C\ **Taylor &Francis**
Taylor &Francis Group

# Asymptotic properties of neural network sieve estimators

Xiaoxi Shen a , Chang Jiang b , Lyudmila Sakhanenko c and Qing Lu (!) b

•Texas State University, San Marcos, TX, USA; b university of Florida, Gainesville, FL, USA c; Michigan State University, East Lansing, MI, USA

**ABSTRACT**

Neural networks have become one of the most popularly used methods in machine learning and artificial intelligence. Due to the universal approximation theorem, a neural network with one hidden layer can approximate any continuous function on compact support as long as the number of hidden units is sufficiently large. Statistically, a neural network can be classified into a nonlinear regression framework. However, if we consider it parametrically, due to the unidentifiability of the parameters, it is difficult to derive its asymptotic properties. Instead, we consider the estimation problem in a nonparametric regression framework and use the results from sieve estimation to establish the consistency, the rates of convergence and the asymptotic normality of the neural network estimators. We also illustrate the validity of the theories via simulations.

## 1. Introduction

With the success of machine learning and artificial intelligence in research and industry, neural networks have become popularly used methods nowadays. Many newly developed machine learning methods are based on deep neural networks and have achieved great classification and prediction accuracy. We refer interested readers to Goodfellow et al. (2016) for more background and details. In classical statistical learning theory, the consistency and the rate of convergence of the empirical risk minimisation principle are of great interest. Many upper bounds have been established for the empirical risk and the sample complexity based on the growth function and the Vapnik- Chervonenkis dimension (see, e.g., Vapnik 1998; Anthony and Bartlett 2009; Devroye et al. 2013). However, few studies have focused on the asymptotic properties for neural networks. As Thomas J. Sargent said, 'artificial intelligence is actually statistics, but in a very gorgeous phrase, it is statistics '. So it is natural and worthwhile to explore whether neural networks possess nice asymptotic properties. As if they do, it may be possible to conduct statistical inference based on neural networks. Throughout this paper, we will focus on the asymptotic properties of neural networks with one hidden layer.

One of the most important properties of neural networks is that they are universal approximants (Hornik et al. 1989), which means any continuous function on a compact support can be approximated arbitrarily well by a neural network with one hidden layer. So it seems natural to consider it as a nonparametric regression problem and approximate the underlying function class through a class of neural networks with one hidden layer. Many other series-based estimators such as splines and wavelets also possess similar approximation properties and have been extensively studied in the literature. Chen (2007) provides a comprehensive review for those methods. For nonparametric regression problems, random design and fixed design are the two main frameworks. Many existing literature on neural networks focus on random design (e.g. Chen and White 1999; Gyorfi et al. 2002). On the other hand, general theories on nonparametric regression under fixed design have been well studied in van de Geer (2000). Therefore, it is still worthwhile to study neural networks under fixed design.

Specifically, consider the following nonparametric regression mode l:

$$y_i = f_0(x_i) + \epsilon_i,$$

where $E_1, \ldots, E_n$ are i.i.d. random variables defined on a complete probability space $(Q, A, IJD)$ with $\mathbb{E}[E] = 0$, $Var[E] = a^2$ and $\mathbb{E}[|E|^{\bar{H}}] < \infty$ for some $A > 0$; $x_i, \ldots, x_n \in X \subset \mathbb{R}^d$ are vectors of covariates with $X$ being a compact set in $\mathbb{R}^d$ and $f_0$ is an unknown function needed to be estimated. We assume that $f_0 \in F$, where $F$ is the class of continuous functions with compact supports. Clearly, $f_0$ minimises the population criterion function

$$Q_n(f) = \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2\right]$$

$$= \frac{1}{n}\sum_{i=1}^{n}(f(x_i) - f_0(x_i))^2 + \sigma^2.$$

A least squares estimator of the regression function can be obtained by minimising the empirical squared error loss $Q_n(f)$:

$$\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{Q}_n(f) = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2.$$

However, if the class of functions $F$ is too rich, the resulting least squares estimator may have undesired properties, such as inconsistency (Shen and Wong 1994; Shen 1997; van de Geer 2000). Instead, we can optimise the squared error loss over some less complex function space $F_n$, which is an approximation of $F$ while the approximation error tends to 0 as the sample size increases. In the language of Grenander (1981), such a sequence of function classes is known as a sieve. More precisely, we consider a sequence of function classes,

$$F_1 \; F_2 \cdots F_n \; F_{n+1} \cdots F,$$

approximating $F$ in the sense that $\bigcup_1 F_n$ is dense in $F$. In other words, for each $f \in F$, there exists $\pi_n f \in F_n$ such that $d(f, \pi_n f) \to 0$ as $n \to \infty$, where $d(\cdot, \cdot)$ is some pseudometric defined on $F$. With some abuse of notation, an approximate sieve estimator $\hat{f}_n$ is

defined to be

$$Qnlfn) \therefore= \inf_{jE:Fn} Qn<J) + Op(T/n), \qquad (1)$$

where $TJn - 0$ as $n - oo$. We refer interested reader to Chen (2007) for a thorough discussion on sieve extremum estimators.

Throughout the rest of the paper, we focus on the sieve of neural networks with one hidden layer and sigmoid activation function. Specifically, we let

$$\mathcal{F}_{r_n} = \left\{ \alpha_0 + \sum_{j=1}^{r_n} \alpha_j \sigma \left( \gamma_j^T x + \gamma_{0,j} \right) : \gamma_j \in \mathbb{R}^d, \alpha_j, \gamma_{0,j} \in \mathbb{R}, \right.$$

$$\left. \sum_{j=0}^{r_n} |\alpha_j| \leq V_n \text{ for some } V_n > 4 \text{ and } \max_{1 \leq j \leq r_n} \sum_{i=0}^{d} |\gamma_{i,j}| \leq M_n \text{ for some } M_n > 0 \right\}, \qquad (2)$$

where $rn, Vn, Mn$ ↑ $oo$ as $n - oo$. For theoretical simplification, we impose bounded-ness assumption on the weights of the neural networks in $F_{,.,}$ which is related to the £1-regularisation when fitting a neural network. Such a method has been discussed in previous literatures (e.g. White 1989, 1990). In those papers, the consistency of the neural network sieve estimators has been established under random designs. However, there are few results on the asymptotic distribution of the neural network sieve estimators, which will be established in this paper. In terms of the rate of convergence, Chen and Shen (1998) obtained rate of convergence $Op((\frac{1}{10})^{\frac{4}{<i+i( 2dJ>}})$ for neural network sieve estimators. Later on, Chen and White (1999) improved the convergence rate to $O p(( _{10}.n-)^{4} \frac{l+2a/(d+l)}{d+a /(d+in} )$, where $a$ relates to the smoothness of the true function $Jo$ and in their paper, a central limit theorem for smooth functional of the estimated function is also provided. In Chen et al. (2001), rate of convergence was also obtained for stationary ,B- m ixing data. One important characteristic of the aforementioned rate of convergence is that as $d - oo$, both rates become $Op((_{1} n -) \perp f\frac{4}{2})$. Similar results can also be found in Barron (1994). It is well known from Stone (1982) that local smoothing methods suffer from the curse of dimensionality. But such phenomenon seems to vanish for approximate sieve extremum estimators based on neural networks. Bauer and Kohler (2019) also developed general theory and conditions to justify that neural networks can be used to circumvent the issue of curse of dimensionality. Hornik et al. (1989) showed that $\mathrm{Un} Frn$ is dense in F under the sup-norm But when considering the asymptotic properties of the sieve estimators, we use the pseudo-norm $\|/\| = \frac{1}{n} I \frac{1}{7} \frac{1}{1} f- (x;)$ (see Proposition 2.1 in the supplementary material) defined on $F$ and $F,•$.

With the increasing popularity of deep learning, recent research also starts to focus on the statistical propert ies of deep neural networks. For example, Schmidt -Hieber (2020) provided the rate of convergence for sparse deep neural network with Rectified Linear Unit (ReLU) activation function under the assumption that the underlying function is a composition of functions in some Holder space. Farrell et al. (2020) and Farrell et al. (2021) established the rate of convergence for deep neural network estimators when the underlying function belongs to a unit ball in a Sobolev space. It is worth pointing out that in these

two papers, no restrictions on the boundedness of weights in neural networks are imposed. The rate of convergence for deep ReLU neural networks has also been developed in Fabozzi et al. (2019) and Kohler and Langer (2021). It is worth mentioning that the results in Fabozzi et al. (2019) were developed based on the original version of this manuscript on ArXiv. We believe that the ultimate goal of developing these theories is to perform statistical inference based on neural networks for real-world problems and the results discussed in this paper may provide a starting point for further developments. For instance, Chen and White (1999) developed asymptotic normality for neural network sieve extremum estimators. However, to apply their result, a calculation of a series of covariance is essential, which may be hard to accomplish in practice. Recently, Horel and Giesecke (2020) developed a significance test based on neural networks. However, the theories in that paper are difficult to apply and verify in practice. Using similar techniques to be discussed in this paper, Shen et al. (2021) developed a goodness-of-fit test based on neural networks.

The remaining paper is organised as follows. In Section 2, we discuss the existence of neural network sieve estimato rs. The weak consistency and rate of convergence of the neural network sieve estimators will be established in Sections 3 and 4, respectively. Section 5 focuses on the asymptotic distribution of the neural network sieve estimators. Simulation results are presented in Section 6.

*Notations:* Throughout the rest of this paper, bold font alphabetic letters and Greek letters are vectors. C(X) is the set of continuous functions defined on $X$. The symbol means 'bounded above up to a universal constant' and $an \sim bn$ means $\quad$ -+ 1 as $n$ -+ oo. For a pseudo-metric space $(T, d)$, N(E, $T$, $d$) is its covering number, which is the minimum number of E-balls needed to cover $T$. Its natural logarithm is the entropy number and is denoted by H(E, $T$, $d$).

## 2. Existence

A natural question to ask is whether the sieve estimator based on neural networks exists. Before addressing this question, we first study some properties of $Fnr$• Proposition 2.1 shows that the sigmoid function is a Lipschitz function with Lipschitz constant $L = 1/4$.

**Proposition 2.1:** *A sigmoid function a* (z) $= fiZ'/(1 + f?')$ *is a Lipschizt function on* JR *with Lipschitz constant* 1/4.

The second proposition provides an upper bound for the envelope function $sup_{fE:F}$ ,,. $|f|$.

**Proposition** 2.2: *For each fixed n,*

$$\sup_{fE:F,n} \|f\|oo :S Vn.$$

Now we quote a general result from White and Wooldridge (1991) for readers who are not familiar with the theories of sieve extremum estimators. The theorem tells us that under some mild conditions, there exists a sieve approximate estimator and such an estimator is also measurable.

Theorem 2.1 (Theorem 2.2 in White and Wooldr idge (1991)): *Let* (Q, A, IP') *be a complete probability space and let* (8, $p$) *be a pseudo-metric space. Let* {8n} *be a sequence*

*of compact subsets of* e. *Let* Qn: *Q* x *en-+* i *be A® B(en)/B(i.)-measurable, and suppose that for each* w E Q, *Qn(w,•) is lower semicontinuous on en, n=1,2, Then for each n = 1, 2, . . ., there exists 0n : Q -+ en, A/ B(en)-measurable such that for each* w E Q, *Qn(W,0n(w)) = infeEBn Qn(w,0).*

Note that

$$Q_n(f) = \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}(f_0(x_i) + \epsilon_i - f(x_i))^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}(f(x_i) - f_0(x_i))^2 - 2\frac{1}{n}\sum_{i=1}^{n}\epsilon_i(f(x_i) - f_0(x_i)) + \frac{1}{n}\sum_{i=1}^{n}\epsilon_i^2.$$

Since the randomness only comes from E;'s, it is clear that $Q_n$ is a measurable function and for a fixed w, $Q_n$ is continuous inf Therefore, to show the existence of the sieve estimator, it suffices to show that $Fr''$ is compact in C(X), which is proved in the following lemma.

**Lemma 2.1:** *Let X be a compact subset of d. Then for each fixed n, $Frn$ is a compact set.*

**Proof:** For each fixed n, let Bn = [ao,. . . ,a ,", Yo, 1, . . . , yo,rn, Y [ , . . . ,y ]T belong to [- V n, Vnl'''+' x [- Mn, Mn l' " (d+I) := e n. If n is fixed, en is a bounded closed set and hence it is a compact set in rn(d+2)_+l Consider a map

$$H : (\Theta_n, \| \cdot \|_2) \to (\mathcal{F}_{r_n}, \| \cdot \|_n)$$

$$Bn \longmapsto H(Bn) = ao + \sum_{j=l}^{T_n} Lap \quad (rJ_x + Y_{0,j})$$

Note that $F,'' = H(en)-$ Therefore, to show that $F,''$ is a compact set, it suffices to show that $H$ is a continuous map due to the compactness of *en*. Let 8 l,n,82,n E en, then

IIH ( B1,n) - H (B2,n)II

$$\text{;; } E \left[ a_l'' + t aJ' a(J''' x, + vJ)') - al'' - t aJ' a(,!''' x, + yJJ') \right]$$

$$'' \text{ ;; } E [ b_i'' - \hat{a} 'l + t H'' a(,)'' x, + rJ;') - aJ' a(,!'''''' + ,m|]$$

$$;t \left[ \text{la}!_. \quad \text{al"I} + t \; \text{1aj'1'1 la} (r?''_x + ,J,?) -_a \; (rJ''_x + ,\text{ml} + \right.$$

$$\left. )_{-a} \;_J \text{lla} (r\backslash >_{\text{T}_x} + \quad \text{C 1}) \right]_{0,J}^{2}$$

$$-<; t \left[ t_{\text{la}^0 - a/''\text{I}+} \quad \cdot t_{,l(rJ^0} \quad r/') \;'^{+}\text{Ir},' - r,\text{TI}] \right]'$$

$$-< \left[ t_{\text{la})1' - aJ''\text{I}\pm} \quad (1 \; \text{V} \quad \text{llxllo}\emptyset) \; t_{\text{llrJ}!_.} \quad r/'t + \text{lrJ}) - ,\text{JJ'lr} \right.$$

$$:\text{S } \underline{(n} \; (1 \; \text{V !!x!}_{\infty})) \;^{2} \quad [rn(d + 1)]^{119!} ,n - 92,n!\text{I} .$$

Hence,  for any E > 0, we choose $8 = \text{E}/((1 \quad \text{V} \; \underline{\|x\|_{\infty}} \; )J \; rn(d + 1)),$. When $1191,n - 92,n112 < 8,$ we have

$$!1H \; (91 \; , n) - H \; (92,n) \; \text{lln} <_{\text{E},}$$

which implies that $H$ is a continuous map and hence $Fr,.$ is a compact set for each fixed $n$. ∎

Asa corollary of Lemma 2.1 and Theorem 2.1, we can easily obtain the existence of sieve estimator.

**Corollary** 2.1: *Based on the notations above, for each $n = 1, 2, \ldots$, there exists fn : Q - Fr,., A/B(Frn)-measurable such that $Q\,nifn(w)) = \inf \; \text{Effi} \; Q \; n(/)$.*

## 3. Consistency

In this section, we are going to show the consistency of the neural network sieve estimator. The consistency result leans heavily on the following Uniform Law of Large Numbers.

**Lemma 3.1:** *Under the assumption of*

$$[rn(d + 2) + 1]\text{V} \; log(V \, n[rn(d + 2) + 1] = o(n), \quad \text{as } n - \;_{00,}$$

*we have*

$$\sup_{f \in \mathcal{F}_{r_n}} |\mathbb{Q}_n(f) - Q_n(f)| \xrightarrow{p^*} 0, \quad \text{as } n \to \infty.$$

*Proof:* For any $\delta > 0$, we have

$$\mathbb{P}^*\left(\sup_{f\in\mathcal{F}_{rn}} |\mathbb{Q}_n(f) - Q_n(f)| > \delta\right)$$

$$= \mathbb{P}^*\left(\sup_{\forall \mathcal{F}_{rn}} \left!-\frac{1}{n}\sum_{i=1}^{t} \mathbb{E}f - a^2 - 2\frac{1}{n}\sum_{i=1}^{t} \mathbb{E};\{f(x;) - fo(x;)\} > \delta\right)\right.$$

$$\mathbb{P}'\left(\frac{1}{n}\sum_{i=1}^{t} \mathbb{E}^2;- a^2 > \frac{\delta}{2}\right) + \mathbb{P}^C\left(\sup_{\mathbb{E}\mathcal{F}_{rn}} \left|\frac{1}{n}\sum_{i=1}^{n} \mathbb{E};(j(x;) - fo(x;))\right| > \frac{\delta}{4}\right)$$

$$:= (I) + (II).$$

It follows from the Weak Law of Large Numbers that (I) $\to$ 0. Now, we are going to evaluate (II). By using the Markov's inequality, (II) $\to$ 0 holds if

$$\mathbb{E}^*\left[\sup_{f\in:\mathcal{F}rn} \left!-\frac{1}{n}\sum_{i=1}^{t} \mathbb{E};(f(x;) - fo(x;))\right] \to 0, \quad \text{as } n \to \infty.$$

Note that $\mathbb{E}[\mathbb{E}] = 0$ and each $f \in \mathcal{F}r''$ has its corresponding parameterisation $\theta n$. Since $\theta n$ is in a compact set, there exists a sequence $\theta nk$, $\to \theta n$ as $k \to \infty$ with $\theta n,k \in Q(n(d+Z)+1n$ $([-\sqrt{n}, \sqrt{n}Y''+1 \times [-Mn, MnY''d +])$. Each $\theta n,k$ Corresponds to a function $fic \in \mathcal{F}r$. Based on continuity, we have $fic(x) \to f(x)$ for each $x \in X$. From Example 2.3.4 in van der Vaart and Wellner (1996), we know that $\mathcal{F}r$. is P-measurable. Based on symmetrisation inequality, we have

$$\mathbb{E}^*\left[\sup_{f\in\mathcal{F}rn} \frac{1}{n}\sum_{i=1}^{t} \mathbb{E};(f(x;) - fo(x;))\right] \to 2\mathbb{E}.,\mathbb{E}\left[\sup_{f\in:\mathcal{F},.} \left!-\frac{1}{n}\sum_{i=1}^{t} ;\mathbb{E};\{f(x;) - fo(x;)\}\right],$$

where $1, \ldots, n$ are i.i.d. Rademacher random variables independent of $\mathbb{E}I, \ldots, \mathbb{E}n$. Based on the Strong Law of Large Numbers, there exists $N1 > 0$, such that for all $n \geq N1$,

$$\frac{1}{n}\sum_{i=1}^{t} \mathbb{E}f < a^2 + 1, \quad \text{a.s.}$$

For fixed $\mathbb{E}I, \ldots, \mathbb{E}n$, $I:=?1 \; ;\mathbb{E};(j(x;) - fo(x;))$ is a sub-Gaussian process indexed by $f \in \mathcal{F}rn$. Suppose that $(\delta, C, \mu,)$ is the probability space on which $1, \ldots, n$ are defined and let $Y(f,w) = I:=?1 \; ;(w)\mathbb{E};(j(x;) - fo(x;))$ with $f \in \mathcal{F}r''$ and $w \in \delta$. As we have shown above, we have $fie f$ and by continuity, $Y(/k,w) \to Y(f,w)$ for any $w \in \delta$. This shows that $\{Y(f,w), f \in \mathcal{F},.\}$ is a separable sub-Gaussian process. Hence Corollary 2.2.8 in van der Vaart and Wellner (1996) implies that there exists a universal constant $K$ and for any $J; \in \mathcal{F}rn$ with $n \geq N1$,

$$\mathbb{E}_\xi\left[\sup_{f\in\mathcal{F}_{rn}} \left|\frac{1}{n}\sum_{i=1}^{n} \xi_i\epsilon_i(f(x_i)) - f_0(x_i))\right|\right]$$

$$\leq \mathbb{E}_\xi\left[\left|\frac{1}{n}\sum_{i=1}^{n} \xi_i\epsilon_i(f_n^*(x_i)) - f_0(x_i))\right|\right] + K\int_0^\infty \frac{\overline{\log N\,(\tfrac{1}{2}rJ, \mathcal{F}r., d)}}{n}\,d17$$

$$\leq \mathbb{E}_\xi \left[ \left| \frac{1}{n} \sum_{i=1}^n \xi_i \epsilon_i (f_n^*(x_i) - f_0(x_i)) \right| \right]$$

$$+ K \int_0^{2V_n} \sqrt{\frac{\log N C \vdots + TJ, Frn \cdot ||\cdot||_{00})}{n}} \, dTJ,$$

where *forf,g* $\in$ *Frn•*

$$d(f,g) = \left( \phantom{xxx} \mathbb{E}; (j(x;) - g(x;))^2 \right)^{1/2}$$

sothat the last inequality follows by noting that $\sup_{1 < \text{ffn},} ||f||_{00} \quad V_n$ and

$$d(f,g) \quad ||f - g||_{00} \left( \phantom{xxx} \mathbb{E} f \right)^{1/2}$$

We then evaluate these two terms. For the first term, for $n \quad N1$, by Cauchy-Schwarz inequality, wehave

$$\mathbb{E}i \left[ t_{S,,,<J:(x;) - /0(x;))} \right];( t_{/r} \left( t_{/f:(x;) - /0(x;))r} \right)^2$$

$$\underline{Ja^2 \pm 1} \sup_{x \in X} |f;(x) - f_0(x)|, \quad \text{a.s.}$$

By choosing $J; = :rcrJo$ and using the universal approximation theorem introduced by Hornik et al. (1989), we know that $\sup_{x \in X} |f;(x;) - f_0(x_i)| \longrightarrow 0$ as $n \longrightarrow \infty$. Therefore, for any $\{ > 0$, there exists $N2 > 0$, such that for all $n \quad N2,$

$$\sup_{x \in X} [f;(x;) - f_0(x;)| < \frac{.k}{a+1}$$

By choosing $n \quad N1 \vee N2$, we get

$$\mathbb{E}_\xi \left[ \left| \frac{1}{n} \sum_{i=1}^n \xi_i \epsilon_i (f_n^*(x_i) - f_0(x_i)) \right| \right] < \zeta \quad \text{a.s.}$$

For the second term, we use the same bound from Theorem 14.5 in Anthony and Bartlett (2009) as we did in the proof of Lemma 2:

$$N\left( \frac{1}{2\sqrt{\sigma^2+1}} \eta, \mathcal{F}_{r_n}, ||\cdot||_\infty \right) \leq \left( \frac{8\sqrt{\sigma^2+1}e[r_n(d+2)+1]\left(\frac{1}{4}V_n\right)^2}{\eta\left(\frac{1}{4}V_n - 1\right)} \right)^{r_n(d+2)+1}$$

$$:= \tilde{B}_{r_n,d,V_n} \eta^{-[r_n(d+2)+1]},$$

where $B_{r,d,V_n} = (2-JaT'+Te[rn(d + 2) + 1]V / (V_n - 4))'(d + \frac{2}{4})+1$. Let

$$B_{rn,dY_n} = \log B_{rn,d,V_n} - [rn(d+2)+1]$$

$$= [rn(d+2)+1] \log \left(1 + \frac{2 a^2 + 1 e[rn(d+2)+1]V}{n V_4}\right) - 1$$

$$< 2[r_n(d+2)+1] \log \frac{[rn(d+ \frac{2}{4} + 1 ]V}{V_n - 4}, \quad \text{for all } n \geq N_1 \vee N_3,$$

wh ere N3 is chosen to satisfy $rn(d+2)+1 ::: 2 a^2 [r(d+1)]/4$. The last inequality then follows b noting that $V_n - V_n + 4 > 0$ for al $V$ so that $\log \frac{[rn(d+1)]V}{V_n-4} > \log \frac{vu-ri (n V_n)}{V_n-4}$

$Y \log \{2 a^2 + 1)$. We also have

$$H\left(\frac{JdT'+117,F,..,|\cdot|_{\infty}}{2 a^2 +1}\right) = \log B_{n,d,V_n} + [rn(d+2)+1] \log \frac{1}{17}$$

$$:S B_{n,,d,V,.} + [r n(d+2)+1]\frac{1}{17}$$

$$\leq B_{r_n,d,V_n}\left(1+\frac{1}{\eta}\right),$$

an d hence for all $n ::: N_1 \vee N_3$,

$$\int_0^{H^{1/2}\left(\frac{1}{2a+1} \frac{17,F,...,\|\cdot\|_{\infty}}{}\right)} d17$$

$$< B_{r,,d,Vn}^{1/2} \int_0^{\{v.} \left(1 + \frac{1}{17}\right)^{1/2} d 17$$

$$= B!.,,v.[1 1 (1+ \frac{3}{4} /2 d17 + 1 /2\nu ,.(1+\frac{3}{4} /2 d17]$$

$$:S B!.,,v. [2 \int_0^1 fo^{-17} 1/^2 d17 + 2(2Vn-1)]$$

$$:S 4 r,B^{1/2} V -\nu_L T_{n,d,V,.} n,$$

which implies that

$$\frac{H^C\left(\frac{2+1 17, Frn>|\cdot|_{\infty})}{n}\right)}{n} d17 :S 8 \frac{[r n(d+2)+1] \log r. (d t^2 \|V}{n}$$

$$\sim 8 \frac{[rn(d+2)+1]VJlog(Vn[rn(d+2)+1])}{n}$$

where the last part follows by noting that $\log v_{4} \sim \log V_n$. Under the assumption given in the Lemma, there exists $N_4 > 0$, such that for all $n :::\_ N_4$, we have

$$\sqrt{\frac{[rn(d+2)+1] \log(Vn[rn(d+2)+1])}{n}} \leq \frac{\{}{8}$$

Therefore, by choosing $n \geq$ N1 $\vee$ N2 $\vee$ N3 $\vee$ N4, we get

$$\mathbb{E}_\xi \left[ \sup_{f \in \mathcal{F}_{r_n}} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \epsilon_i (f(x_i) - f_0(x_i)) \right| \right] < 2\zeta \text{ a.s.,}$$

i.e. $\mathbb{E}$ [supfEF,. ¼ L⁷ 1 ;E;(j(x;) - Jo(x;))I] ---+ 0 a.s .. Moreove r, b ased on what we have shown, for a sufficiently large *n,* we have

$$\mathbb{E} \quad :f.. \, t \quad ;E;(f(xi) - Jo(x;)) ] S \#+l \qquad ilnr,Jo - Jolioo$$

$$+ \, 4 \, 2 \, KB \, {}^{112}_{rn,d,Vn} \, {-n} \, {}^{ll2} vn \, {---} + \, 0, \qquad \text{as ..}$$

Since

$$\mathbb{E}_\epsilon \left[ \sqrt{\sigma^2 + 1} \| \pi_{r_n} f_0 - f_0 \|_\infty + 4\sqrt{2} KB^{1/2}_{r_n, d, V_n} n^{-1/2} V_n \right]$$
$$= \sqrt{\sigma^2 + 1} \| \pi_{r_n} f_0 - f_0 \|_\infty + 4\sqrt{2} KB^{1/2}_{r_n, d, V_n} n^{-1/2} V_n \to 0 < \infty,$$

by using the Generalised Dom inated Convergence Theorem, we know that

$$\mathbb{E}^* \left[ \sup_{f \in \mathcal{F}_{r_n}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(x_i) - f_0(x_i)) \right| \right] \leq 2 \mathbb{E}_\epsilon \mathbb{E}_\xi \left[ \sup_{f \in \mathcal{F}_{r_n}} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \epsilon_i \left( f(x_i) - f_0(x_i) \right) \right| \right] \to 0,$$

which completes the proof. ∎

Based on the above lemmas, we are ready to state the theorem on the consistency of neural network sieve estimato rs.

**Theorem** 3.1: *Under the notation given above, if*

$$[rn(d + 2) + 1] V log(Vn[rn(d + 2) + 1] = o(n), \quad \text{as } n \text{---+ } 00, \qquad (3)$$

*then*

$$\| fn - Jo \|_{in} \xrightarrow{p} 0.$$

***Proof:*** Since Q is continuous at $Jo \in F$ and $Q(fo) = a^2 < oo$, for any $E > 0$, we have

$$\inf_{f : \| f - f_0 \|_n \geq \epsilon} Q_n(f) - Q_n(f_0) = \inf_{f : \| f - f_0 \|_n \geq \epsilon} \frac{1}{n} \sum_{i=1}^n (f(x_i) - f_0(x_i))^2 \geq \epsilon^2 > 0.$$

Hence, based on Lem ma 1, Lemma 3 and Corollary 2.6 in White and Wooldridge {1991), we have

$$\| fn - Jo \|_{in} \xrightarrow{p} 0.$$

∎

Remark 3.1: Wediscuss the condition (3) in Theorem 3.1 via somesimple examples here. If $<x_j = 0(1)$ for $j = 1, \ldots, r_n$, then $V_n = O(r_n)$ and

$$[r_n(d + 2) + 1)V \log(V_n[r_n(d + 2) + 1)) = O(r \log r_n)-$$

T herefore , a possible growth rate for the number of hidden units in a neural network is $T_n = o((n/\log n)^{1}1^3)$. On the other hand, if we have a slow growth rate for the number of hidden units in the neural network, such as $r_n = \log V_n$, then we have

$$[r_n(d + 2) + 1]V_n^2 \log(V_n[r_n(d + 2) + 1]) = \mathcal{O}((V_n \log V_n)^2).$$

Hence, a possible growth rate for the upper bound of the weights from the hidden layer to the output layer is $V_n = o(n^1 1^2 / \log n)$.

## 4. Rate of convergence

To obtain the rateof convergence for neural network sieves, we applyTheorem 3.4.1 in van der Vaart and Wellner (1996).

**Theorem 4.1:** *Based on the above notations, if*

$$\eta_n = \mathcal{O}\left(\min\{\|\pi_{r_n}f_0 - f_0\|_n^2, r_n(d + 2)\log(r_n V_n(d + 2))/n, r_n(d + 2)\log n/n\}\right),$$

*then*

$\|f_n - f_0\|_n$

$$= O_p\left(\max\left\{\|\pi_{r_n}f_0 - f_0\|_n, \sqrt{\frac{r_n(d + 2)\log[r_n V_n(d + 2))}{n}}, \sqrt{\frac{r_n(d + 2)\log n}{n}}\right\}\right).$$

**Proof:** Use the same bound from Theorem 14.5 in Anthony and Bartlett (2009), we have

$$\log N(\eta, \mathcal{F}_{r_n}, \|\cdot\|_n) \le \log N(\eta, \mathcal{F}_{r_n}, \|\cdot\|_\infty) \le \log\left(\frac{4e[r_n(d + 2) + 1]\left(\frac{1}{4}V_n\right)^2}{\eta\left(\frac{1}{4}V_n - 1\right)}\right)^{r_n(d+2)+1}$$

$$= [r_n(d + 2) + 1]\log\frac{\tilde{C}_{r_n,d,V_n}}{\eta},$$

where $c_{n,d} V_n = e_{r_n\,d}\,$ $> e$. Then from Lemma 3.8 in Mendelson (2003), for $\sigma < 1$,

$$\int_0^\delta \sqrt{\log N(\eta, \mathcal{F}_{r_n}, \|\cdot\|_n)}\, d\eta \le [r_n(d + 2) + 1]^{1/2}\int_0^\delta \sqrt{\log\frac{\tilde{C}_{r_n,d,V_n}}{\eta}}\, d\eta$$

$$\lesssim [r_n(d + 2) + 1]^{1/2}\delta\sqrt{\log\frac{\tilde{C}_{r_n,d,V_n}}{\delta}}$$

$$:= \phi_n(\delta).$$

Define h: $\delta \mapsto \langle/Jn(o)/\delta^\circ = [rn(d+2)+1]^{1/2}\underline{\theta^{\perp}}\cdot\theta J \log t, "\{v"}$. Since for $0 < \delta < 1$ and $1 < a < 2$

$$h'(;) \quad [,. (d+2)+1]^{>/>}((1-a)a-0 \quad \overline{\frac{-C,.d,v.}{\lg}} \quad \frac{1}{\delta} \quad \frac{\delta^2}{2} \quad \frac{C,.d,v.1}{c_{r,.}d.v.} \quad \frac{1}{1,^2} \quad \cdot{}^{.1/2}_{\text{og}} \quad \frac{C,.d,\psi}{1,}$$

$$[,. (d+2)+ \quad tJ'i' \left(\left(\middle| \quad -a)a- \quad \middle|_{\text{og}} \frac{C_{\bullet}^{\bullet}d,v.}{\delta^-} - \frac{1}{2}\middle|_{\text{og}} \cdot{}^{112}\frac{C_{g,}d,y.}{\delta^-}\right.\right.$$

$$< 0,$$

$\delta \mapsto \langle/Jn(o)/\delta^\circ$ is decreasing on $(0, oo)$. Let $Pn ;S 1!1'r.fo - foll;;-^{1}$. Note that

$$P\langle/Jn\underline{(:n)} = Pn[rn(d+2)+1]^{112}\underline{\log\frac{112}{}}(PnCr,,,d,Vn)$$

$$= [rn(d+2)+1]^{112}\underline{PnJlog}\,Pn + \log C,,,,d,v.$$

and

$$\log \tilde{C}_{r_n,d,V_n} = 1 + \log\frac{[r_n(d+2)+1]V_n^2}{V_n-4} \lesssim \log\frac{[r_n(d+2)+1]V_n^2}{V_n-4}$$

$$\sim \log[rnVn(d+2)],$$

we have

$$p'fi\langle/Jn\underline{(:n)};S - v'n\{\}\,rn(d+2)p'fi\,(logpn + \log[rnVn(d+2)1)\,;Sn.$$

Therefore, for

$$\overset{\cdot}{\underset{Pn}{Pn}}\,\,_{mt'n}\quad I\left(\frac{n}{-------------------}{rn(d+2)\log[rnVn(d+2)]}\right)^{/2}\left(\,..........\overset{n}{..........}\,\right)^{1/2} \Rightarrow$$

we have $P\overset{2}{n}\not{c}\,n(\underset{Pn}{.1..})\,;S\,..jn.$ Based on these observation, Lemma 1, Lemma 2 in the Supplementary Materials and Theorem 3.4.1 in van der Vaart and Wellner (1996) imply that

Ilfn - Jl',.f olln

$$= Op\left(\max\middle| 1!Jl',.fo - \,follm \quad rn(d+2)\,logn[rnVn(d+2)]\,\underline{;rn(d+n2)\,logn}\,)\right).$$

By using the triangle inequality, we can further get

Ilfn - Jolin :'.:: Ilf n - rr,Jolin + Iirr,Jo - Jolin

$$= Op\left(\max \text{Iirr,Jo} - Jolin, \quad \overline{\frac{rn(d+2)\log[rnVn(d+2))}{n}} \quad \overline{rn(d+:)\,logn}\,)\right).$$

■

**Remark 4.1:** Recall that a sufficient condition to ensure consistency is $rn(d + 2)$ $V \log[rnVn(d + 2)] = o(n)$. Under such a condition, $rn(d + 2)\log[rnVn(d + 2)]$ _:s n, the rate of convergence can be simplified to

$$\text{llfn - } Jolin = Op\left( \max \text{ lin,,Jo - Jolin, } r\;nd + n2)\; \text{lo g } n \right).$$

If we assume $Jo \in F$ where $F$ is the space of functions with finite firstabsolute moments of the Fourier magnitude distributions, i.e.

$$F = \left\{ t : Rd - R : J(x) = \int \exp\{iaTX\}\; d\mu, f(a), \text{IIILJ li1} \right.$$

$$:= \int_{\max(1 \text{iall1}, 1)} d i\mu, J1(a) : Sc \Big\},\tag{4}$$

where $/.Lf$ is a complex measure on $Rd$. IILJ I denotes the total variation of $/.Lf$, i.e. $[\mu, I(A) = \sup Z::_{,1} / \mu,(An) I$ and the supremum is taken over all measurable partitions $\{An\}_1$ of $A$. $1 \text{iali1} = I : f_1$ la il for $a = [a1,\ldots, adjY \in Rd$. Theorem 3 in Makovoz (1996) shows that $On := \text{llfo- } n,,Jolin ; S r-;;_{-1/(}^{112}{}^{2}_{d)}$. Therefore, if we let $d$ fixed and $Pn = 0,;;-1$ and $Vn = V$ in the proof of Theorem 4.1, $On$ must also satisfy the following inequality:

$$P </>\_n \;;S\; Pnr^{12}\log^{11\,2}(PnCmd,,v,.)\;;S,In$$

$$=\}\quad p\;rn\;\log pn + p\;rn\;\log r n\;;S\;n$$

$$=\}\quad rn^{1+1}{}_3\;rnlogr\;n\;;S\;n.$$

One poss ible choice of $rn$ to satisfy such condition is $rn ;,c::: (n/\log n)r+a$. In such a case, we obtain

$$\text{llfn - } Joli\; n = Op\left( \frac{n}{\log n} \right)^{- 4(1\frac{1}{4}+1/\,2d))},$$

which is the same rate obtained in Chen and Shen (1998). It is interesting to note that in the case where $d = 1$, we have $\text{llfn - } Joli\; n = O p((n/\log n\text{-})^{1}1^3)$. Such rate is close to the $Op(\text{-}n^{1}1^3)$, which is the convergence rate in nonparametric least square problems when the class of functions considered has bounded variation in $R$ (see Example 9.3.3 in van de Geer (2000)). As shown in Proposition 3 in the supplementary material, $F,.$ is a class of functions with bounded variation in $R$. Therefore, the convergence rateweobtained makes sense.

## 5. Asymptotic normality

To establish the asymptotic normality of sieve estimator for neural networks, we follow the idea in Shen (1997) and start by calculatingthe Gateaux derivative of the empirical

criterion function $Qn(f) = _n{}^{-1} \sum_{1} (y_i - f(x_i))^2$,

$$Q'_{n,f_0}[f - f_0] = \lim_{t \to 0} \frac{1}{t} \left[ \frac{1}{n} \sum_{i=1}^{n} (y_i - f_0(x_i) - t(f(x_i) - f_0(x_i)))^2 - \frac{1}{n} \sum_{i=1}^{n} (y_i - f_0(x_i))^2 \right]$$

$$= -\frac{2}{n} \sum_{i=1}^{n} E; (j(x;) - f_0(x;)).$$

Then the remainder of first-order functional Taylor series expansion is

$$Rn[f - f_0] = Qn(f) - Qn[f_0] - "nJ_0 [f - f_0] = \frac{1}{n} \sum_{i=1}^{n} (j(x_i) - f_0(x;))^2 = \|f - f_0\|_n^2.$$

As will be seen in the proof of asymptotic normality, the rate of convergence for the empirical process $\{_n{}^{-1/2} \sum_{1} E;(j(x;) - f_0(x;)) : f \in Fr_n\}$ plays an important role. Here we establish a lemma, which will be used to find the desired rate of convergence.

**Lemma 5.1:** *Let $X_i, \ldots, X_n$ be independent random variables with $X_i \sim P_i$. Define the empirical process $\{V_n(f)\}$ as*

$$v_n(f) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} [f(X_i) - P_i f].$$

*Let $F_n = \{f : \|f\|_\infty \leq \ldots : V_n\}$, $E > 0$ and a $\sup_f E F_n \cdot n^{-1} \sum_{1} Var[f(X;)]$ be arbitrary. Define to by $H(to, F_n, \| \cdot \|_\infty) = \frac{1}{4} \frac{1}{l(M, n, a)}$, where $1/l(M, n, a) = M^2/[2a(1 + \frac{}{})]$. If*

$$H(u, \mathcal{F}_n, \| \cdot \|_\infty) \leq A_n u^{-r}, \tag{5}$$

*for some $0 < r < 2$ and $u \in (0, a]$, where a is a small positive number, and there exists a positive constant $K_i = K(r, E)$, $i = 1, 2$ such that*

$$M \geq K_1 A_n^{\frac{2}{r+2}} V_n^{\frac{2-r}{r+2}} n^{\frac{r-2}{2(r+2)}} \vee K_2 A_n^{1/2} \alpha^{\frac{2-r}{4}},$$

*we have*

$$IP^* (\sup_{EF_n} |v_n(f)| > M) S S \exp\{-(1 - E) 1/l(M, n, a)\}.$$

*Proof:* The proof of the lemma is similar to the proof of Corollary 2.2 in Alexander (1984) and the proof of Lemma 1 in Shen and Wong (1994). Since $H(u, F_n, \| \cdot \|_\infty) S A_n^{-ur}$ for

some O < r < 2, we have

$$I(s,t) := \int \int H^{112}(u, Fn, \|\cdot\|_{\infty}) \, du \, S \, 2(2-r)^{-1} AJ \, t^{1-}f$$

Based on the assumption of

$$Ant^{-r} 2: H(to, J^{r''}m_{11\cdot11}) =^{E} ifr(M, n, a),$$

$$_0 \qquad _{00} \qquad 4$$

we have *to* S $[':4:)$l r/ . Note that ifr(M, *n,a*) 2: M $^2$/(4a) if MS *3,Jria/Vn* and 2(,Jria $+$ MVn/3) S *4MVn/3* if *M* 2: *3,Jria/Vn* and hence *ifr(M,n,a)* 2: *3,JriM/(4Vn)*. In summary,

$$\begin{array}{c} M\,2/(4a) \qquad \text{if } M < 3,Jria/Vn, \\ 1/t(M,n,a) \; 2:\{ \quad 3,JriM/(4Vn) \qquad \text{if } M\,2:3,/na/Vn. \end{array}$$

Therefore, if *M* 2: *3,/na/Vn,*

$$2B - 3/2J \left( \underset{64,/n'}{\underline{EM}} \, t \right) < 2\mathcal{E} - 3/2(2-r)^{-1} A1/2/\quad 5:$$
$$_0 \qquad\qquad\qquad - \qquad\qquad n \quad 0$$

$$\le 2^9 \epsilon^{-3/2}(2-r)^{-1}\left(\frac{4}{\epsilon}\right)^{\frac{1}{r}-\frac{1}{2}} A_n^{1/r}\left(\frac{3}{4V_n}\sqrt{n}M\right)^{\frac{1}{2}-\frac{1}{r}}$$

$$= \tilde{K}_1 A_n^{1/r} V_n^{\frac{1}{r}-\frac{1}{2}} n^{\frac{1}{4}-\frac{1}{2r}} M^{\frac{1}{2}-\frac{1}{r}},$$

where $\tilde{K}_1 = 2^9\epsilon^{-3/2}(2-r)^{-1}(\frac{4}{\epsilon})^{\frac{1}{r}-\frac{1}{2}}(\frac{3}{4})^{\frac{1}{2}-\frac{1}{r}}$. Hence

$$2^8\epsilon^{-3/2}I\left(\frac{\epsilon M}{64\sqrt{n}}, t_0\right) < M \Leftrightarrow \tilde{K}_1 A_n^{1/r} V_n^{\frac{1}{r}-\frac{1}{2}} n^{\frac{1}{4}-\frac{1}{2r}} M^{\frac{1}{2}-\frac{1}{r}} < M$$

$$\Leftrightarrow \tilde{K}_1 A_n^{1/r} V_n^{\frac{1}{r}-\frac{1}{2}} n^{\frac{r-2}{4r}} < M^{\frac{1}{r}+\frac{1}{2}}$$

$$\Leftrightarrow M > K_1 A_n^{\frac{2}{r+2}} V_n^{\frac{2-r}{r+2}} n^{\frac{r-2}{2(r+2)}},$$

where K1 $= kF$. On the other hand, if *M* < *3,/na/Vn,*

$$2\,8E - 3/2r \left( \underset{64,/n'}{\underline{EM}} \, t \right) < 2\,\mathcal{Q} - 3/2\,2\,2 - r\,)\cdot 1A\;1/\mathcal{Q}4\quad f$$
$$_0 \qquad\qquad\qquad - \qquad\qquad n \quad 0$$

$$\le 2^9\epsilon^{-3/2}(2-r)^{-1}\left(\frac{4}{\epsilon}\right)^{\frac{1}{r}-\frac{1}{2}} A_n^{1/r}\left(\frac{M^2}{4\alpha}\right)^{\frac{1}{2}-\frac{1}{r}}$$

$$= \tilde{K}_2 A_n^{1/r} M^{1-\frac{2}{r}} \alpha^{\frac{1}{r}-\frac{1}{2}},$$

where $\tilde{K}_2 = 2^9\epsilon^{-3/2}(2-r)^{-1}(\frac{4}{\epsilon})^{\frac{1}{r}-\frac{1}{2}}(\frac{1}{4})^{\frac{1}{2}-\frac{1}{r}}$. Hence

$$2^8\epsilon^{-3/2}I\left(\frac{\epsilon M}{64\sqrt{n}}, t_0\right) < M \Leftrightarrow \tilde{K}_2 A_n^{1/r} M^{1-\frac{2}{r}} \alpha^{\frac{1}{r}-\frac{1}{2}} < M$$

$$\Leftrightarrow \tilde{K}_2 A_n^{1/r} \alpha^{\frac{2-r}{2r}} < M^{\frac{2}{r}}$$

$$\Leftrightarrow M > K_2 A_n^{1/2} \alpha^{\frac{2-r}{4}},$$

where K $= K l_{,}^{r/2}$. In conclusion, if $M$ 2: $K A + v; +$ $\bullet - 2 \atop n![r+2i}$ $\vee$ $K_2 A^{1/2}{}_{aT,}^{2-,}$ then $2^8_E$ ${}^3$ ${}^2 J(\frac{1}{64} M_{,r} t)$ $< M$ By Theorem 2.1 in Alexander (1984), we have the desired reruk ∎

As a Corollary to Lemma 5.1, we can show that the supremum of the empirical process $\{n \ {}^{112} I=:7$ , $E_i(j(x;) - fo(xi))$ $: f \in Fr.\}$ converges to O in probabil ity.

**Coro llary** 5. 1: *Let Pn satisfy Pnllfn - fo lln $=$ Op(l) and Fr. be the class of neural network sieves as defined in* (2). *Then under the conditions*

(Cl) $rn(d + 2)Vn \log[rnVn(d + 2)] = o(n^1 1^4)$;
(C2) $np -;;^2/V = o\{l\}$,

*we have*

$$\sup_{\|f- fo\|n:S.P;;^1 JE:F,n} \left| \frac{1}{n} \sum \frac{\mathbb{I}}{\cdot \cdot} E_;(j - fo)(x;) \right| = op\{1\}.$$

**_Proof_** To estab lish the desired result, we apply the truncation device.

$$JP>* \left( \sup_{\|f-folln::'.Pn^1 JEJ'',n} \left| \sum_{l=I} t \ E_i(j - fo)(x;) \right| \geq M \right)$$

$$\leq \mathbb{P}^* \left( \sup_{\|f-fo\|_n \leq \rho_n^{-1} f \in \mathscr{F}_{rn}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i \mathbb{I}_{\{|\epsilon_i| \leq V_n\}} (f - f_0)(x_i) \right| \gtrsim M \right)$$

$$+ JP>* \left( \sup_{\|f- Jolln.9;;^1 J E:F,} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i \mathbb{I}_{\{|\epsilon_i| > V_n\}} (f - f_0)(x_i) \right| \gtrsim M \right)$$

$$:= (I) + (I I) .$$

For (I), we can apply Lemma 6 directly. Note that IElllfl l:S.V.}if - fo)(x) I :S $Vn<Vn +$ llfolloo) ;S $V$ since llfolloo < oo and for O < $T/< 1$,

$$\log N (TJ,F r. , ll \cdot lloo) :S \log \left( \frac{4e \ [rn(d+ 2) + 1] (1/4Vn)}{TJ(_4 Vn - 1} \right)^{2) r \ ,d +2) +l}$$

$$:S [rn(d + 2) + 1] (\log Cr. ,d, Vn + t-1)$$

$$= Cr,,,d, Vn (1 + t)$$

$$:S 2Cr,,,d,Vn - \frac{1}{T/},$$

where $\tilde{C}_{rn,d,V.} = \dfrac{e[rn(d+2)+1]V;}{Vn-4}$ d an

$$C_{r_n,d,V_n} = [r_n(d+2)+1]\log \bar{C}_{r_n,d,V_n} - [r_n(d+2)+1] \sim r_n(d+2)\log[r_n V_n(d+2)].$$

Therefore, Equation (5) in the main text is satisfied with $r = 1$ and $An = 2Cr\,d\,v$ . Following from Lemma 6, for $M$ $c2^l v$ $V$ $^{13}$-$n$ $1^6 v c$ $^{ll}v$ $a \mid 4$, we have (l) $f$: s e; p{- (1 - E)i/r(M, n, a)} and hence

$$\sup_{\|f-f_0\|\le \rho_n^{-1},f\in\mathcal{F}_{r_n}} \left| \frac{1}{\sqrt{n}}\epsilon_i \mathbb{I}_{\{|\epsilon_i|\le V_n\}}(f-f_0)(x_i)\right| = \mathcal{O}_p\left(\frac{C_{r_n,d,V_n}^{2/3} V_n^{2/3}}{n^{1/6}}\right).$$

From (Cl),

$$\frac{C_{r_n,d,V_n}^{2/3} V_n^{2/3}}{n^{1/6}} \sim \left(\frac{r_n(d+2)V_n\log[r_n V_n(d+2)]}{n^{1/4}}\right)^{2/3} = o_p(1).$$

For(II), by using the Cauchy- Schwarz inequality, we have

$$\left|\frac{1}{n}\sum_{i=1}^n \epsilon_i\mathbb{I}_{\{|\epsilon_i|>V_n\}}(f-f_0)(x_i)\right| \le \left(\frac{1}{n}\sum_{i=1}^n \epsilon_i^2\mathbb{I}_{\{|\epsilon_i|>V_n\}}\right)^{1/2}\|f-f_0\|_n.$$

Then it follows from the Markov inequality that

$$(II) \le \mathbb{P}\left(\left(\frac{1}{n}\sum_{i=1}^n \epsilon_i^2\mathbb{I}_{\{|\epsilon_i|>V_n\}}\right)^{1/2}\rho_n^{-1} \gtrsim Mn^{-1/2}\right) \lesssim M^{-2}n\rho_n^{-2}\mathbb{E}[\epsilon^2\mathbb{I}_{|\epsilon|>V_n}]$$

$$\lesssim M^{-2}n\rho_n^{-2}\frac{\mathbb{E}[|\epsilon|^{2+\lambda}]}{V_n^\lambda}.$$

Based on condition (C2), we have *(II)* - 0, and

$$\sup_{\|f-f_0\|_n\le \rho_n^{-1},f\in\mathcal{F}_{r_n}} \left|\frac{1}{n}\sum_{i=1}^n \epsilon_i\mathbb{I}_{\{|\epsilon_i|>V_n\}}(f-f_0)(x_i)\right| = o_p(1).$$

Combining the results we obtained above, we get

$$\sup_{\|f-f_0\|_n\le \rho_n^{-1},f\in\mathcal{F}_{r_n}} \left|\frac{1}{\sqrt{n}}\sum_{i=1}^n \epsilon_i(f-f_0)(x_i)\right| = o_p(1).$$

**Remark** 5.1 : Co ndi tion (C2) can be further simplifiedusing the resultsfrom Theorem 4.1. If

$$T/n = 0 \, (\min\{IIJrr, Jo \text{-} foll, rn(d+2)\log(rnVn(d+2))/n, rn(d+2)\log n/\,n\}) \,,$$

then

$$\rho_n^{-1} \asymp \max\left\{\|\pi_{r_n}f_0 - f_0\|_n, \sqrt{r_n(d+2)\log[r_n V_n(d+2)]/n}, \sqrt{r_n(d+2)\log n/n}\right\}.$$

It followsfrom condition (Cl) that

$$\rho_n^{-1} \asymp \max\left\{\|\pi_{r_n}f_0 - f_0\|_n, \sqrt{r_n(d+2)\log n/n}\right\}.$$

For simplicity, we assume that $p_{;,}^{;1} \times \underline{\sqrt{rn(d+2)\log n/n}}$, which holds for functions hav-ing finite first absolute moments of the Fourier magnitude distributions as discussed at the end of Section 4.4. Then in this case,

$$n\rho_n^{-2}/V_n^{\lambda} \asymp r_n(d+2)\log n/V_n^{\lambda},$$

so that condition (C2) becomes $rn(d+2) \, \log n/V!{\to} 0$.

Now we are going to establish the asymptotic normality for neural network estimators. For/ $\in$ ff $\in$ $Frn$ : llf - Jo lin S $p_{;,}^{;1}$}, we consider a local alternative

$$\tilde{f}_n(f) = (1 - \delta_n)f + \delta_n(f_0 + \iota), \tag{6}$$

where O S $8n = TJ!^{12} = o(-n^{-1}1^2)$ is chosen such that $Pn8n = o(1)$ and$t(x) = 1$.

**Theorem 5.1 (Asymptotic Normality):** *Suppose that* O S $T/n = o(-n^{-1})$ *) and conditions* (Cl) *and*(C2) *in Corollary 5.1 hold. We further assume that the following two conditions hold:*

(C3) $\sup_{f\in\mathcal{F}_{rn}:\|f-f_0\|_n\leq\rho_n^{-1}} \|\pi_{r_n}\tilde{f}_n(f) - \tilde{f}_n(f)\|_n = \mathcal{O}_p(\rho_n\delta_n^2);$

(C4) $\sup_{f\in\mathcal{F}_{rn}:\|f-f_0\|_n\leq\rho_n^{-1}} \frac{1}{n}\sum_{i=1}^{n} \epsilon_i(\pi_{r_n}\tilde{f}_n(f)(x_i) - \tilde{f}_n(f)(x_i)) = \mathcal{O}_p(\delta_n^2),$

*then*

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \left[\hat{f}_n(x_i) - f_0(x_i)\right] \xrightarrow{d} \mathcal{N}(0, \sigma^2).$$

**Proof** The ma in ideaof the proof is to usethe functional Taylor series expansion for $Qn(j)$ and then carefully bound each term in the expansion. For any $f \in$ ff $\in$ $Fr,,$ : llf - $Jo_{lln}$ S

*Pñ* [1],

$$\mathbb{Q}_n(f) = \mathbb{Q}_n(f_0) + \mathbb{Q}'_{n,f_0}[f - f_0] + R_n[f - f_0]$$

$$= \frac{1}{n}\sum_{i=1}^{n}\epsilon_i^2 - \frac{2}{n}\sum_{i=1}^{n}\epsilon_i(f(x_i) - f_0(x_i)) + \frac{1}{n}\sum_{i=1}^{n}(f(x_i) - f_0(x_i))^2. \qquad (7)$$

Note that

$$llfn(f) - Jolin = 11(1 - on)\}n + Onlfo + l) - Jolin$$

$$= 11(1 - On)ifn - Jo) + Ontll\ n$$

$$S\ (1 - On\ )llfn - Joli\ n + On,$$

and since $On = o(n^{-1/2})$, we can know that with probability tending to 1, lifn(f) - Jolin S $p;;$ Then replacing $f$ in (7) by $]n$ and $n:,,Jn(f)$, we get

$$<r:lnlfn) = \underset{n}{\overset{n}{\underset{i=1}{L}}}Ef - \underset{n}{\overset{n}{\underset{i=1}{L}}}E;lfn(x;) - fo(x;)) + llfn - foll$$

$$\mathbb{Q}_n(\pi_{r_n}\tilde{f}_n(f)) = \frac{1}{n}\sum_{i=1}^{n}\epsilon_i^2 - \frac{2}{n}\sum_{i=1}^{n}\epsilon_i(\pi_{r_n}\tilde{f}_n(f)(x_i) - f_0(x_i)) + \|\pi_{r_n}\tilde{f}_n(f) - f_0\|_n^2.$$

Subtracting these two equations yields

$$<r:lnlfn) = <r:ln(n:,,Jn(f)) + \overset{n}{\underset{l=1}{L}}Ei\ (\ n:,,Jn(f)(x;) - ]n(x;))$$

$$+\ llfn - Joli\ n - lirr;.fn(f) - J\mathring{o}lin-$$

Now note that

$$\|\pi_{r_n}\tilde{f}_n(f) - f_0\|_n^2 = \|\pi_{r_n}\tilde{f}_n(f) - \tilde{f}_n(f) + \tilde{f}_n(f) - f_0\|_n^2$$

$$= \|\pi_{r_n}\tilde{f}_n(f) - \tilde{f}_n(f) + (1 - \delta_n)(\hat{f}_n - f_0) + \delta_n\iota\|_n^2$$

$$S\ (l\ \cdot\ On)^2\ llfn - Jol\hat{i}n + 2(1 - On)\}n - fo,Onl + o\dot{n}$$

$$+\ 2(1 - On)lln:,,Jn(f) - ]n(f)llnllfn - Jolin$$

$$+\ 2on lrrr,Jrf) - \bar{f}n(f)lln + lln:,,\ Jn(f) - \bar{f}n(f)ll^2\ n,$$

where the last inequality follows from the Cauchy-Schwarzinequality.Since

$$\underset{l=1}{\overset{n}{E;(}}\ n:,,Jn(f)(x;) - Jn(x;)) = \overset{n}{\underset{l=1}{L}}Ej(\ n:,,Jnlf)(xi) - Jnlf)(xi) + Jn(f)Cxi) - Jn(x;))$$

$$= \overset{n}{\underset{l=1}{L}}Ei\ (\ n:,,Jn(f)(Xi) - ]nlf)(Xi))$$

$$-\ \frac{2}{n}\ On\overset{n}{\underset{i\ 1}{L}}\ Ei\{\ ]n(x;) - fo(xi)) - \frac{2}{n}\ On\overset{n}{\underset{i=1}{L}}\ Ei,$$

by the definition of $J_n$, we have

$$- O p(\delta n)\, S \inf_{f E:\bar{F}_{,,,}} Qn(J) - Qnlfn)\, \hat{S}\, Qn(7rr,Jn(J)) - Qnlfn)$$

$$S\, (\!(\,\cdot\,\, \delta n)^2\, llfn -\, Joli\hat{n} -\, llf\, n -\, Joli\hat{n} + 2(1 -\, \delta n)\delta n\, \overset{lr}{yn} -\, fo,L\!)$$

$$+\, 2(1 - \delta n)llfn - Jolinlin,Jnlf) - Jn(J)lin$$

$$+\, 2\delta nli7rr,\overline{Jn(J)} -\, \overline{Jn(J)}lin\, +\, li lrr,\overline{Jn(J)} -\, \overline{Jn(J)}li\overset{2}{n}$$

$$-\, \overset{n}{\underset{1=1}{LE}};(n,Jnlf)(x;) -\, Jnlf)(x;))+\, 8n\, \overset{n}{\underset{1=1}{E}};(!\, n(x\,;) - Jo(x;))$$

$$+\, -\overset{2}{\underset{n}{8n}}\, \overset{n}{\underset{i=1}{LE}};+\, Op(\delta\,)$$

$$S\, 8\hat{n}llfn -\, Joli\hat{n} + 2(1 -\, \delta n)\delta n\, \overset{lr}{yn} -\, Jo,\!)$$

$$+\, 2(1 - \delta n)llfn - Jolinli1r,Jn(J) - ]n(J)lin$$

$$+\, 2\delta nll7r,;Jn(J)- Jn(J)lln+\, lin,,Jnlf) -\, Jnlf)li\hat{n}$$

$$-\, t\, \underset{z=1}{E};(\, 1r,Jn(J)(x;) - ]n(J)(x;))\, +\, 8n\, t\, \underset{z=1}{E};(!n(x;) - Jo(x;))$$

$$+\frac{2}{n}\delta n \sum_{i=1}^{n} \epsilon_i + \mathcal{O}_p(\delta_n^2), \tag{8}$$

where the last inequality follows by noting that $(1 - \delta n)2 - 1 = -2\delta n + 8\, S\, \delta$. From the condition (Cl), we can get

$$[rn(d + 2) + 1]V \log[rnVn(d + 2) + 1\,]$$

$$S\, ([rn(d + 2) + 1\,]Vn \log[rnVn(d + 2) + 1])4 = o(n).$$

Combining with Theorem 2, we obtain that $llf\, n - Jo_{11}n = Op(1\,)$ and hence $\delta;llfn - Jo_{11} = op(\delta;)$. From condition (C3), we have

$$2(1 - \delta n)\, llfn - fo\, ll\, n\, lln\, ,J\, nlf\,) - Jn\, (J)\, lln\, .S\, 2\, llfn - Jo\, lin\, lln\, ,J\, nlf\,) - Jn\, (J)\, lln$$

$$= Op\, (p;^{1}Pn\delta\,) = O\, p(\delta\,).$$

Similarly, since $Pn\delta n = o(l)$, we have

$$2\delta n lln\, ,,\overline{J}\,\, n(J) -\, \overline{Jn(J)}lln =\, Op(\delta n \cdot Pn\overset{2}{\delta}n) =\, Op(\overset{2}{\delta}n)$$

$$lln\, ,,\overline{J}\,\, n(J)\, -\, \overline{Jn}\, lf\,)\, ll\overset{2}{n} = Op(Pn\overset{2}{\delta}n) =\, o\, p(\delta\overset{2}{n}).$$

Based on condition (C4), we know that

$$\frac{2}{n}\sum_{i=1}^{n} \epsilon_i \left(\pi_{r_n}\tilde{f}_n(f) - \tilde{f}_n(f)\right) = \mathcal{O}_p(\delta_n^2),$$

and from Corollary 5.1, we also have

$$2_{;;}\delta n^{"}L, \quad \mathrm{E;} \{[n(x_;) - fo(x_;)\} = Op(\delta n \cdot n^{-1/2}).$$

It follows from these observations that

$$-2(1 - \delta_n)\left(\hat{f}_n - f_0, \delta_n\iota\right) + \frac{2\delta_n}{n}\sum_{i=1}^{n}\epsilon_i \le \mathcal{O}_p(\delta_n^2) + o_p(\delta_n^2) + o_p(\delta_n \cdot n^{-1/2}),$$

which implies that

$$-(1 - \delta_n)\left(\hat{f}_n - f_0, \iota\right) + \frac{1}{n}\sum_{i=1}^{n}\epsilon_i \le \mathcal{O}_p(\delta_n) + o_p(n^{-1/2}) = o_p(n^{-1/2}).$$

By replacing $t$ with $-t$, we can obtain the same result and hence

$$\left|\left(f/n - fo, t\right) - \sum_{l=1}^{n}\overset{n}{L}\mathrm{E;} s \right| (1 - \delta n)\,f/n - fo, t) - \sum_{l=1}^{n}\overset{n}{L}\mathrm{E;}\right| + \delta n \,|f\hat{h} - Jo, i)|$$

Therefore,
$$\le o_p(n^{-1/2}) + \delta_n\|\hat{f}_n - f_0\|_n$$
$$= o_p(n^{-1/2}).$$

$$\left(f\hat{n} - fo\right)\cdot , t - \frac{1}{;;}\sum_{l=1}^{n}\mathrm{E;} + op(n^{-1/2}),$$

and the desired result follows from the classical Central Limit Theorem. ∎

Let us focus on the conditions given in the theorem. Note that if (Cl) holds, we have

$$r_n(d+2)V_n^2 \log[r_n V_n(d+2)] \le [r_n(d+2)]^4 V_n^4 \left(\log[r_n V_n(d+2)]\right)^4 = o(n),$$

so it is a sufficient condition to ensure the consistency of the neural network sieve esti-mator. As in Remark 3.1, we consider some simple scenarios here. If $Vn = O(rn)$ , then $Tn(d+2)Vn \log[rnVn(d+2)] = O(\log Tn)$ so that a possible growth rate for $Tn$ is $Tn = o(n^1 1^8 / (\log n)2)$. On the other hand, ifrn = log$Vn$, then $Tn(d+2)Vn \log[rnV n(d+2)] = O(Vn(\log Vn)2)$ and a possible growth rate for $Vn$ is $Vn = o(n^{-1}1^4 /(\log n)^2)$. Thus, in both cases, the growth rate required for the asymptotic normality of neural network sieve esti-mator is slower than the growth rate required for the consistency as given in Remark 3.1. One explanation is that due to the Universal Approximation Theorem, a neural network with one hidden layer can approximate a continuous function on compact support arbi-trarily well if the number of hidden units is sufficiently large. Therefore, if the number of hidden units is too large, the  neural network sieve estimator $\hat{f}n$ may be very close to the best projector of the true function Jo in $Fr$. so that the error $\mathrm{E}\hat{\iota}$  $[fn(x_;) - fo(x_;)]$ could becloseto zero, resulting a small variation. Byallowing slower growth rateof the number of

hidden units can increasethe variations of I: $=\!\!\exists_1$ [fn(x;) - fo(xi)], which makes the asymptotic normality more reasonable. On the otherhand, condition (C3) and condition (C4) are similar conditions as in Shen (1997), which are known for conditions on approximation error. These conditions indicate that the approximation rate of a single layer neural network cannot be too slow, otherwise it may require a huge number of samples to reach the desired approximation error. Therefore, the conditions in the theorem can be considered as a trade-off between bias and variance.

Theorem 5.1 can be used directly for hypothesis testing of neural network with one hidden layer if we know the variance of the random error $a^2$. In practice, this is rarely the case. To perform hypothesis testing when $a^2$ is unknown, it is natural to find a good estimator of $a^2$ and use a 'plug-in' test statistic. A natural estimator for $a^2$ is

$$\hat{\sigma}_n^2 = \frac{1}{n}\sum_{i=1}^n \left(y_i - \hat{f}_n(x_i)\right)^2 = \mathbb{Q}_n\left(\hat{f}_n\right).$$

We then need to establish the asymptotic normality for the statistic $a\text{-}.Jn$ $L=7\,I$ [fn(x;) - fo(x;)].

**Theorem** 5.2 **(Asymptotic Normality for Plug-in Statistic):** *Suppose that Jo* E *C(X), where X* C *lIlld is a compact set and* O _:'s $T/n = o(\text{-}n^1)$. *Then under the conditions as stated in Theorem* 5.1, *we have*

$$O\text{-}n\backslash \quad \overset{n}{\quad} [Jn(Xi) - fo(x;)] \overset{.}{1}\ N(0,1).$$

*Proof:* Note that

$$a_;^{\cdot} = Qnlfn) = \overset{}{\underset{t=1}{\sum}} (>,i\text{-}fn(X;)r = \overset{}{\underset{t=1}{\sum}} (to(x;) + {}_{Ej\text{-}}\ fn(X;)r$$

$$= \text{-}^{In}\underset{n\quad i=l}{\sum} (!n(x;) - fo(xi))^2 - \text{-}^{2n}\underset{n\quad i=l}{\sum}\!E; (!n(x;) - fo(xi)) + \text{-}^{I\,n}\underset{n\quad i=l}{\sum}\!Et$$

$$= {}^1_;;\ \underset{t=1}{\sum} L,\ E^2i\text{-}\ {}^2_{;;}\ \underset{t=I}{\underset{E;}{\sum}} \left(fn(x;) - fo(Xi)\right) + \text{llfn} - Joll^2n$$

Based on the rate of convergence of $\hat{f}n$ we obtained in Theorem 4.1 and condition (Cl), we know that

$$lvn - !0[ = 0;\left(\max\,\{\,\text{lln,,Jo} - \text{foll}\,,\,\underline{rn(d+;)logn}\,\}\right).$$

Under (C3), $\text{llrr,,Jo} - \text{foll} = o(p\,o) = o(\text{-}n^{1}\!{}^{1^2})$ and under (Cl), we have

$$\underline{c\,n(d+\,:)\,b\,g\ n)}\ :'s\ {}_o(\,\underline{n}\,^{1\,4}_{L}\,_{:og\,n\,}) $$

$$\overset{=}{\text{-}}\ o\ \frac{(\text{logn})}{n\,3/4}\ \overset{112}{\text{-}}\ \overset{}{\text{-}}\overline{o(n}\ ), $$

which implies that $\|f_n - f_0\| = o_p(n^{-1/2})$. Moreover, by the same arguments as in the proof of Theorem 5.1, we can show that

$$2 \frac{1}{n} \sum_{i=1}^{n} \epsilon_i \left( \hat{f}_n(x_i) - f_0(x_i) \right) = o_p(n^{-1/2}).$$

Therefore

$$\mathbb{Q}_n(\hat{f}_n) = \frac{1}{n} \sum_{i=1}^{n} \epsilon_i^2 + o_p(n^{-1/2}).$$

Based on the Weak Law of Large Numbers, we know that $\frac{1}{n} \sum_{i=1}^{n} \epsilon_i^2 = \sigma^2 + o_p(1)$. Therefore ,

$$\hat{\sigma}_n^2 = \mathbb{Q}_n(\hat{f}_n) = \sigma^2 + o_p(1),$$

and it follows from the Slutsky's Theorem and Theorem 5.1in the main text, we obtain

$$\frac{1}{\hat{\sigma}_n \sqrt{n}} \sum_{i=1}^{n} \left[ \hat{f}_n(x_i) - f_0(x_i) \right] = \frac{\sigma}{\hat{\sigma}_n} \frac{1}{\sigma \sqrt{n}} \sum_{i=1}^{n} \left[ \hat{f}_n(x_i) - f_0(x_i) \right] \xrightarrow{d} \mathcal{N}(0,1). \qquad \blacksquare$$

## 6. Simulation studies

In this section, simulations were conducted to check the validity of the theoretical results obtained in the previous sections. The consistency of the neural network sieve estimators was examined under various simulation scenarios. Finally, we evaluated the asymptotic normality of the neural network sieve estimators . For illustration purpose, we only include the simulations where the dimension of the covariates is 1. More simulations for the multivariate cases are given in the supplementary materials.

### 6.1. Consistency for neural network sieve estimators

In this simulation, we are going to check the consistency result from Section 3 and the validity of the assumption made in Theorem 3.1. Based on our construction of the neural network sieve estimators, in each sieve space $\mathcal{F}_{r_n}$, there is a constraint on the $l_1$ norm for $a: \|a\|_0 \leq V_n$. So finding the nearly optimal function in $\mathcal{F}_{r_n}$ for $\mathbb{Q}_n(f)$ is in fact a constrained optimisation problem. A classical way to conduct this optimisation is through introducing a Lagrange multiplier for each constraint. Nevertheless, it is usually hard to find an explicit connection between the Lagrange multiplier and the upper bound in the inequality constraint. Instead, we use the subgradient method as discussed in section 7 in Boyd and Mutapcic (2008). The basic idea is to update the parameter $a_0, \ldots, a_{r_n}$ through

$$\alpha_i^{(k+1)} = \alpha_i^{(k)} - \delta_k g^{(k)}, \quad i = 0, \ldots, r_n,$$

where $\delta_k > 0$ is a stepsize and $\delta_k$ is chosen to be $0.1/\log(e + k)$ throughout the simulation, which is known as a nonsummable diminishing step size rule. $g_k$ is a subgradient of the

objective or the constraint function $Lj \ O \ |a| \cdot$ 𝒱n at a <k). More specifically, we take

$$
\begin{cases}
\partial_{\alpha^{(k)}} \mathbb{Q}_n(f) & \text{if } \sum_{j=0}^{r_n} |a_j| \quad Vn \\
\partial_{\alpha^{(k)}} \sum_{j=0}^{r_n} |\alpha_j| & \text{if } \sum_{j=0}^{r_n} |a_j| > Vn,
\end{cases}
$$

**Table 1.** Comparison of errors $\|f_n - f_0\|$ and the least square errors $Q_n(f_n)$ after 20,000 iterations under different sample sizes.

| Sample sizes | Neural network | | Sine | | Piecewise continuous | |
|---|---|---|---|---|---|---|
| | $\|f_n - f_0\|$ | $Q_n(f_n)$ | $\|f_n - f_0\|$ | $Q_n(f_n)$ | $\|f_n - f_0\|$ | $Q_n(f_n)$ |
| 50 | 3.33E-2 | 0.519 | 6.04E-2 | 0.513 | 6.20E-1 | 1.124 |
| 100 | 2.79E-2 | 0.552 | 3.04E-2 | 0.587 | 3.20E-1 | 0.920 |
| 200 | 6.0SE-3 | 0.500 | 1.05E-2 | 0.501 | 2.51E-1 | 0.786 |
| 500 | 8.1SE-3 | 0.484 | 1.19E-2 | 0.499 | 3.26E-1 | 0.769 |
| 1000 | 3.02E-3 | 0.475 | 1.54E-2 | 0.480 | 2.98E-2 | 0.489 |
| 2000 | 2.88E-3 | 0.500 | 9.72E-3 | 0.506 | 1.69E-2 | 0.515 |



**Figure 1.** Comparison of the true function and the fitted function for three different types of non-linear functions. The top panel shows the scenario when the true function is a single layer neural network; the middle panel shows the scenario when the true function is a sine function, and the bottom panel shows the scenario when the true function is a continuous function having a non-differentiable point. As we can see from all the cases, the fitted curve becomes closer to the truth as the sample size increases.
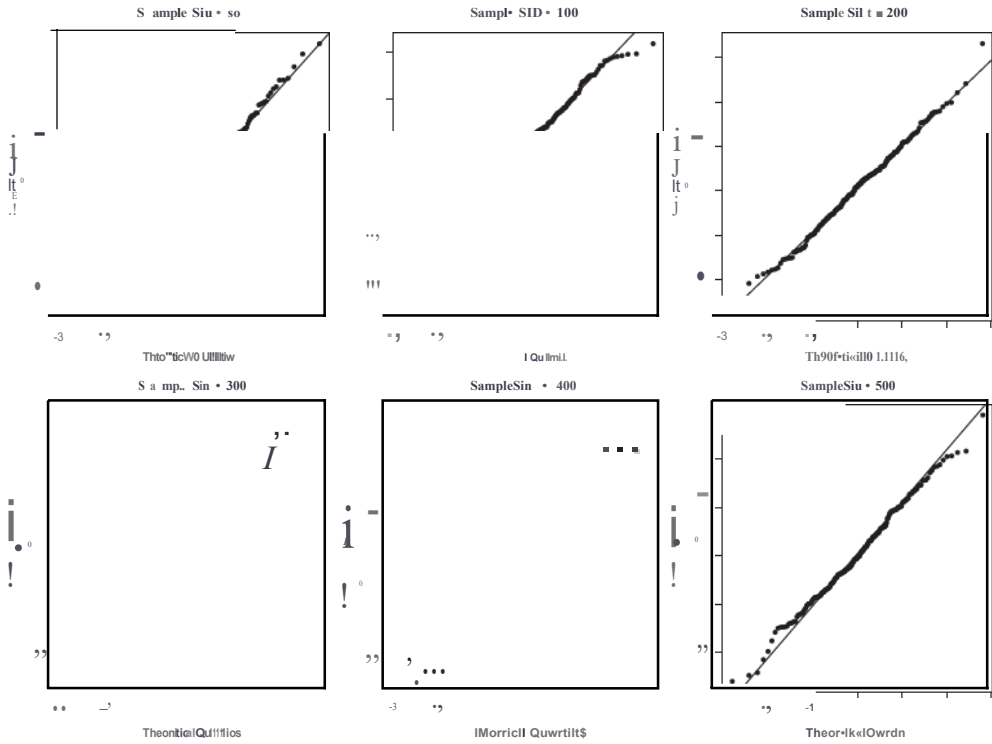
**Figure 2.** Normal Q-Q plot for $n^{-1/2} L = l$ [fn(x1) - fo(x1)] various sample sizes. Thetrue function fo is a single-layerneural network with two hidden units asdefined in (10).

The updating equations of $y_1, \ldots y'''$, $Yo,1, \ldots, Yo,r,$, remain the same as those in the classical gradient descent algorithm.

We simulated the response through the following model:

$$y_i = fo(x_i) + \varepsilon_i, \quad i = 1, \ldots, n, \tag{9}$$

where the total sample size $n$ varies from 50 to 2000, $x_t, \ldots ,x_n \sim$ i.i.d. $N(0, 1)$, $\varepsilon_1, \ldots, \varepsilon_n \sim$ i.i.d. $N(0,0.7^2)$. For the truefunction $Jo(x)$,weconsidered the following three functions:

(1) A neural network with a single hidden layer and two hidden units:

$$fo(x_i) = -1 - a(2x_i + 1) + a(-x_i + 1).$$

(2) A trigonometric function:

$$fo(x) = \sin(x) + \cos(x+1) \tag{11}$$

(3) A continuous function havinga non-differential point

$$f_0(x) = \begin{cases} -2x & \text{if } x \leq 0 \\ \sqrt{x}\left(x - \dfrac{1}{4}\right) & \text{if } x > 0. \end{cases} \tag{12}$$
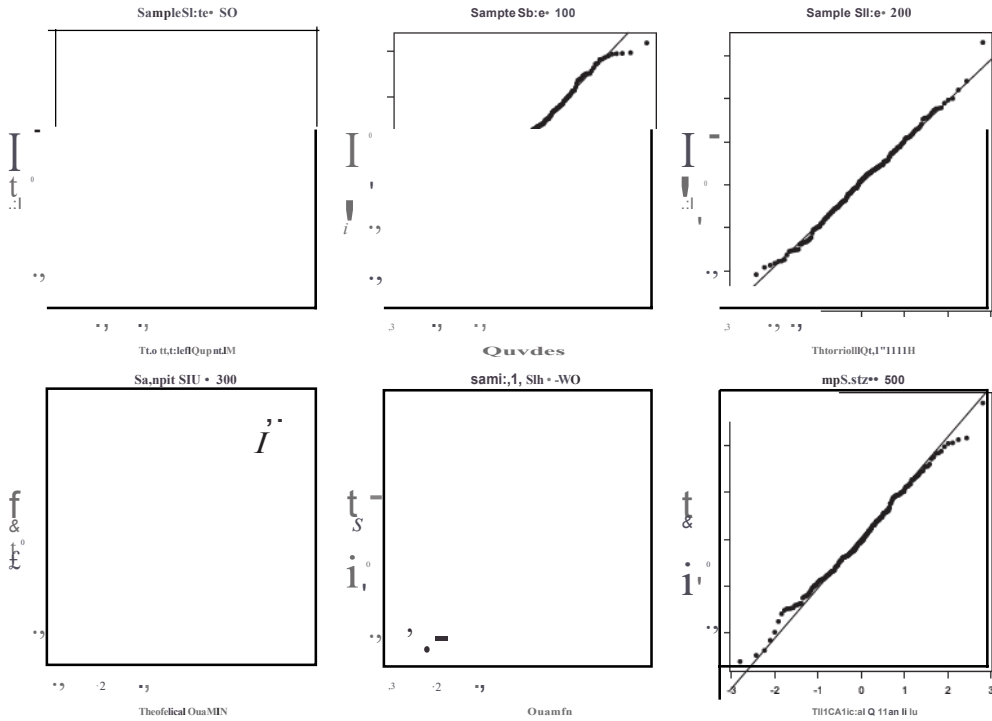
**Figure 3.** Normal Q- Q plot for-n $^{1/2}$ $I$: $_1$ *[fn(x1) - fo(x1)]* various sample sizes. The true function fo is a trigonometric function as defined in (11).

We then trained a neural network using the subgradient method mentioned above and set the number of iterations used for fitting as 20,000. We chose the growth rate on the number of hidden units $rn = n^{1/4}$ and the upper bound for .f.$_1$ norm of the weights and bias from the hidden layer to the output layer $Vn = 10n^{1/4}$. Such choice satisfies the condition mentioned in Remark 3.1 and hence satisfies the condition in Theorem 3.1. We compared the errors llf n - fo ll and the least square errors *Qnifn)* under different sample sizes. The results are summarised in Table 1.

As we can see from Table 1, the errors llf n - foll overall has a decreasing pattern as the sample size increases. There are some cases where the error becomes a little bit larger when the sample sizes increases (e.g. the errors using 500 samples in all scenarios is larger than those errors using 200 sample). One explanation is that the number of hidden units increases from 3 (for 200 samples) to 4 (for 500 samples) under our simulation setup, which adds variation to the estimation performance. Overall, the table shows that the estimated function]n is indeed consistent in the sense that llfn - f oll n = *o;(l)*. Figure 1 plots the fitted functions and the true function, from which we can straightforwardly visualise the result more and draw the conclusions.

## 6.2. *Asymptotic normality for neural network sieve estimators*

The last part of the simulation focuses on the asymptotic normality derived in Theorem 5.1. We still considered the same types of true functions as described in Section 6.1 but sampled
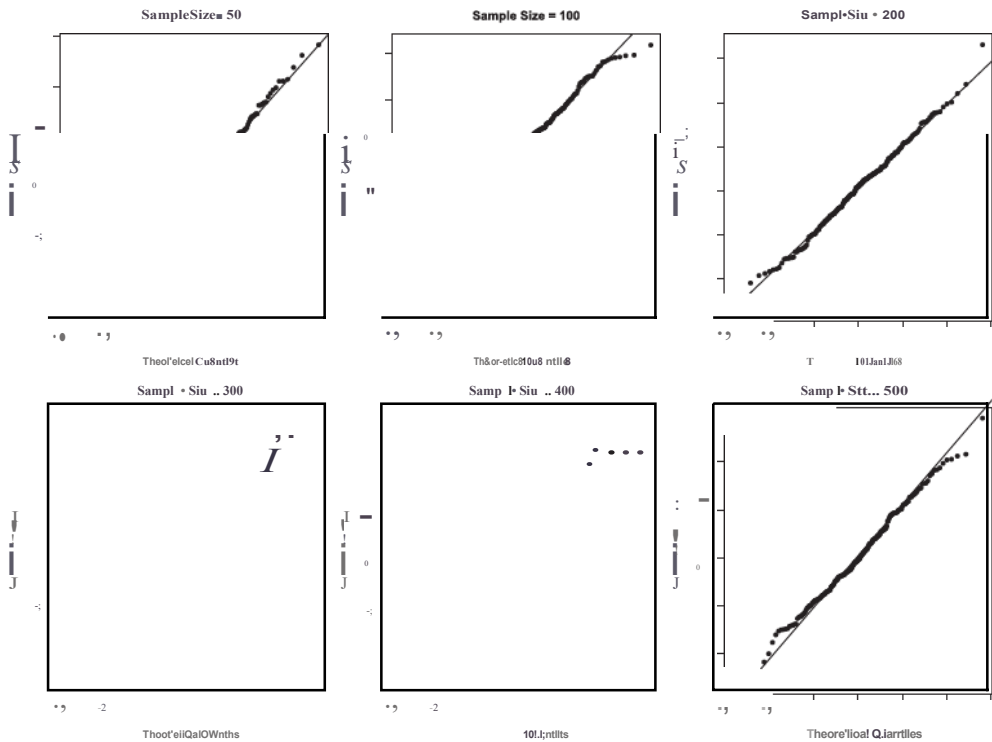
**Figure 4.** Normal Q- Q plot for $n-{}^{1}1^{2} L =l$ *[fn( x1) - fo(x1)l* varioussamplesizes. Thetrue function fo is a continuous function having a non-differential point as definedin (12).

the random errors from the standard normal distribution. In this simulation, we still used the subgradient method to obtain the fitted model. Thenumberofiterationsused for fitting was set at 20,000. What isdifferent from Section 6.1 is the growth rates for $r_n$ and $V_n$ set in this simulation. Asmentioned in Section 5, the growth rates required for asymptotic normality are slower than those required for consistency. Therefore, we chose $r_n = n^1 1^8$ and $V_n = 10n^1 1^{10}$. Such choice satisfies the condition (Cl) inTheorem 5.1. To get the normal Q- Q plot for $-n^{112} I =: 7_1$ *lfn(Xi) - fo(xi)]*, we repeated the simulation 200 times.

Figures 2 to 4 are the normal Q- Q plots under different nonlinear functions and various sample sizes. From the figures, we found that the statistic $n^{112} E?$ *lfn(Xi) - fo(xi)]* fit the normal distribution pretty well under all simulation scenarios. It is also worth to note that the Q- Q plots looks similar under all simulation scenarios. This is what we would expect since the limiting distribution for the statistic $-n^{1/2} E?$ *lfn(Xi) - fo(xi)]* is $N(0,1)$ under all scenarios. Another implication we can obtain from the Q- Q plots is that the statistic $-n^{1/2} I: =?$ *i lf n( Xi ) - fo(xi )]* is robust to the choice *offo*. Therefore, as long as the true function *Jo* is continuous, $N(O,1)$ is a good asymptotic distribution for $-n^{1/2} I: ?= $ *[fn(xi) - Jo(xi)]*, which facilitates hypothesis testing.

Besides the Q- Q plots, we also conducted the normality tests to check whether $-n^{1/2}$ $I:=7_1$ *[fn(Xi) - fo(xi)]* follows the standard normal distribution. Specifically, we usedthe Shapiro- Wilks test and the Kolmogorov- Smim ov test to perform the normality test. Table 2 summarises the p-values for both normality tests. As we observed from Table 2, in

**Table 2.** Summary of results from the Shapiro- Wilks test and the Kolmogorov- Smirnov test. We use 'NN', 'TRI' and 'ND' to denote a neural network described in (10), a trigonometric function described in **(11)** and a continuous function having a non-differential point described in (12), respectively.

| Sample sizes | Shapiro- Wilks test | | | Kolmogorov- Smirnov test | | |
|---|---|---|---|---|---|---|
| | **NN** | TRI | ND | **NN** | TRI | ND |
| 50 | 0.878 | 0.884 | 0.881 | 0.584 | 0.597 | 0.595 |
| 100 | 0.098 | 0.095 | 0.095 | 0.472 | 0.508 | 0.484 |
| 200 | 0.940 | 0.944 | 0.944 | 0.731 | 0.719 | 0.708 |
| 300 | 0.884 | 0.888 | 0.872 | 0.976 | 0.986 | 0.973 |
| 400 | 0.514 | 0.525 | 0.513 | 0.670 | 0.754 | 0.708 |
| 500 | 0.768 | 0.778 | 0.768 | 0.733 | 0.769 | 0.733 |

all cases, we failed to reject that $n^{-1/2} \sum_i [n(x;) - J_0(x;)]$ follows the standard normal distribution.

## 7. Discussion

We have investigated the asymptotic properties, including the consistency, rate of convergence and asymptotic normality for neural network sieve estimators with one hidden layer. While in practice, the number of hidden unites is often chosen ad hoc, it is important to note that the conditions in the theorems provide theoretical guidelines on choosing the number of hidden units for a neural network with one hidden layer to achieve the desired statistical properties. The validity of the conditions made in the theorems has also been checked through simulation results. Theorems 5.1 and 5.2 depend on the knowledge of the underlying function $J_o$, which is typically unknown in practice. Therefore, if we assume $J_o$ has some certain form, the results can be applied and served as preliminary work for conducting hypothesis testing on $H_o : J_o = h_o$ for a fixed function $h_o$. On the other hand, since multiple functions can lead to the same value of $n^{-1} \sum_1 f_o(x;)$, the test may not be power. The asymptotic normality results are crucial in developing more sophisticated significance test methods for neural networks (Shen et al. 2022).

The work conducted in this paper mainly focuses on sieve estimators based on neural networks with one hidden layer and standard sigmoid activation function. The work presented in this paper can be extended in several ways. The main theorems in this paper depend heavily on the covering number or the entropy number of the function class consisting of neural network with one hidden layer. Theorem 14.5 in Anthony and Bartlett (2009) provides a general upper bound for the covering number of a function class consisting of deep neural networks with Lipchitz continuous activation functions. Therefore, it is possible to extend our results discussed in this paper to a deep neural network with Lipchitz continuous activation functions. It is also worthwhile to investigate asymptotic properties of other commonly used deep learning models such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs).

On the other hand, although homoscedasticity is assumed in the previous discussions, it can be relaxed to take heteroscedasticity into consideration To see this, if $\mathbb{E}[\varepsilon] = \phi^2 \sigma^2(x;)$, then under the assumptions that $0 < \sigma < S$ for some $> 0$ and $L$ $\mathbb{E}[1/\sigma^2 <$ oo, the proof of Lemma 3.1 can go through. The only modifications to be made are to use Kolmogorov Strong Law of Large Numbers to show (I) $0$ and to change a $^2$ to $\sigma^2 \phi^2$

later on. Therefore, the consistency result stillholds under heteroskedasticitywith the aforementioned two assumptions satisfied. Moreover, after a clear examination on the proof of Theorems 4.1, 5.1 and 5.2, it is easy to see that only the consistency part is involved with heteroskedasticity. Therefore, these results still hold under the aforementioned two assumptions.

When we train a deep neural network, we usually need to face an overfitting issue. In practice, regularisation is frequently used to reduce overfitting. Another natural extension of the work discussed in this paper is to modify the loss function by involving some regularisation terms. By taking regularisation into account, we believe the theories could have a much broader application in real-world scenarios.

## Acknowledgments

## Disclosure statement

No potential conflict of interest was reported by the author{s).

## Funding

## ORCID

*Qing Lu* G http://orcid.org/0000-0002-7943-966X

## References

Alexander, K.S. {1984 ), 'Probability Inequalities for Empirical Processes and a Law of the Iterated Logarithm', *The Annals of Probability,* 12(4), 1041- 1067.

Anthony, M., and Bartlett, P.L. {2009 ), *Neural Network Learning: Theoretical Foundations,* Cambridge University Press. https://www.stat.berkeley.edu/ bartlett/ nnl/index .html

Barron, A.R. {1994), 'Approximation and Estimation Bounds for Artificial Neural Networks', *Machine Learning,* 14(1), 115- 133.

Bauer, B., and Kohler, M.{2019), 'On Deep Learning As a Remedy for the Curse of Dimensionality in Nonparametric Regression', *The Annals of Statistics,* 47(4), 2261- 2285.

Boyd, S., and Mutapcic, A. (2008), 'Subgradient Methods (Notes for EE364B Winter 2006-07, Stanford University)'.

Chen, X. (2007), 'Large Sample Sieve Estimation of Semi-nonparametric Models', *Handbook of Econometrics,* 6, 5549- 5632.

Chen, X., Racine, J., and Swanson, N.R. (2001), 'Semiparametric Arx Neural-network Models with an Application to Forecasting Inflation', *IEEE Transactions on Neural Networks,* 12(4), 674- 683.

Chen, X., and Shen, X. (1998), 'Sieve Extremum Estimates for Weakly Dependent Data', *Econometrica,* 66(2), 289- 314.

Chen, X., and White, H. {1999), 'Improved Rates and Asymptotic Normality for Nonparametric Neural Network Estimators', *IEEE Transactions on Information Theory,* 45(2), 682-691.

Devroye, L., Gyorfi, L., and Lugosi, G. (2013), *A Probabilistic Theory of Pattern Recognition,* Vol. 31. New York, NY: Springer Science & Business Media.

Fabozzi, F.J., Fallahgoul, H., Franstianto, V., and Loeper, G. (2019), 'Towards Explaining Deep Learning: Asymptotic Properties of relu ffn Sieve Estimators', Available at SSRN 3499324.

Farrell, **M.H .,** Liang, T., and Misra, S. (2020), 'Deep Learning for Individual Heterogeneity: an Automatic Inference Framework', Preprint arXiv:2010.14694

Farrell, M.H., Liang, T., and Misra, S. (2021), 'Deep Neural Networks for Estimation and Inference', *Econometrica,* 89(1), 181- 213.

Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016 ), *Deep Learning,* Vol. 1, MIT Press Cambridge. https://www.deeplearningbook.org/

Grenander, U. (1981), *Abstract Inference,* New York: Wily.

Gyorfi, L., Kohler, M., Krzyzak, A., and Walk, H. (2002), *A Distribution-FreeTheory of Nonparametric Regression,* Vol. 1. New York, NY: Springer .

Hore!, E., and Giesecke., K. (2020), 'Significan ce Tests for Neural Networks', *Journal of Machine Learning Research,* 21(227), 1- 29.

Hornik, K., Stinchcombe, M., and White, H. (1989), 'Multilayer Feedforward Networks Are Universal Approximators', *Neural Networks,* 2(5), 359- 366.

Kohler, M., and Langer, S. (2021), 'On the Rate of Convergence of Fully Connected Deep Neural Network Regression Estimates', *The Annals of Statistics,* 49(4), 2231- 2249.

Makovoz, Y. (1996), 'Random Approximants and Neural Networks', *Journal of Approximation Theory,* 85(1), 98- 109.

Mendelson, S. (2003), 'A few notes on statistical learning theory', in *Advanced Lectures on Machine Learning,* Springer, pp. 1- 40. https://link.springer.com/chapter/10.1007/3-540-36434-X_1

Schmidt-Hieber, J. (2020), 'Nonparametric Regression Using Deep Neural Networks with Relu Activation Function', *The Annals of Statistics,* 48(4), 1875- 1897.

Shen, X. (1997), 'On Methods of Sieves and Penalization', *The Annals of Statistics,* 25(6), 2555- 2591.

Shen, X., Jiang, C., Sakhanenko, L., and Lu, Q. (2021), 'A Goodness-of-fit Test Based on Neural Network Sieve Estimators', *Statistics & Probability Letters,* 174, 109100.

Shen, X., Jiang, C., Sakhanenko, L., and Lu, Q. (2022), 'A Sieve Quasi-likelihood Ratio Test for Neural Networks with Applications to Genetic Association Studies', Preprint arXiv:2212.08255.

Shen, X., and Wong, W H. (1994), 'Convergence Rate of Sieve Estimates', *The Annals of Statistics,* 580 - 615. htt ps:// www.jstor .org/ stable/2242281

Stone, C.J. (1982), 'Optimal Global Rates of Convergence for Nonparametric Regression', *TheAnnals of Statistics,* 1040 - 1053. http s://www .jstor.o rg/stable/2240707

van de Geer, S. (2000), *Empirical Processes in M-estimation,* Vol. 6, Cambridge University Press. https://www.cambridge.org/ us/academic/ subjects/statistics-probability/statistical-theory-and - methods/empirical-processes-mestimation?format=PB&isbn=9780521123259

van der Vaart, A.W., and Wellner, J.A. (1996), *Weak Convergence and Empirical Processes.* New York: Springer.

Vapnik, V. (1998 ), *Statistical Learning Theory,* Vol. 3, New York: Wiley.

White, H. (1989), 'Learning in Artificial Neural Networks: a Statistical Perspective', *Neural Compu - tation,* 1(4), 425- 464.

White, H. (1990), 'Connectionist Nonparametric Regression: Multilayer Feedforward Networks Can Learn Arbitrary Mappings', *Neural Networks,* 3(5), 535- 549.

White, H., and Wooldridge, J. ( 1991 ), 'Some results on sieve estimation with dependent observations', in *Nonparametric and SemiparametricMethods in Economics,* ed. W Barnett, J. Powell, and G. Tauchen, New York: Cambridge University Press, pp. 459- 493.