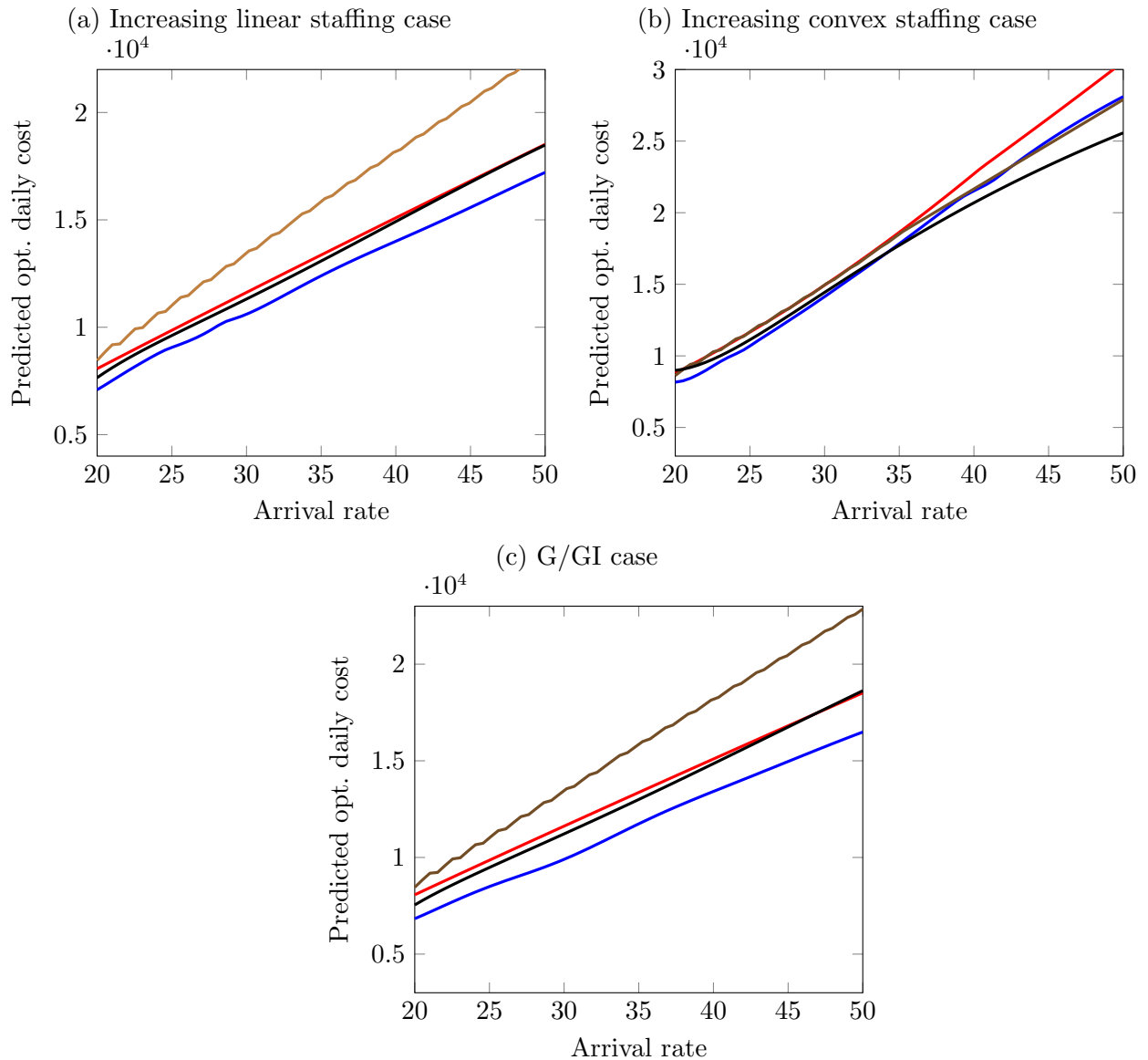


Figure EC.8 The illustration of the predicted costs of each methodology in the queuing setting.

EC.6. Example

In this section, we look at an example on how to apply our framework in details.² Consider an example with one scalar covariate x and one scalar decision variable z where $f_*(x, z) = (x - z)^2$.

We simulate a large data set with 50,000 data points where x and z are uniformly and independently drawn from $[-10, 10]$. We train a DNN with $L = 3$, $H = 3$ and Swish activation function. The fitted DNN has the following weight matrices and constant terms.

² See <https://github.com/saman-lagzi/Data-driven-Optimization-with-Neural-Networks> for a detailed Python implementation of Example EC.6 along with the generated dataset.

$$\mathbf{W}_1 = \begin{bmatrix} -1.7124 & 1.7100 \\ 0.0081 & -0.0080 \\ 1.7235 & -1.7209 \end{bmatrix}, \mathbf{W}_{1z} = \begin{bmatrix} 1.7100 \\ -0.0080 \\ -1.7209 \end{bmatrix}, \mathbf{b}_1 = \begin{bmatrix} -1.6824 \\ 1.3022 \\ -1.7184 \end{bmatrix},$$

$$\mathbf{W}_2 = \begin{bmatrix} 4.6518 & 1.9831 & 4.5889 \\ 0.2762 & -1.0031 & -0.2200 \\ -1.4228 & 3.8292 & -1.4879 \end{bmatrix}, \mathbf{b}_2 = \begin{bmatrix} 2.3721 \\ -2.2213 \\ 3.6026 \end{bmatrix},$$

$$\mathbf{W}_3 = \begin{bmatrix} 1.4862 & -1.1508 & -5.0347 \\ -0.1134 & -0.1684 & -0.1325 \\ 4.0344 & -0.6634 & -1.0397 \end{bmatrix}, \mathbf{b}_3 = \begin{bmatrix} 2.4532 \\ 0.1484 \\ 2.0276 \end{bmatrix},$$

$$\mathbf{W}_4 = \begin{bmatrix} 2.1921 \\ 0.8658 \\ 3.2429 \end{bmatrix}, b_4 = 0.9145.$$

First, we note that the derivative of the fitted DNN with respect to the decision variable z at any given point, calculated through Proposition 1 is identical to the value calculated using `torch.autograd` function in PyTorch. For example, at input point $[1, 1]$ both methods return a gradient of $\frac{\partial f_{\hat{\theta}}}{\partial z}(1, 1) = 0.183$. Moreover, the point $x = 0$ translates into 0.0014 when standardized for input to the fitted DNN, and using both Proposition 1 and the `torch.autograd` in PyTorch we have $\hat{z}(0.0014) = -0.0055$ where $\frac{\partial f_{\hat{\theta}}}{\partial z}(0.0014, -0.0055) = 0$.

To find $\hat{z}(0.0014) = \arg \min_z f_{\hat{\theta}}(0.0014, z)$, we use the `scipy.optimize` package from Python with a starting point of $z_0 = -1.96$. We provide the gradient $\frac{\partial f_{\hat{\theta}}(0.0014, z)}{\partial z}$, using both our approach highlighted in Proposition 1 and the `torch.autograd` function in PyTorch. While both methods lead to $\hat{z}(0.0014) = -0.0055$, the execution time is 0.10 second for the first method and 0.12 second for the latter. The faster execution time based on Proposition 1 is well within expectation as it provides the closed form formula for $\frac{\partial f_{\hat{\theta}}(0.0014, z)}{\partial z}$ based on the weight matrices and constant terms of $f_{\hat{\theta}}$, however, `torch.autograd` is a generalized method and needs to apply the chain rule,

$$\text{layer by layer to } f_{\hat{\theta}} \text{ to find } \frac{\partial f_{\hat{\theta}}(0.0014, z)}{\partial z}.$$