

Randomized function approximation

John Darges

**Department of Mathematics
North Carolina State University**

December 4, 2023

Learning vs. Approximation

Suppose we have a data set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, $(\mathbf{x}_i, y_i) \in \mathbb{R}^{d \times 1}$

Approximation

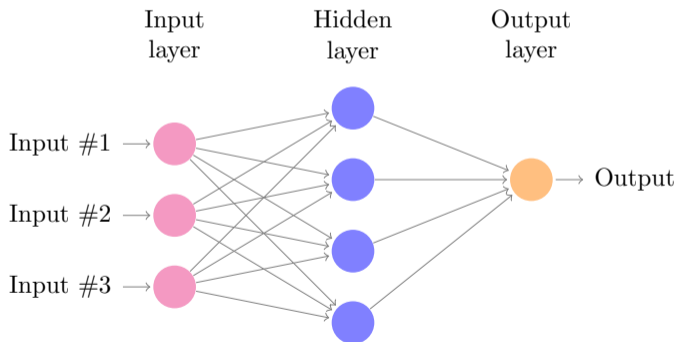
- Data comes from a known model
 $y_i = F(\mathbf{x}_i)$
- Inputs follow a known distribution $\mathbf{x} \sim \mathcal{K}$
- Can choose what is in our data set, but data may be expensive to generate

Learning

- Model is unknown but inputs/outputs are labeled
- Data follows some unknown distribution
 $(\mathbf{x}, y) \sim \mathcal{D}$
- Have access to some set of data (which is i.i.d.)

Feedforward neural networks

Artificial neural networks (ANNs) have very general structure (we focus on single layer neural networks)



Feedforward neural networks

Single layer NN has the form

$$F_{\text{NN}}(\mathbf{x}) = \sum_{k=1}^M \alpha_k \sigma(W_j^\top \mathbf{x} + b_j) = \boldsymbol{\alpha}^\top \sigma(\mathbf{W}^\top \mathbf{x} + \mathbf{b}) \quad (1)$$

- $\mathbf{W} \in \mathbb{R}^{d \times M}$ is hidden layer weight matrix, $\mathbf{W} = [W_1 \ \dots \ W_M]$
- $\boldsymbol{\alpha} \in \mathbb{R}^M$ outer weight vector
- $\mathbf{b} \in \mathbb{R}^M$ bias vector
- $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ activation function (with universal approximation property¹)

To train, solve nonlinear least squares problem

$$\min_{\mathbf{W}, \boldsymbol{\alpha}, \mathbf{b}} \sum_{i=1}^N (F_{\text{NN}}(\mathbf{x}_i; \mathbf{W}, \boldsymbol{\alpha}, \mathbf{b}) - y_i)^2 \quad (2)$$

¹M. Leshno and V. Ya Lin and A. Pinkus and S. Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. 1993.

Kernels

- Kernels are functions $K : \mathbf{X} \times \mathbf{X} \rightarrow \mathbb{R}$ (we will consider $\mathbf{X} = \mathbb{R}^d$)
- Some kernels are induced by a feature map $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$
- \mathcal{H} is a Hilbert space of certain functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- Induced kernel is

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi_{\mathbf{x}}, \phi_{\mathbf{x}'} \rangle_{\mathcal{H}} \quad (3)$$

- Feature space \mathcal{H} is a Reproducing kernel Hilbert space

$$f(\mathbf{x}) = \langle f, K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}, \quad f \in \mathcal{H}, \mathbf{x} \in \mathbb{R}^d \quad (4)$$

- Under right conditions, some function in RKHS can model our data

Kernel ridge regression

- Find $g \in \mathcal{H}$ so that $\langle g, \phi_{\mathbf{x}_i} \rangle_{\mathcal{H}} \approx y_i$, have a linear least squares problem

$$\min \sum_{i=1}^N (\langle g, \phi_{\mathbf{x}_i} \rangle_{\mathcal{H}} - y_i)^2 = \sum_{i=1}^N (\langle f, K(\cdot, \mathbf{x}_i) \rangle_{\mathcal{H}} - y_i)^2 \quad (5)$$

- Construct kernel matrix $\mathbf{K} \in \mathbb{R}^{d \times d}$ where $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$, which is SPD
- Kernel ridge regression solves regularized least squares problem

$$\min_{\boldsymbol{\alpha}} \sum_{i=1}^N \left(\sum_{j=1}^N \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) - y_i \right)^2 + \frac{\lambda^2}{2} \|\boldsymbol{\alpha}\|^2 \quad (6)$$

- Solution: $\boldsymbol{\alpha} = (\mathbf{K} + \lambda^2 I)^{-1} \mathbf{y}$ gives kernel machine

$$F_{\text{KRR}}(\mathbf{x}) = \sum_{j=1}^N \alpha_j K(\mathbf{x}, \mathbf{x}_j) = \boldsymbol{\alpha}^{\top} K(\mathbf{x}) \quad (7)$$

Connection between Kernels and ANN

- Consider a feature map induced by an activation function

$$\phi_{\omega}(\mathbf{x}) = \sigma(\omega^{\top} \mathbf{x} + b) = \sigma(\omega^{\top}(\mathbf{x}, 1)) \quad (8)$$

- We focus on Hilbert spaces with the L^2 inner product
- Functions in the RKHS look like

$$F(\mathbf{x}) = \int \alpha(\omega) \sigma(\omega^{\top}(\mathbf{x}, 1)) d\omega \quad (9)$$

Randomization

- Randomized algorithms have become prevalent in numerical linear algebra
- Improve efficiency of smaller problems and feasibility of large problems
- Many ways to introduce to introduce randomness and still guarantee good results (almost surely)

Random weight neural networks

- Recall single layer NN

$$F_{\text{NN}}(\mathbf{x}) = \sum_{k=1}^M \alpha_k \sigma(W_j^\top \mathbf{x} + b_j) = \boldsymbol{\alpha}^\top \sigma(\mathbf{W}^\top \mathbf{x} + \mathbf{b}) \quad (10)$$

- Instead of training over all $(d + 2)M$ parameters, just randomly sample hidden layer weights and biases
- Only need to optimize output weights

$$\min_{\boldsymbol{\alpha}} \sum_{i=1}^N (F_{\text{NN}}(\mathbf{x}_i; \mathbf{W}', \boldsymbol{\alpha}, \mathbf{b}') - y_i)^2 = \|\mathbf{H}\boldsymbol{\alpha} - \mathbf{y}\|^2 \quad (11)$$


- Here $\mathbf{H}_{ij} = \sigma(W_j^\top \mathbf{x}_i + b_j)$
- Choices of activation function and sampling distribution matter!

History of random weight NNs

- Broomhead first introduces idea of turning training to linear least squares problem²
- Schmidt neural networks³
- Barron gives $\mathcal{O}(1/n^{2/d})$ convergence rates for sigmoidal networks⁴

²D.S. Broomhead and D. Lowe. Multivariable Functional Interpolation and Adaptive Networks. 1988.

³W.F. Schmidt, M.A. Kraaijveld, R.P.W. Duin. Feedforward neural networks with random weights. 1992.

⁴A. Barron. Universal Approximation Bounds for Superpositions of a Sigmoidal Function. 1993. 

Random vector functional link

- Random vector functional link (RVFL)⁵
Approximates functions with compact support, weights sampled uniformly
Average asymptotic convergence and generalization bounds⁶
- Corrected theorems given in⁷

⁵Y.-H. Pao, G.-H. Park, D. Sobajic. Learning and generalization characteristics of the random vector functional-link net. 1994.

⁶B. Igel'nik, and Y.-H. Pao. Stochastic choice of basis functions in adaptive function approximation and the functional-link net. 1995.

⁷D. Needell, A. Nelson, R. Saab, P. Salanevich. Random Vector Functional Link Networks for Function Approximation on Manifolds. 2022.

Extreme learning machine

- Extreme learning machine (ELM) includes neural network and radial basis function versions⁸
- Prove universal approximation when activation function is bounded and sampling distribution is continuous⁹
- Claim much more broad approximation capabilities¹⁰
- No convergence or generalization guarantees

⁸G.-B. Huang and D. Wang and Y. Lan. Extreme learning machines: A review. 2011.

⁹G.-B. Huang, L. Chen, and C.-K. Siew. Universal Approximation Using Incremental Constructive Feedforward Networks with Random Hidden Nodes. 2006.

¹⁰G.-B. Huang, L. Chen. Convex incremental extreme learning machine. 2007

Kernel methods the random feature way

- Consider a kernel $K(\mathbf{x}, \mathbf{x}') = \int \phi_{\mathbf{x}}(\boldsymbol{\omega})\phi_{\mathbf{x}'}(\boldsymbol{\omega})d\boldsymbol{\omega}$
- For a large data set, impractical to compute (and store) full kernel matrix \mathbf{K} .
- Instead approximate kernel matrix by a rank one approximation
- Random features¹¹ use Monte Carlo sampling

$$\mathbf{K} \approx \sum_{k=1}^K \mathbf{z}_k \mathbf{z}_k^\top, \quad \mathbf{z}_k = [\phi_{\mathbf{x}_1}(\boldsymbol{\omega}_k) \quad \dots \quad \phi_{\mathbf{x}_N}(\boldsymbol{\omega}_k)]^\top \quad (12)$$

- $\mathcal{O}(\sqrt{n} \log(n))$ features give $\mathcal{O}(1/\sqrt{n})$ bounds¹²

¹¹A. Rahimi, B. Recht. Random Features for Large-Scale Kernel Machines. 2007.

¹²A. Rudi, L. Rosasco. Generalization Properties of Learning with Random Features. 2017.

Random bases

- Random basis expansion¹³ does not work with the kernel, instead the feature map
- With neural network style feature maps, they are equivalent to random weight neural networks
- Can take advantage of RKHS theory and functional analysis¹⁴
- RKHS should be dense in space of continuous functions
- We should be able to express the following by a series

$$\int f(\omega)\phi_x(\omega)d\omega \quad (13)$$

- Can recreate claim of Huang '06 using function analysis¹⁵

¹³A. Rahimi, B. Recht. Uniform approximation of functions with random bases. 2008.

¹⁴F. Bach. On the Equivalence between Kernel Quadrature Rules and Random Feature Expansions. 2017.

¹⁵Y. Sun, A. Gilbert, A. Tewari. On the Approximation Properties of Random ReLU Features. 2019. 

Random bases with structure

- Do random bases/neural networks have any advantages over random features?
- We have broad choices for activation functions
- Sampling distribution for ω has many choices, too
- By clever sampling, can impose function structure (interactions/main effects) on

$$F(\mathbf{x}) = \sum_{k=1}^M \alpha_k \sigma(W_j^\top \mathbf{x} + b_j) \quad (14)$$

- In¹⁶ and¹⁷ use sparse sampling ($W \sim X \cdot Y$, X continuous RV and Y Bernoulli RV) to impose structure

¹⁶A. Hashemi, H. Schaeffer, R. Shi, U. Topcu, G. Tran, R. Ward. Generalization bounds for sparse random feature expansions. 2023.

¹⁷J. Darges, A. Alexanderian, P. Gremaud. Extreme learning machines for variance-based global sensitivity analysis. 2023.

References

- 1 M. Leshno and V. Ya Lin and A. Pinkus and S. Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. 1993.
- 2 C.E. Rasmussen, C.K.I. Williams. Gaussian Processes for Machine Learning. 2006.
- 3 R. Schaback, H. Wendland. Kernel techniques: From machine learning to meshless methods. 2006.
- 4 P. Martinsson, J.A. Tropp. Randomized numerical linear algebra: Foundations and algorithms 2020.
- 5 D.S. Broomhead and D. Lowe. Multivariable Functional Interpolation and Adaptive Networks. 1988.
- 6 W.F. Schmidt, M.A. Kraaijveld, R.P.W. Duin. Feedforward neural networks with random weights. 1992.
- 7 A. Barron. Universal Approximation Bounds for Superpositions of a Sigmoidal Function. 1993.

References

- 1 Y.-H. Pao, G.-H. Park, D. Sobajic. Learning and generalization characteristics of the random vector functional-link net. 1994.
- 2 B. IgelNIK, and Y.-H. Pao. Stochastic choice of basis functions in adaptive function approximation and the functional-link net. 1995.
- 3 D. Needell, A. Nelson, R. Saab, P. Salanevich. Random Vector Functional Link Networks for Function Approximation on Manifolds. 2022.
- 4 G.-B. Huang, L. Chen, and C.-K. Siew. Universal Approximation Using Incremental Constructive Feedforward Networks with Random Hidden Nodes. 2006.
- 5 G.-B. Huang and D. Wang and Y. Lan. Extreme learning machines: A review. 2011.
- 6 A. Rahimi, B. Recht. Random Features for Large-Scale Kernel Machines. 2007.
- 7 A. Rahimi, B. Recht. Uniform approximation of functions with random bases. 2008.
- 8 F. Bach. On the Equivalence between Kernel Quadrature Rules and Random Feature Expansions. 2017.

References

- 1 Y. Sun, A. Gilbert, A. Tewari. On the Approximation Properties of Random ReLU Features. 2019.
- 2 A. Rudi, L. Rosasco. Generalization Properties of Learning with Random Features. 2017.
- 3 A. Hashemi, H. Schaeffer, R. Shi, U. Topcu, G. Tran, R. Ward. Generalization bounds for sparse random feature expansions. 2023.
- 4 J. Darges, A. Alexanderian, P. Gremaud. Extreme learning machines for variance-based global sensitivity analysis. 2023.
- 5 F. Liu and X. Huang and Y. Chen, J. K. Suykens. Random Features for Kernel Approximation: A Survey on Algorithms, Theory, and Beyond. 2022.
- 6 M. Nguyen, N. Mücke. Random feature approximation for general spectral methods. 2023.
- 7 P.N. Suganthan, R. Katuwal. On the origins of randomization-based feedforward neural networks. 2021.
- 8 W. Cao, X. Wang, Z. Ming, J. Gao. A review on neural networks with random weights. 2018.