

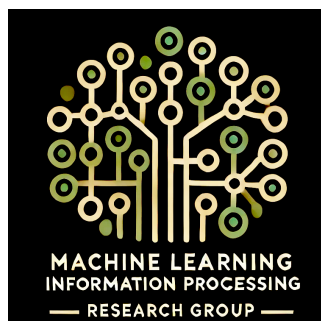


ICLR
International Conference On
Learning Representations

Asymptotic Analysis of Two-Layer Neural Networks after One Gradient Step under Gaussian Mixtures Data with Structure

Samet Demir & Zafer Dogan

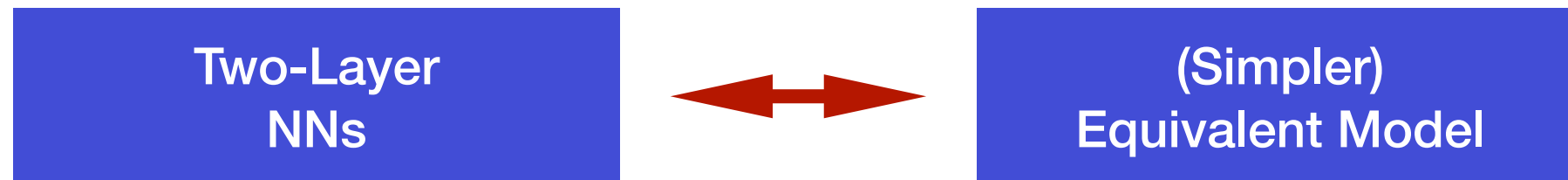
MLIP Research Group, KUIS AI Center, Koç University



**KOÇ
UNIVERSITY**

Theory of Neural Networks via Equivalent Models

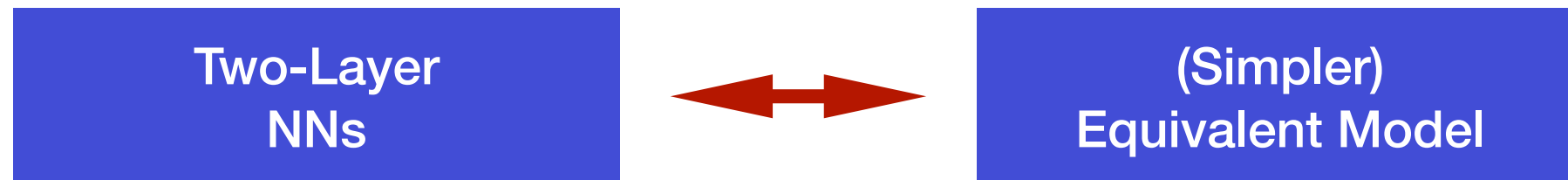
- We focus on the analysis of generalization performance of two-layer neural networks (NNs) for supervised learning via equivalent models.



Generalization (and training) errors of the two models are equivalent.

Theory of Neural Networks via Equivalent Models

- We focus on the analysis of generalization performance of two-layer neural networks (NNs) for supervised learning via equivalent models.

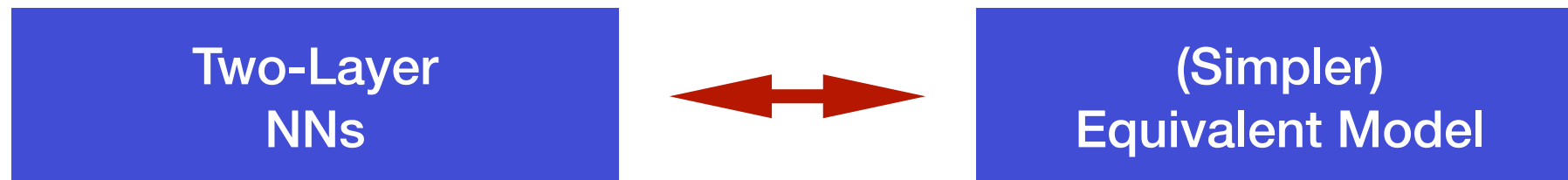


Generalization (and training) errors of the two models are equivalent.

- We consider the asymptotically proportional limit:
 - input dimension, number of neurons, and sample size diverge with finite ratios

Theory of Neural Networks via Equivalent Models

- We focus on the analysis of generalization performance of two-layer neural networks (NNs) for supervised learning via equivalent models.



Generalization (and training) errors of the two models are equivalent.

- We consider the asymptotically proportional limit:
 - input dimension, number of neurons, and sample size diverge with finite ratios
- Limitations of the literature
 - Lack of feature learning (e.g., random feature model) [1]
 - Limited data assumption (e.g., standard Gaussian inputs) [1,2,3,4]

[1] Hu and Lu. Universality laws for high-dimensional learning with random features. IEEE Trans. Inf. Theory, 69(3):1932–1964, Mar. 2023.

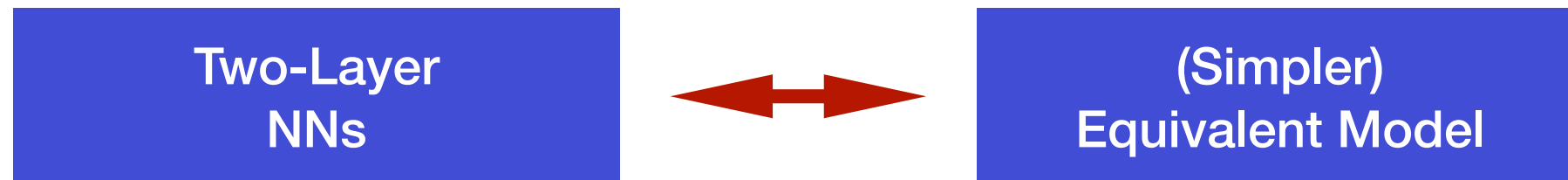
[2] Dandi et al. How two-layer neural networks learn, one (giant) step at a time, arXiv:2305.18270, 2023.

[3] Cui et al. Asymptotics of feature learning in two-layer networks after one gradient-step. ICML 2024.

[4] Moniri et al. A theory of non-linear feature learning with one gradient step in two-layer neural networks. ICML 2024.

Theory of Neural Networks via Equivalent Models

- We focus on the analysis of generalization performance of two-layer neural networks (NNs) for supervised learning via equivalent models.



Generalization (and training) errors of the two models are equivalent.

- We consider the asymptotically proportional limit:
 - input dimension, number of neurons, and sample size diverge with finite ratios
- Limitations of the literature
 - Lack of feature learning (e.g., random feature model) [1]
 - Limited data assumption (e.g., standard Gaussian inputs) [1,2,3,4]

[1] Hu and Lu. Universality laws for high-dimensional learning with random features. IEEE Trans. Inf. Theory, 69(3):1932–1964, Mar. 2023.

[2] Dandi et al. How two-layer neural networks learn, one (giant) step at a time, arXiv:2305.18270, 2023.

[3] Cui et al. Asymptotics of feature learning in two-layer networks after one gradient-step. ICML 2024.

[4] Moniri et al. A theory of non-linear feature learning with one gradient step in two-layer neural networks. ICML 2024.

Data: Gaussian Mixtures with Structured Covariance

- Our goal: a realistic data assumption

Data: Gaussian Mixtures with Structured Covariance

- Our goal: a realistic data assumption
- Gaussian Mixtures Data
 - Motivation: Real-world data can be better modeled with mixtures instead of single Gaussian.
 - e.g., classification problems

$$\textbf{Input: } x \sim \sum_{j=1}^c \rho_j \mathcal{N}(\mu_j, \Sigma_j), \quad \textbf{Label: } y := \sigma_* (\xi^T x, c)$$

Data: Gaussian Mixtures with Structured Covariance

- Our goal: a realistic data assumption
- Gaussian Mixtures Data
 - Motivation: Real-world data can be better modeled with mixtures instead of single Gaussian.
 - e.g., classification problems

$$\textbf{Input: } x \sim \sum_{j=1}^c \rho_j \mathcal{N}(\mu_j, \Sigma_j), \quad \textbf{Label: } y := \sigma_* (\xi^T x, c)$$

- Structured Covariance
 - Motivation: Real-world data includes low-dimensional structure.
 - e.g., MNIST and CIFAR10 have intrinsic dimensions of approx. 15 and 35, resp. [5]

$$\Sigma_c = \mathbf{I}_n + \sum_{i=1}^{d_c} \theta_{c,i} \gamma_{c,i} \gamma_{c,i}^T$$

Structure (spikes)

Feature Learning by One Gradient Step

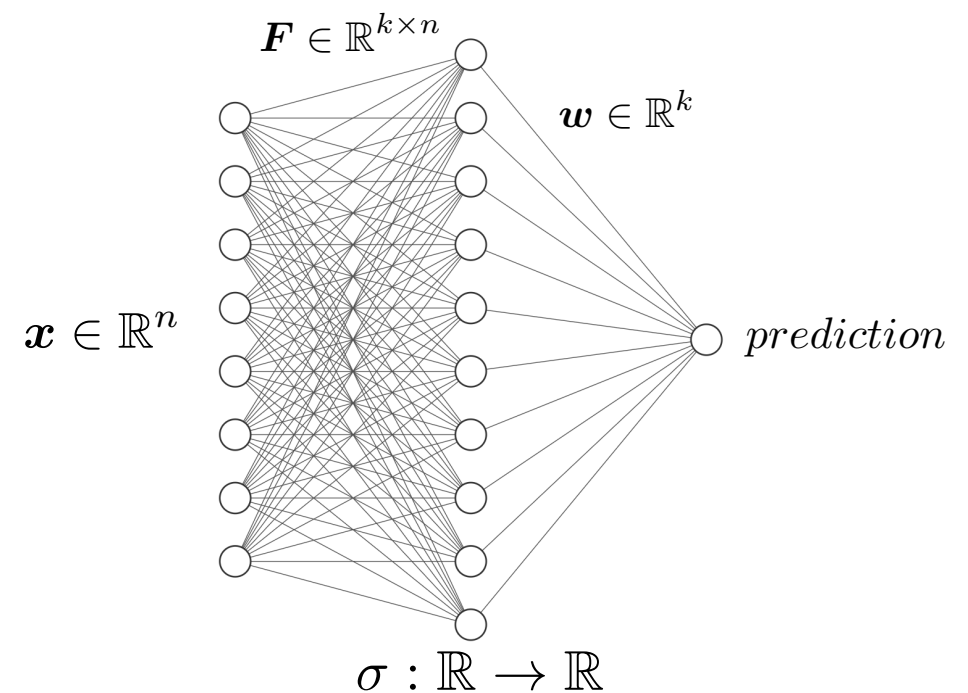
- Consider two-layer NNs for the described learning problem
- To simplify the analysis, train the model with two stages [6]

- i) Training of the first layer:
 - One step gradient descent

$$\hat{\mathbf{F}} := \mathbf{F} + \eta \mathbf{G}$$

learning rate: $\eta > 0$

gradient matrix: \mathbf{G}



Feature Learning by One Gradient Step

- Consider two-layer NNs for the described learning problem
- To simplify the analysis, train the model with two stages [6]

- i) Training of the first layer:

- One step gradient descent

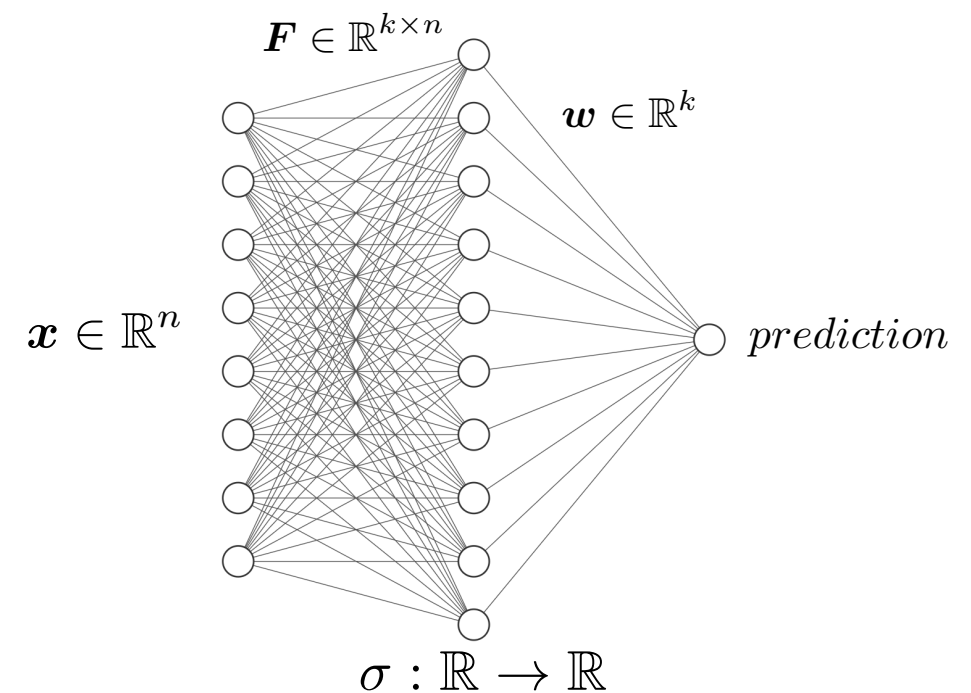
$$\hat{\mathbf{F}} := \mathbf{F} + \eta \mathbf{G}$$

learning rate: $\eta > 0$

gradient matrix: \mathbf{G}

- ii) Training of the second layer:

- Ridge regression
- $$\hat{\mathbf{w}} := \arg \min_{\mathbf{w} \in \mathbb{R}^k} \frac{1}{2m} \sum_{i=1}^m \left(y_i - \frac{1}{\sqrt{k}} \mathbf{w}^T \sigma(\hat{\mathbf{F}} \mathbf{x}_i) \right)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$



Feature Learning by One Gradient Step

- Consider two-layer NNs for the described learning problem
- To simplify the analysis, train the model with two stages [6]

- i) Training of the first layer:

- One step gradient descent

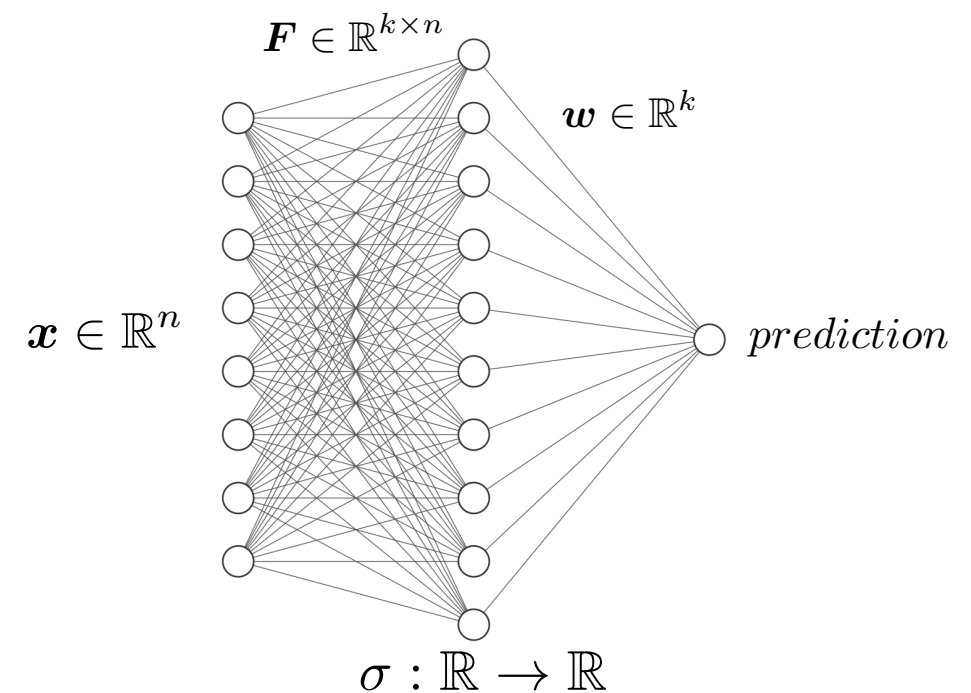
$$\hat{\mathbf{F}} := \mathbf{F} + \eta \mathbf{G}$$

learning rate: $\eta > 0$

gradient matrix: \mathbf{G}

- ii) Training of the second layer:

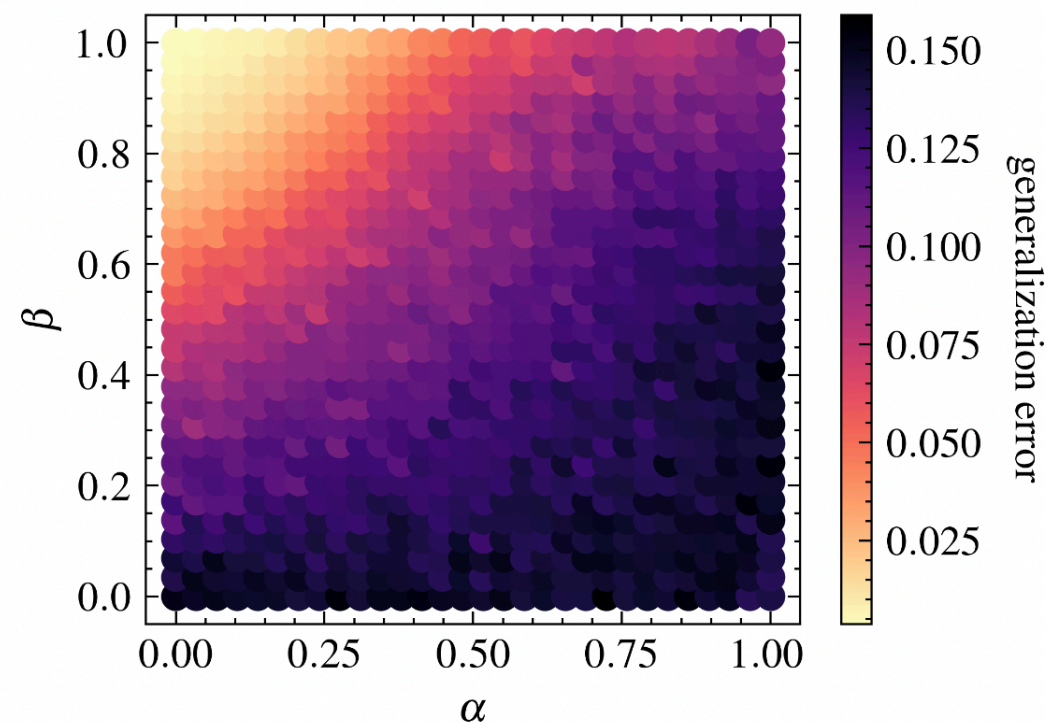
- Ridge regression
- $$\hat{\mathbf{w}} := \arg \min_{\mathbf{w} \in \mathbb{R}^k} \frac{1}{2m} \sum_{i=1}^m \left(y_i - \frac{1}{\sqrt{k}} \mathbf{w}^T \sigma(\hat{\mathbf{F}} \mathbf{x}_i) \right)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$



- Generalization error: $\mathbb{E}_{(\mathbf{x}, y)} \left[\left(y - \frac{1}{\sqrt{k}} \hat{\mathbf{w}}^T \sigma(\hat{\mathbf{F}} \mathbf{x}) \right)^2 \right]$

Scalings: data spread and learning rate

- Scalings of learning rate η and data spread $\|\Sigma\|$ shape the generalization
- To control these, we define
 - A “strength parameter” $\beta \in [0, 1]$ governing joint scaling: $\eta\|\Sigma\| \asymp n^\beta$
 - A “weighting parameter” $\alpha \in [0, 1]$ controlling individual scalings:
 - data spread $\|\Sigma\| \asymp n^{\beta(1-\alpha)}$ and learning rate $\eta \asymp n^{\beta\alpha}$



Theorem: Conditional Gaussian Equivalence

- *(Informal Statement) The following two feature maps are equivalent in terms of generalization (and training) errors:*

- *Original feature map* $\phi(\mathbf{x}) := \sigma(\hat{\mathbf{F}}\mathbf{x})$

- *A conditional feature map defined as*

$$\hat{\phi}(\mathbf{x}; c, \kappa_c) := \underbrace{\nu(c, \kappa_c)}_{\text{mean}} + \underbrace{\Psi(c, \kappa_c)}_{\text{cross-covariance}} \mathbf{z}^\perp + \underbrace{\Phi(c, \kappa_c)^{1/2}}_{\text{remaining covariance}} \mathbf{g} \quad \text{with} \quad \mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_n) \quad \text{and} \quad \mathbf{z}^\perp \text{ derived from } \mathbf{x}$$

mean cross-covariance remaining covariance

**This is an equivalent conditional Gaussian feature map,
conditioned on**

- (i) c : mixture component index,**
- (ii) κ_c : alignment of the input sample with the spikes in data covariance and gradient.**

Theorem: Equivalent Hermite (Polynomial) Model

- *(Informal Statement) If strength parameter satisfies $\frac{l-2}{l-1} < \beta < \frac{l-1}{l}$, then, we can replace original activation function with the following one without any change on the generalization (and training) performance.*

$$\hat{\sigma}_l(x) := \left(\sum_{j=0}^{l-1} \frac{1}{j!} h_j H_j(x/b) \right) + h_l^* z \quad \text{with} \quad z \sim \mathcal{N}(0, 1)$$

(finite-order) Hermite expansion remainder term to match the variance

Note: $H_j : \mathbb{R} \rightarrow \mathbb{R}$ denotes j -th (probabilist's) Hermite polynomial.

Theorem: Equivalent Hermite (Polynomial) Model

- *(Informal Statement) If strength parameter satisfies $\frac{l-2}{l-1} < \beta < \frac{l-1}{l}$, then, we can replace original activation function with the following one without any change on the generalization (and training) performance.*

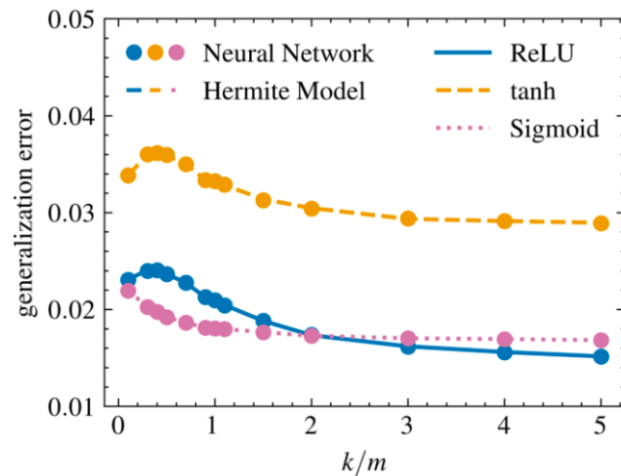
$$\hat{\sigma}_l(x) := \left(\sum_{j=0}^{l-1} \frac{1}{j!} h_j H_j(x/b) \right) + h_l^* z \quad \text{with } z \sim \mathcal{N}(0, 1)$$

(finite-order) Hermite expansion remainder term to match the variance

Note: $H_j : \mathbb{R} \rightarrow \mathbb{R}$ denotes j -th (probabilist's) Hermite polynomial.

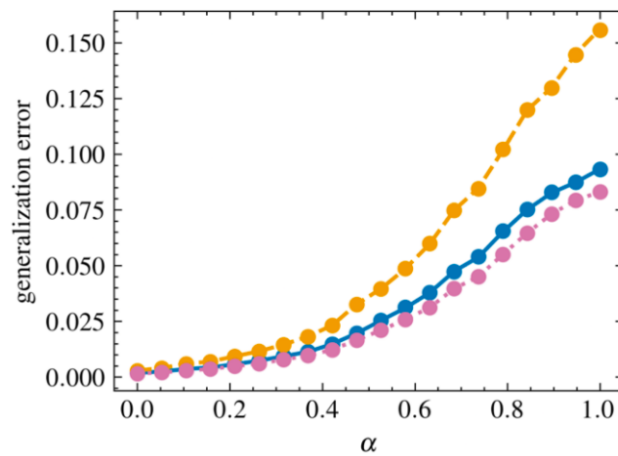
- Equivalent “Hermite Model”: $\frac{1}{\sqrt{k}} \mathbf{w}^T \hat{\sigma}_l(\hat{\mathbf{F}} \mathbf{x})$

Simulation: Impacts of Data Spread and Learning Rate



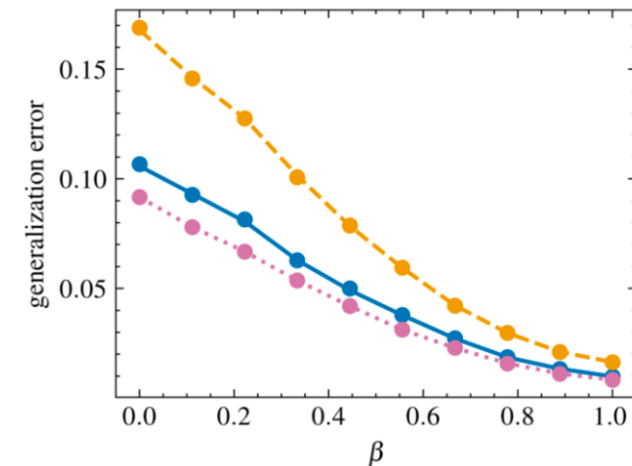
(a) Impact of k/m

($\beta = 3/4^-$ and $\alpha = 1/2$)



(b) Impact of α

($\beta = 3/4^-$ and $k/m = 1$)



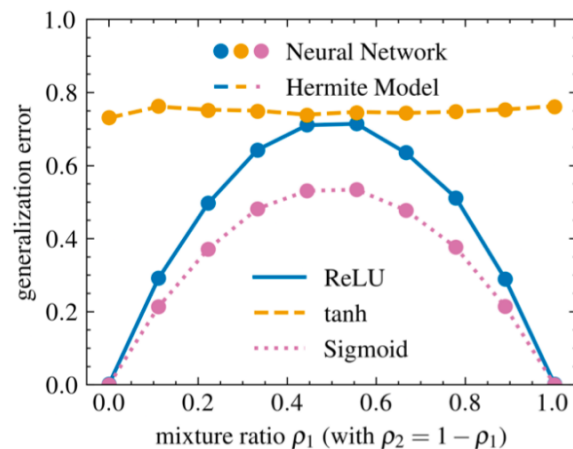
(c) Impact of β

($\alpha = 1/2$ and $k/m = 1$)

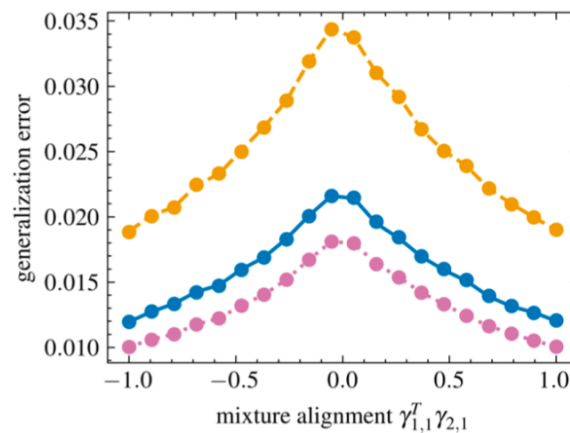
■ Reminder: data spread $\|\Sigma\| \asymp n^{\beta(1-\alpha)}$ and learning rate $\eta \asymp n^{\beta\alpha}$

- (a, b, c) The generalization errors of the NNs and the Hermite model closely align.
- (b) High data spread leads to better performance compared to high learning rate.
- (c) Larger strength parameter results in improved generalization in general.

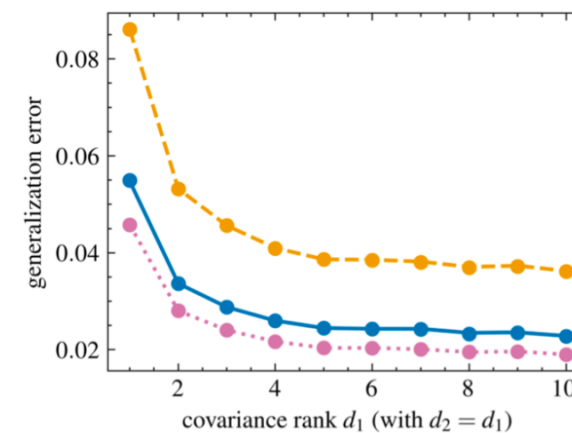
Simulation: Impacts of Mixture Properties



(a) Impact of mixture ratio
for classification



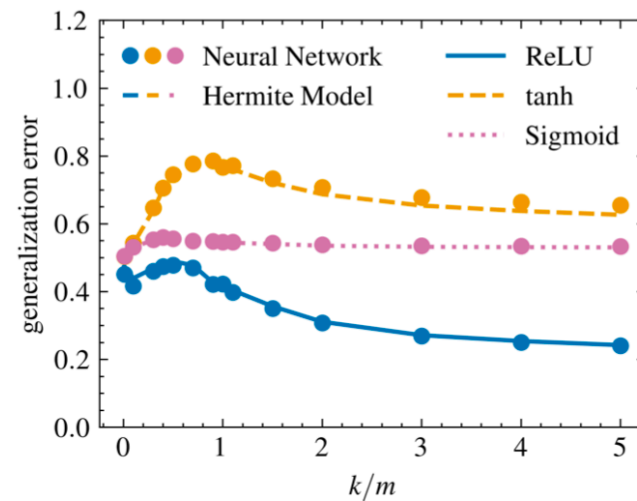
(b) Impact of mixture
alignment for regression



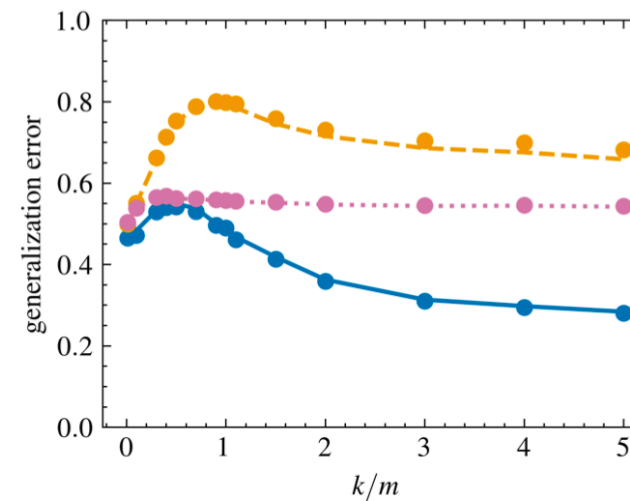
(c) Impact of covariance
rank for regression

- The generalization errors of the NNs and the Hermite model closely align.
- Mixture properties significantly affect the generalization errors.

Real Data: Fashion-MNIST Classification



(a) Class 0 vs. Class 1
(T-shirt/top vs. Trouser)



(b) Class 5 vs. Class 9
(Sandal vs. Ankle boot)

- The generalization errors of the NNs and the Hermite model closely align.

Summary

- **Takeaway**: Data distribution (data spread) impacts the generalization performance of neural networks together with the feature learning.
- Under Gaussian Mixtures data assumption and with feature learning via one gradient step, we found simpler models equivalent to two-layer NNs:
 - A conditional Gaussian model,
 - A polynomial model formed by Hermite polynomials.