

# Towards Explainable AI: Significance Tests for Neural Networks

**Kay Giesecke**

*Advanced Financial Technologies Laboratory*  
Stanford University

`people.stanford.edu/giesecke/  
fintech.stanford.edu`

Joint work with Enguerrand Horel (Stanford)

- Neural networks underpin many of the best-performing AI systems, including speech recognizers on smartphones or Google's latest automatic translator
- The tremendous success of these applications has spurred the interest in applying neural networks in a variety of other fields, including finance, economics, operations, marketing, medicine, and many others
- In finance, researchers have developed several promising applications in risk management, asset pricing, and investment management

- First wave: single-layer nets
  - Financial time series: White (1989), Kuan & White (1994)
  - Nonlinearity testing: Lee, White & Granger (1993)
  - Economic forecasting: Swanson & White (1997)
  - Stock market prediction: Brown, Goetzmann & Kumar (1998)
  - Pricing kernel modeling: Bansal & Viswanathan (1993)
  - Option pricing: Hutchinson, Lo & Poggio (1994)
  - Credit scoring: Desai, Crook & Overstreet (1996)
- Second wave: multi-layer nets (deep learning)
  - Mortgages: Sirignano, Sathwani & Giesecke (2016)
  - Order books: Sirignano (2016), Cont and Sirignano (2018)
  - Portfolio selection: Heaton, Polson & Witte (2016)
  - Returns: Chen, Pelger & Zhu (2018), Gu, Kelly & Xiu (2018)
  - Hedging: Halperin (2018), Bühler, Gonon & Teichmann (2018)
  - Optimal stopping: Becker, Cheridito & Jentzen (2018)
  - Treasury markets: Filipovic, Giesecke, Pelger, Ye (2019)
  - Real estate: Giesecke, Ohlrogge, Ramos & Wei (2019)
  - Insurance: Wüthrich and Merz (2019)

- The success of NMs is largely due to their amazing approximation properties, superior predictive performance, and their scalability
- A major caveat however is **model explainability**: NMs are perceived as black boxes that permit little insight into how predictions are being made
- Key inference questions are difficult to answer
  - Which input variables are statistically significant?
  - If significant, how can a variable's impact be measured?
  - What's the relative importance of the different variables?

This issue is not just academic; it has slowed the implementation of NIS in financial practice, where regulators and other stakeholders often insist on model explainability

- Credit and insurance underwriting (regulated)
- Transparency of underwriting decisions
- Investment management (unregulated)
- Transparency of portfolio designs
- Economic rationale of trading decisions

- We develop a **pivotal test** to assess the statistical significance of the input variables of a NN
  - Focus on single-layer feedforward networks
  - Focus on regression setting
- We propose a **gradient-based test statistic** and study its asymptotics using nonparametric techniques
  - Asymptotic distribution is a mixture of  $\chi^2$  laws
- The test enables one to address key inference issues:
  - Assess statistical significance of variables
  - Measure the impact of variables
  - Rank order variables according to their influence
- Simulation and empirical experiments illustrate the test

- Regression model  $Y = f_0(X) + \epsilon$
- $X \in \mathcal{X} \subset \mathbb{R}^d$  is a vector of  $d$  feature variables with law  $P$
- $f_0 : \mathcal{X} \rightarrow \mathbb{R}$  is an unknown deterministic  $C^1$ -function
- $\epsilon$  is an error variable:  $\epsilon \perp\!\!\!\perp X, \mathbb{E}(\epsilon) = 0, \mathbb{E}(\epsilon^2) = \sigma^2 < \infty$
- To assess the significance of a variable  $X_j$ , we consider sensitivity-based hypotheses:

$$H_0 : \lambda_j := \int_{\mathcal{X}} \left( \frac{\partial f_0(x)}{\partial x_j} \right)^2 p_{\mu}(x) dx = 0$$

$$H_A : \lambda_j \neq 0$$

Here,  $\mu$  is a positive weight measure

- A typical choice is  $\mu = P$  and then  $\lambda_j = \mathbb{E} \left[ \left( \frac{\partial f_0(X)}{\partial x_j} \right)^2 \right]$

- Suppose the function  $f_0$  is linear (multiple linear regression)

$$f_0(x) = \sum_{k=1}^p \beta_k x_k$$

Then  $\lambda_j \propto \beta_j^2$ , the squared regression coefficient for  $X_j$ , and the null takes the form  $H_0 : \beta_j = 0$  ( $\leftarrow$  t-test)

- In the general nonlinear case, the derivative  $\frac{\partial f_0(x)}{\partial x_j}$  depends on  $x$ , and  $\lambda_j = \int \left( \frac{\partial f_0(x)}{\partial x_j} \right)^2 p(x)$  is a weighted average



- We study the case where the unknown regression function  $f_0$  is modeled by a single-layer feedforward NN

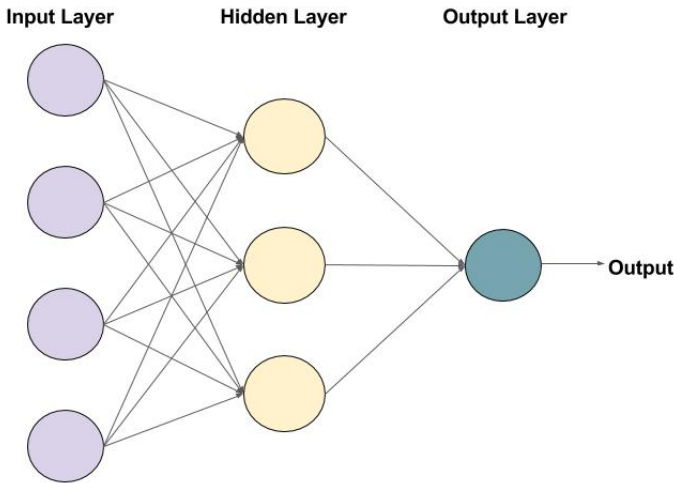
- A **single-layer NN**  $f$  is specified by a bounded activation function  $\psi$  on  $\mathbb{R}$  and the number of hidden units  $K$ :

$$f(x) = b_0 + \sum_{k=1}^K b_k \psi(a_{0,k} + a_{k,\perp} x)$$

where  $b_0, b_k, a_{0,k} \in \mathbb{R}$  and  $a_k \in \mathbb{R}^d$  are to be estimated

- Functions of the form  $f$  are dense in  $C(\mathcal{X})$  (they are *universal approximators*): choosing  $K$  large enough,  $f$  can approximate  $f_0$  to any given precision

Neural network:  $d$  4 features,  $K$  3 hidden units



- We use  $n$  i.i.d. samples  $(Y_i, X_i)$  to construct a sieve M-estimator  $f_n$  of  $f$  for which  $K = K_n$  increases with  $n$
- We assume  $f_0 \in \Theta =$  class of  $C^1$  functions on  $d$ -hypercube  $\mathcal{X}$  with uniformly bounded Sobolev norm

- Sieve subsets  $\Theta_n \subseteq \Theta$  generated by NNS  $f$  with  $K_n$  hidden units, bounded  $L^1$  norms of weights, and sigmoid  $\psi$

- The sieve M-estimator  $f_n$  is the approximate maximizer of the empirical criterion function  $L_n(\mathcal{g}) = \frac{1}{n} \sum_{i=1}^n l(Y_i, X_i; \mathcal{g})$ , where  $l : \mathbb{R} \times \mathcal{X} \times \Theta \rightarrow \mathbb{R}$ , over  $\Theta_n$ :

$$L_n(f_n) \geq \sup_{\mathcal{g} \in \Theta_n} L_n(\mathcal{g}) - o_p(1)$$

- The NN test statistic is given by

$$\chi_n^f = \int x \left( \frac{\partial f_n(x)}{\partial x_j} \right)^2 p(x) \phi[f_n]$$

- We will use the asymptotic ( $n \rightarrow \infty$ ) distribution of  $\chi_n^f$  for testing the null since a bootstrap approach would typically be too computationally expensive

- 1 Asymptotic distribution of  $f_n$
- 2 Functional delta method

- In the large- $n$  regime, due to the universal approximation property, we are actually performing inference on the "ground truth"  $f_0$  (**model-free inference**)

Assume that

- $dP = \nu d\lambda$  for bounded and strictly positive  $\nu$
- The dimension  $K_n$  of the NN satisfies  $K_n^{2+1/d} \log K_n = O(n)$ ,
- The loss function  $l(y, x, g) = -\frac{1}{2}(y - g(x))^2$ .

Then

$$r_n(f_n - f_0) \iff h^*$$

in  $(\Theta, L^2(P))$  where

$$r_n = \left( \frac{\log n}{n} \right)^{\frac{2}{2(2d+1)}}$$

and  $h^*$  is the argmax of the Gaussian process  $\{\mathbb{G}_t^f : f \in \Theta\}$  with mean zero and  $\text{Cov}(\mathbb{G}_s, \mathbb{G}_t) = 4\sigma^2 \mathbb{H}(s(X)t(X))$ .

- $r_n$  is the **estimation rate** of the NN (Chen and Shen (1998)):

$$\mathbb{E}^X[(f_n(X) - f_0(X))^2] = O_p(r_n^{-1})$$

assuming the number of hidden units  $K_n$  is chosen such that

$$K_n^{2+1/d} \log K_n = O(n)$$

- Proof uses empirical process arguments
- Estimation rate implies tightness of  $h_n = r_n(f_n - f_0)$
- Rescaled and shifted criterion function converges weakly to Gaussian process
- Gaussian process has a unique maximum at  $h^*$
- Argmax continuous mapping theorem

## Theorem

*Under the conditions of Theorem 1 and the null hypothesis,*

$$r_n^2 \lambda_j^n \implies \int_{\mathcal{X}} \left( \frac{\partial h^*(x)}{\partial x_j} \right)^2 d\mu(x)$$

## Theorem

Assume  $\mu = P$  so that the test statistic

$$\lambda_j^n = \mathbb{E}_X \left[ \left( \frac{\partial f_n(X)}{\partial x_j} \right)^2 \right].$$

Under the conditions of Theorem 1 and the null hypothesis, the empirical test statistic satisfies

$$r_n^2 n^{-1} \sum_{i=1}^n \left( \frac{\partial f_n(X_i)}{\partial x_j} \right)^2 \implies \mathbb{E}_X \left[ \left( \frac{\partial h^*(X)}{\partial x_j} \right)^2 \right]$$



# Identifying the asymptotic distribution

## Theorem

Take  $\mu = P$ . Let  $\{\phi_i\}$  be an orthonormal basis of  $\Theta$ . If that basis is  $C^1$  and stable under differentiation, then

$$\mathbb{E}_X \left[ \left( \frac{\partial h^*(X)}{\partial x_j} \right)^2 \right] = \frac{B^2}{\sum_{i=0}^{\infty} \frac{\chi_i^2}{d_i^2}} \sum_{i=0}^{\infty} \frac{\alpha_{i,j}^2}{d_i^4} \chi_i^2,$$

where  $\{\chi_i^2\}$  are i.i.d. samples from the chi-square distribution, and where  $\alpha_{i,j} \in \mathbb{R}$  satisfies  $\frac{\partial \phi_i}{\partial x_j} = \alpha_{i,j} \phi_{k(i)}$  for some  $k : \mathbb{N} \rightarrow \mathbb{N}$ , and the  $d_i$ 's are certain functions of the  $\alpha_{i,j}$ 's.

- Truncate the infinite sum at some finite order  $N$
- Draw samples from the  $\chi^2$  distribution to construct a sample of the approximate limiting law
- Repeat  $m$  times and compute the empirical quantile  $Q_{N,m}^*$  at level  $\alpha \in (0, 1)$  of the corresponding samples
- If  $m = m_N \rightarrow \infty$  as  $N \rightarrow \infty$ , then  $Q_{N,m_N}^*$  is a consistent estimator of the true quantile of interest
- Reject  $H_0$  if  $\chi_{\alpha}^f < Q_{N,m_N}^*(1 - \alpha)$  such that the test will be asymptotically of level  $\alpha$ :

$$\mathbb{P}^{H_0}(\chi_{\alpha}^f < Q_{N,m_N}^*(1 - \alpha)) \leq \alpha$$

- 8 variables:

$$X = (X_1, \dots, X_8) \sim U(-1, 1)^8$$

- Ground truth:

$$Y = 8 + X_1^2 + X_2X_3 + \cos(X_4) + \exp(X_5X_6) + 0.1X_7 + \epsilon$$

- where  $\epsilon \sim N(0, 0.01^2)$  and  $X_8$  has no influence on  $Y$
- Training (via TensorFlow): 100,000 samples  $(Y_i, X_i)$
- Validation, Testing: 10,000 samples each
- Out-of-sample MSE:

Model	Mean Squared Error
NN with $K = 25$	$3.1 \cdot 10^{-4} \sim \text{Var}(\epsilon)$
Linear Regression	0.35

## Linear model fails to identify significant variables

Variable	coef	std err	t	$P >  t $
<b>const</b>	<b>10.2297</b>	0.002	5459.250	<b>0.000</b>
1	-0.0031	0.003	-0.964	0.335
2	0.0051	0.003	1.561	0.118
3	-0.0026	0.003	-0.800	0.424
4	0.0003	0.003	0.085	0.932
5	0.0016	0.003	0.493	0.622
6	-0.0033	0.003	-1.035	0.300
<b>7</b>	<b>0.0976</b>	0.003	30.059	<b>0.000</b>
<b>8</b>	-0.0018	0.003	-0.563	<b>0.573</b>

Only the intercept and the linear term  $0.1X_7$  are identified as significant. The irrelevant  $X_8$  is correctly identified as insignificant.

Variable	Test Statistic	Power/Size
1	1.310	1
2	0.332	1
3	0.331	1
4	0.267	1
5	0.480	1
6	0.479	1
7	$1.010 \cdot 10^{-2} (= 0.1^2)$	1
8	$4.200 \cdot 10^{-6}$	$0.13 < 0.05$

- Size: asymptotic distribution tends to underestimate the variance of the finite sample distribution of the test statistic
- Efficiency: gradient (TensorFlow), no re-fitting required
- Robustness: insensitive to correlated feature data

## Application: House price valuation

- **Data:** 120+ million housing sales from county registrar of deed offices across the US (source: CoreLogic)
- **Sample period:** 1970 to 2017
- **Geographical area:** Merced County, CA; 76,247 samples
- **Prediction of  $Y$**  = log sale price
- **Variables  $X$  ( $d = 68$ ):** Bedrooms, Full\_Baths, Last\_Sale\_Amount, N\_Originations, N\_Past\_Sales, Sale\_Month, SqFt, Stories, Tax\_Amount, Time\_Since\_Prior\_Sale, etc.
- Training and gradients via TensorFlow, Adam
- Validation (70/20/10 split):  $K = 150$  nodes,  $L_1$  weight  $10^{-5}$
- Test MSE is 0.45

# Application: House price valuation



## Top 10 significant (5%) variables (out of 68)

<b>Variable Name</b>	<b>Test Statistic</b>
Last_Sale_Amount	1.640
Tax_Land_Square_Footage	1.615
Sale_Month_No	1.340
Tax_Amount	0.383
Last_Mortgage_Amount	0.104
Tax_Assd_Total_Value	0.081
Tax_Improvement_Value_Calc	0.072
Tax_Land_Value_Calc	0.069
Year_Built	0.068
SqFt	0.056
...	...



- We develop a computationally efficient, pivotal significance test for neural networks
  - Assess the impact of feature variables on the prediction
  - Rank variables according to their predictive importance
- This opens up a broader range of applications of NNs in financial practice
- Ongoing work
  - Treatment of NN classifiers and deep networks
  - Cross derivatives for testing interactions between variables
  - Alternative approaches

- Suppose the elements of  $X$  are i.i.d. uniform on  $[-1, 1]$

- Using the Fourier basis, the limiting distribution takes the form

$$B^2 \frac{\sum_{n \in \mathbb{N}^d} \chi_n^{\frac{d}{2}}}{\sum_{n \in \mathbb{N}^d} \frac{d_n^{\frac{d}{2}}}{n_j^2 \pi^2} \chi_n^{\frac{d}{2}}}$$

- $n = (n_1, n_2, \dots, n_j, \dots, n_d)$
- $d_n^{\frac{d}{2}} = \sum_{|\alpha| \leq \lfloor \frac{n}{2} \rfloor + 2} \prod_{k=1}^d (n_k \pi)^{2\alpha_k}$
- $\{\chi_n^{\frac{d}{2}}\}_{n \in \mathbb{N}^d}$  are i.i.d. chi-square variables

- We note that  $\Theta$  is a subspace of the Hilbert space  $L^2(P)$  which admits an orthonormal basis  $\{\phi_i\}_{i=0}^{\infty}$
- If this basis is  $C^1$  and stable under differentiation, i.e. if there are a real  $\alpha_{ij}$  and a mapping  $k : \mathbb{N} \rightarrow \mathbb{N}$  such that

$$\frac{\partial \phi_i}{\partial x_j} = \alpha_{ij} \phi_{k(i)},$$

then there exists an invertible operator  $D$  such that

$$\|f\|_{k,2} = \|Df\|_{2} = \sum_{i=0}^{\infty} d_i^2 \langle f, \phi_i \rangle_{L^2(P)}$$

where the  $d_i$ 's are certain functions of the  $\alpha_{ij}$ 's