

Classifying high-dimensional Gaussian mixtures: Where kernel methods fail and neural networks succeed

Jaehyoung Hong

ICML 2021 Poster

Kernel methods shows competitive performance to Neural network for some task

Jacot *et al* '18 → wide neural networks, with appropriate scaling behave as kernel methods

Matthews *et al* '18

Lee *et al* '18

Garriga *et al* '18

→ On some benchmark tasks Kernel methods do almost as well as neural networks

Do Neural Networks only learn efficiently
if Kernel methods can also learn

Kernel methods did not beat the Neural network for high-dimensional setting

Jacot *et al* '18 → wide neural networks, with appropriate scaling behave as kernel methods

Matthews *et al* '18
Lee *et al* '18
Garriga *et al* '18

→ On some benchmark tasks Kernel methods do almost as well as neural networks

Do Neural Networks only learn efficiently
if Kernel methods can also learn

Chizat & Bach '20 → No! Neural Networks can grasp underlying low dimensional data structure, Kernels cannot

Ghorbani *et al* '19 , '20 → No! Data structure can break the curse of dimensionality for neural networks but not for kernel methods

Mean-Field Limit

Requires very large
hidden layer

✓ Today's topic: NN > Kernel method for $O(1)$ hidden nodes, 2 layers

Compare two-layer neural network(2LNN) and random features

- $\#(\text{sample}) = N, \#(\text{input dimension}) = D, \#(\text{random features}) = P$
- $\text{Input } x = (x_r) \in \mathbb{R}^D, \text{Labels } y \in \{-1, 1\}$: Binary classification (for simplicity)

1. 2LNN ($\phi_\theta^L: \mathbb{R}^D \rightarrow \mathbb{R}$)

x



$$\lambda^k \equiv \frac{1}{\sqrt{D}} \sum_{r=1}^D w_r^k x_r \quad \longrightarrow \quad \phi_\theta^L(x) = \sum_{k=1}^K v^k g(\lambda^k)$$

- ✓ Trainable weight $w = (w_r^k) \in \mathbb{R}^{K \times D}, v = (v^k) \in \mathbb{R}^K$
- ✓ $\frac{K}{D} = O(1)$: Small node setting, $t \equiv \frac{N}{D} = O(1)$: High dimensional setting ($D \rightarrow \infty$)

2. Random Features ($\phi_\theta^R: \mathbb{R}^D \rightarrow \mathbb{R}$)

$$x \quad \longrightarrow \quad u_i \equiv \frac{1}{\sqrt{D}} \sum_{r=1}^D F_{ir} x_r \quad \longrightarrow \quad z_i = \psi(u_i) \quad \longrightarrow \quad \phi_\theta^R(z) = \frac{1}{\sqrt{P}} \sum_{i=1}^P w_i z_i$$

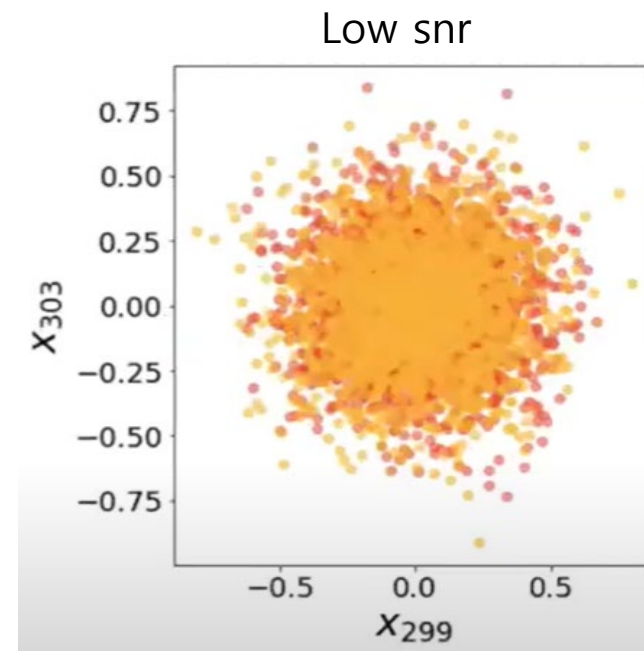
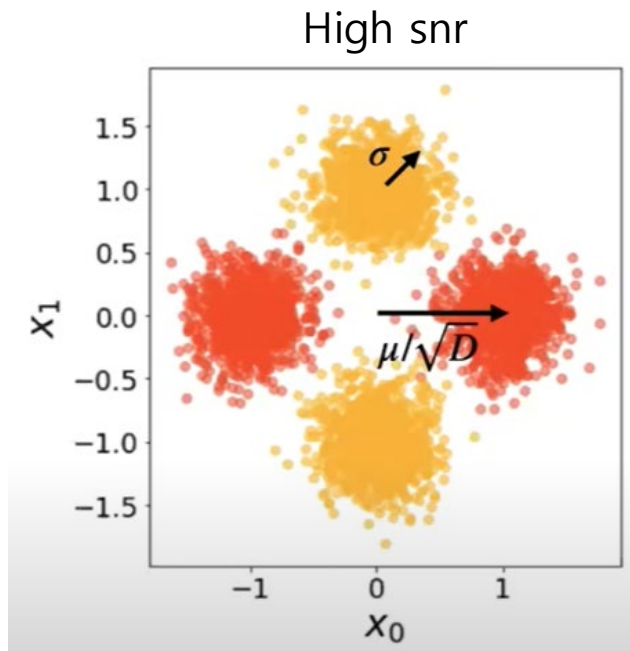
- ✓ Trainable weight $w = (w_i) \in \mathbb{R}^P$
- ✓ $\gamma \equiv \frac{P}{D} = O(1), N > P$

$$\epsilon_c(\theta) = \mathbb{E} \mathbf{H}[-y \phi_\theta(x)]$$

$N = O(D^2)$
needed for RF
while $N = O(D)$
for 2LNN

RF is similar to neural network for High signal to noise ratio(SNR) while worse for Low SNR

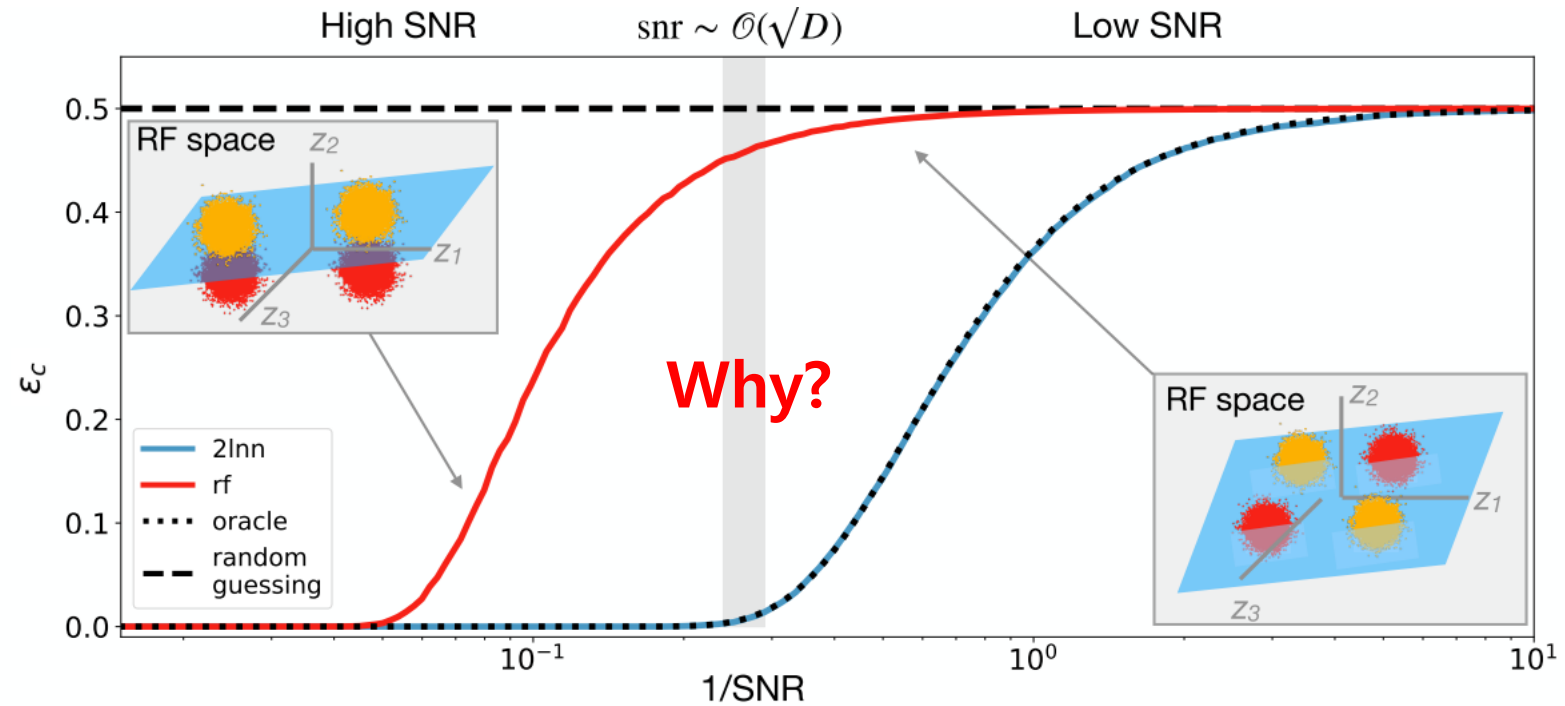
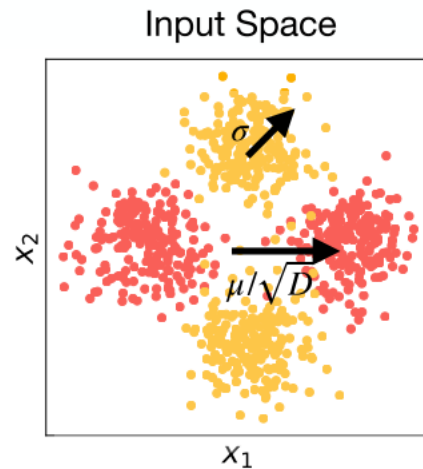
- XOR inputs



- ✓ $snr = |\mu|/\sqrt{D}\sigma$ for input $x_i = \mu_i + \sigma z_i$ ($z_i \sim N(0,1)$), XOR: $\mu_+ = (\pm\sqrt{D}, 0, \dots, 0)$, $\mu_- = (0, \pm\sqrt{D}, 0, \dots, 0)$
- ✓ *oracle*: Known about μ_i and labels x_i to that of nearest μ_i

RF is similar to neural network for High signal to noise ratio(SNR) while worse for Low SNR

- XOR inputs



- ✓ $\text{snr} = |\mu|/\sqrt{D}\sigma$ for input $x_i = \mu_i + \sigma z_i$ ($z_i \sim N(0,1)$), XOR: $\mu_+ = (\pm\sqrt{D}, 0, \dots, 0)$, $\mu_- = (0, \pm\sqrt{D}, 0, \dots, 0)$
- ✓ *oracle*: Known about μ_i and labels to nearest μ_i

Setup for dynamics of 2LNN: Moments of λ^k

- 2LNN for GM classification

✓ Sample (x, y) is drawn from $q(x, y) = q(y)q(x|y)$,
 $q(x, y) = \sum_{\alpha \in \mathcal{S}(y)} \mathcal{P}_\alpha N_\alpha(x), N_\alpha(x) \sim N(\frac{\mu^\alpha}{\sqrt{D}}, \Omega^\alpha)$

✓ Online learning, $\Delta = \sum_{k=1}^K v^k g(\lambda^k) - y$

$$\begin{aligned} dw_i^k &= -\frac{\eta}{\sqrt{D}} v^k \Delta g'(\lambda^k) x_i - \frac{\eta}{\sqrt{D}} \kappa w_i^k, \\ dv^k &= -\frac{\eta}{D} g(\lambda^k) \Delta - \frac{\eta}{D} \kappa v^k, \end{aligned} \quad \begin{array}{l} \text{SGD} \\ L^2\text{-Regularization} \end{array}$$

✓ Static

$$\begin{aligned} \text{pmse}(\theta) &= \mathbb{E}_{q(x, y)} (y - \phi_\theta(x))^2 \\ &= \sum_y \sum_{\alpha \in \mathcal{S}(y)} q(y_i) \mathcal{P}_\alpha \mathbb{E}_\alpha \left[\sum_k v^k g(\lambda^k) - y \right]^2 \end{aligned}$$

Setup for dynamics of 2LNN: Moments of λ^k

- 2LNN for GM classification

$$\begin{aligned}
 \text{Average over the inputs } x \downarrow \\
 \text{pmse}(\theta) = \mathbb{E}_x \left[\left(\sum_{k=1}^K v^k g \left(\frac{w^k x}{\sqrt{D}} \right) - y \right)^2 \right] & \xrightarrow{\text{Zoom into each cluster}} \text{pmse}(\theta) = \sum_{\text{clusters}} \underbrace{\mathbb{E}_x \left[\left(\sum_{k=1}^K v^k g \left(\frac{w^k x}{\sqrt{D}} \right) - y \right)^2 \mid x \sim \mathcal{N} \left(\frac{\mu_\alpha}{\sqrt{D}}, \Sigma_\alpha \right) \right]}_{\text{Expectation over cluster } \alpha} \\
 & \text{pmse}(\theta) = \sum_{\text{clusters}} \underbrace{\mathbb{E}_{\lambda_\alpha} \left[\left(\sum_{k=1}^K v^k g(\lambda^k) - y \right)^2 \mid x \sim \mathcal{N} \left(\frac{\mu_\alpha}{\sqrt{D}}, \Sigma_\alpha \right) \right]}_{\substack{\uparrow \\ \text{Average over} \\ \text{the local fields}}}
 \end{aligned}$$

- ✓ λ^k are jointly Gaussian when averages are evaluated over just a single distribution in the mixture
- ✓ pmse is determined by the moments of λ^k

Setup for dynamics of 2LNN: Moments of λ^k

- 2LNN for GM classification

- ✓ Order parameters: Moments of λ^k

$$M_{\alpha}^k \equiv \mathbb{E}_{\alpha} \lambda^k = \frac{1}{D} \sum_r w_r^k \mu_r^{\alpha},$$

$$Q_{\alpha}^{k\ell} \equiv \text{Cov}_{\alpha} (\lambda^k, \lambda^{\ell}) = \frac{1}{D} \sum_{r,s} w_r^k \Omega_{rs}^{\alpha} w_s^{\ell}.$$

- ✓ Since any average over a Gaussian is fcn of two moments,

$$\lim_{D \rightarrow \infty} \text{pmse}(\theta) \rightarrow \text{pmse}(Q, M, v).$$

- ✓ Pmse of network $\approx (Q, M, v)$ dynamics

Dynamics of weight are determined by order parameters

- Moments of functions of weakly correlated variables has explicit form
- $t = N/D$ can be regarded as continuous time variable for online learning

Let r.v. $x, y \in \mathbb{R}$, jointly Gaussian & weakly correlated:

$$\begin{aligned}
 P(x, y) &= \frac{1}{2\pi\sqrt{\det M_2}} \exp \left[-\frac{1}{2} \begin{pmatrix} x - \bar{x} & y - \bar{y} \end{pmatrix} M_2^{-1} \begin{pmatrix} x - \bar{x} \\ y - \bar{y} \end{pmatrix} \right], \quad M_2 = \begin{pmatrix} C_x & \epsilon M_{12} \\ \epsilon M_{12} & C_y \end{pmatrix} \\
 &= \frac{1}{2\pi\sqrt{C_x C_y}} e^{-\frac{1}{2C_x}(x-\bar{x})^2 - \frac{1}{2C_y}(y-\bar{y})^2} [1 - \epsilon(x - \bar{x})(C_x^{-1}M_{12}C_y^{-1})(y - \bar{y}) + O(\epsilon^2)] \\
 \mathbb{E}_{(x,y)} [f(x)g(y)] &= \mathbb{E}_x [f(x)] \mathbb{E}_y [g(y)] \\
 &\quad + \epsilon \mathbb{E}_x [f(x)(x - \bar{x})] (C_x^{-1}M_{12}C_y^{-1}) \mathbb{E}_y [g(y)(y - \bar{y})] + O(\epsilon^2).
 \end{aligned}$$

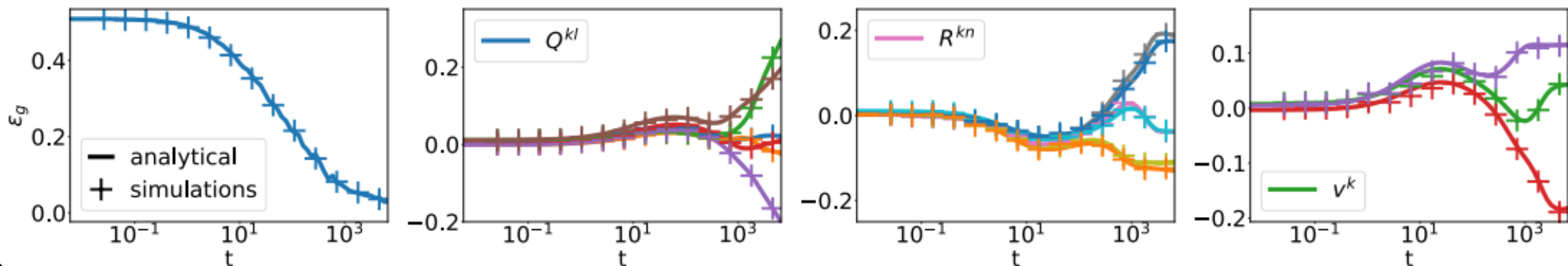
- Dynamics of second layer ($\Omega_\alpha = \Omega$)

$$\begin{aligned}
 \mathbb{E} dv^k &= \sum_{\alpha \in \mathcal{S}(+)} \mathcal{P}_\alpha dv_{\alpha+}^k + \sum_{\alpha \in \mathcal{S}(-)} \mathcal{P}_\alpha dv_{\alpha-}^k, \\
 dv^k &= -\frac{\eta}{D} g(\lambda^k) \Delta - \frac{\eta}{D} \kappa v^k, \\
 dv_\alpha^k &= \frac{\eta}{D} \mathbb{E}_\alpha y_\alpha g(\lambda^k) - \frac{\eta}{D} \sum_j v^j \mathbb{E}_\alpha g(\lambda^k) g(\lambda^j) - \frac{\eta}{D} \kappa v^k.
 \end{aligned}$$

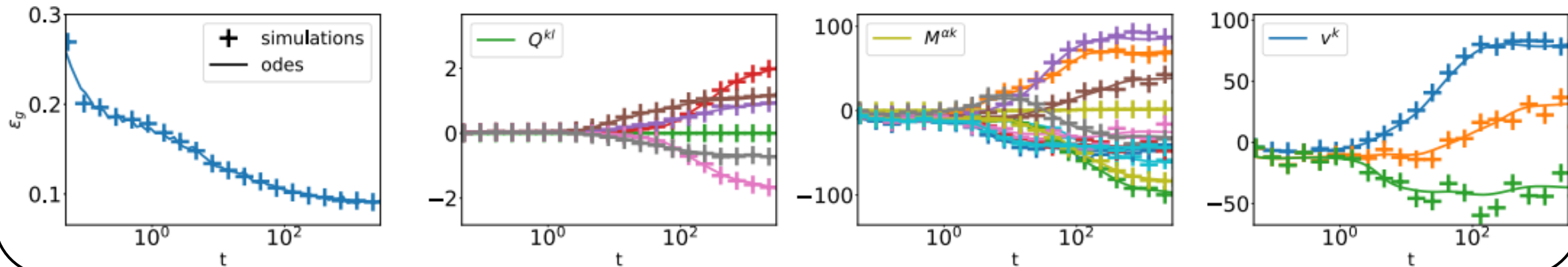
- ✓ Dynamics of weight are determined by order parameters /
Expectation of activation function is estimated by MC methods

Order parameters and weights from ODE are agreed with simulation(Single run of SGD)

- Simulation data(GM mixture)

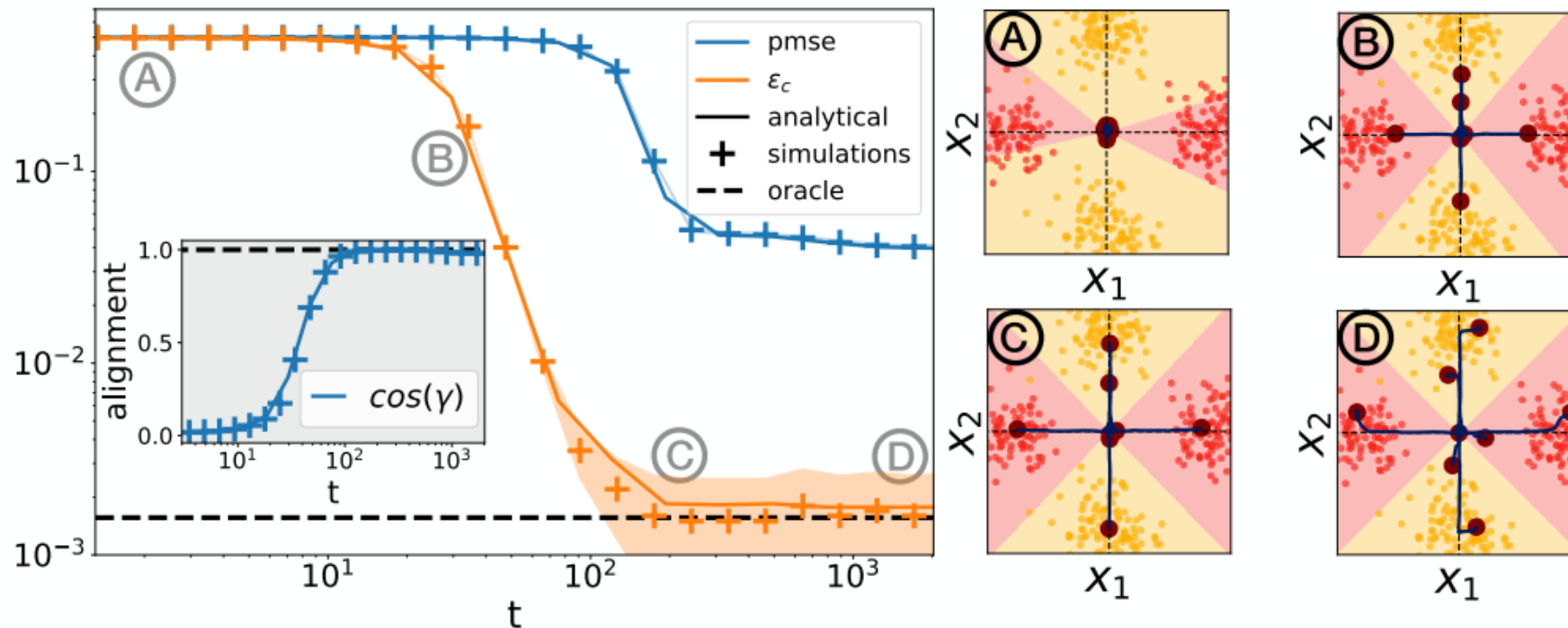


- FashionMNIST, Ω_α is different for α



2LNN learns XOR inputs when first weight vectors approach the four means

- (Left) Error plot (Right): First weight projected onto the sample space (Blood dots)

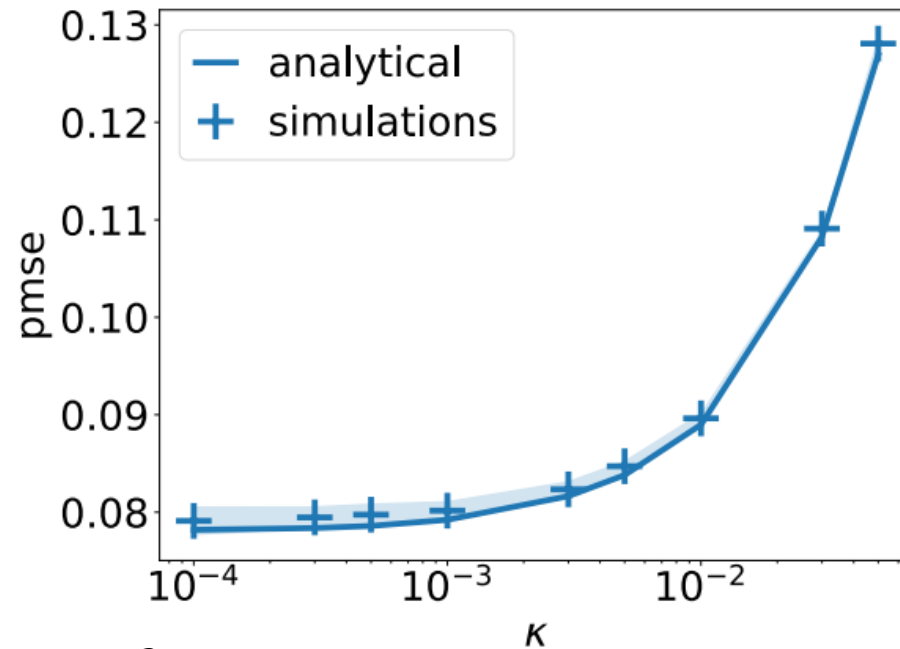


C: First weight
finds four means
(Maximal angle)

D: Decrease pmse

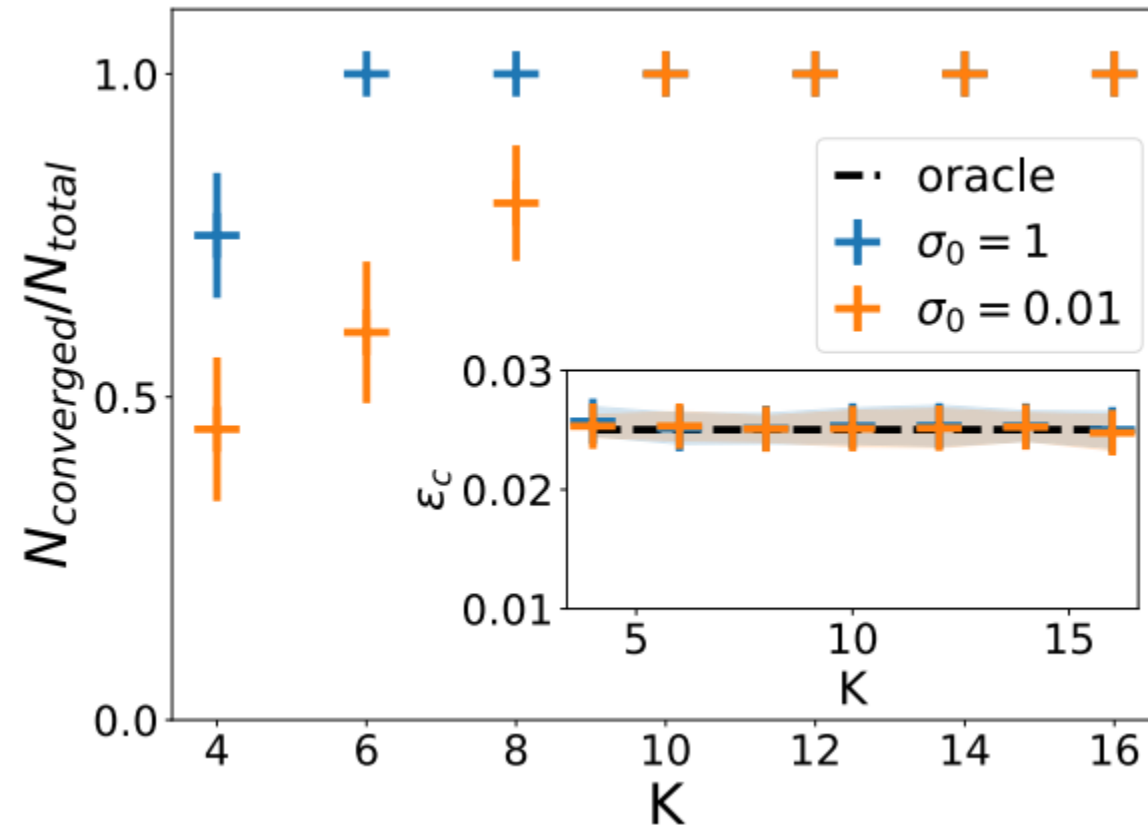
Predicting long-time performance using ansatz also agreed with simulation

- Long time performance of ODE: Asymptotic fixed point
- Scale $K^2 + 4K$ equations for each time step
- Ansatz (e.g. symmetry) reduce parameters to K



- ✓ L^2 -regularization hurt performance
- ✓ $K = 4$ is enough

Over-parametrization do not improve test error, but has acceleration effect



Dynamics of RF

- Assume $N \gg P$, for any finite D, P , running the algorithm upto convergence (\approx Taking the limit $t \rightarrow \infty$)

- ρ_τ are the eigenvalues of the feature's covariance matrix $\Omega_{ij} = \mathbb{E}z_i z_j$, with associated eigenvector Γ_τ . $\tilde{\Phi}_\tau \equiv \sum_{i=1}^P \Gamma_{\tau i} \Phi_i / \sqrt{P}$

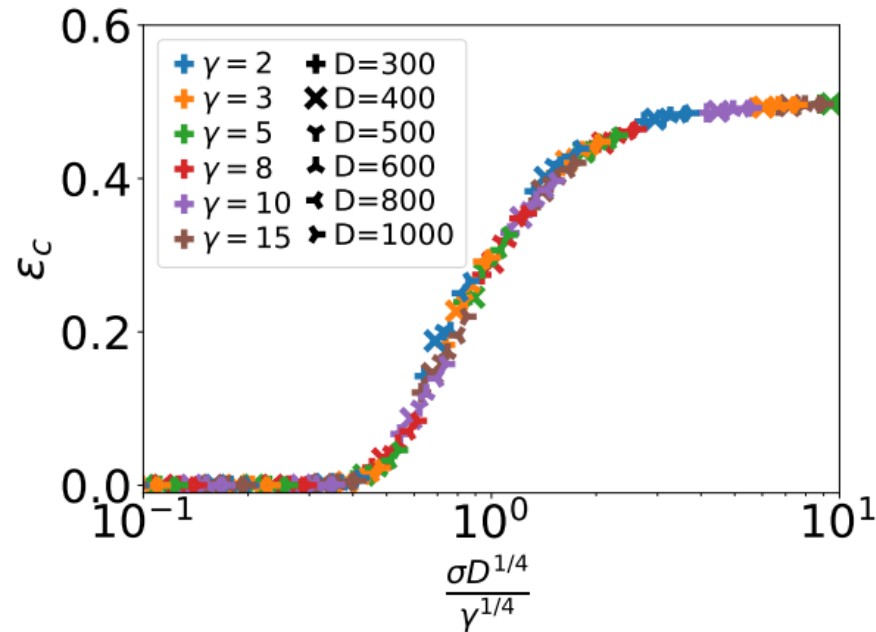
$$\text{pmse}_{t \rightarrow \infty} = \frac{1}{2} \left(1 - \sum_{\tau} \frac{\tilde{\Phi}_\tau^2}{\rho_\tau} \right),$$

$$M_\alpha = \sum_{i=1}^P \frac{\hat{w}_i \mathbb{E}_\alpha[z_i]}{\sqrt{P}}, \quad Q_\alpha = \sum_{i=1}^P \frac{\hat{w}_i \hat{w}_j}{P} \text{Cov}_\alpha(z, z).$$

$$\epsilon_{ct \rightarrow \infty} = \frac{1}{2} \left(1 - \sum_{\alpha} \mathcal{P}_\alpha y \operatorname{erf} \left(\frac{M_\alpha}{\sqrt{2Q_\alpha}} \right) \right)$$

Classification error of RF for various values shows that $P \approx D^2$ features are required

- Classification error when $N \gg P$



✓ Graph shows for good performance $P = O(D^2)$ while $N \gg P$ & $N = O(D)$ makes it impossible

✓ Indeed, classification error is fcn of $\sigma D^{\frac{1}{2}} / \min(N, P)^{1/4}$

✓ If $N = O(D)$ & $\sigma \gg N^{\frac{1}{4}} / D^{1/2}$,
performance degrades to no more than random guess

In Low SNR setting($snr \sim O(1)$), the transformation of the means is only linear

- $\frac{|\mu|}{\sqrt{D}} \sim O(1)$ & $\sigma \sim O(1) \rightarrow \frac{F_{ir}\mu_r}{D} \sim O(\frac{1}{\sqrt{D}})$

$$a \equiv \mathbb{E} \psi(\sigma \zeta), \quad b \equiv \mathbb{E} \zeta \psi(\sigma \zeta), \quad c^2 \equiv \mathbb{E} \psi(\sigma \zeta)^2$$

$$\mathbb{E} z_i = a + b \sum_{r=1}^D \frac{F_{ir}\mu_r}{\sigma D}$$

$$\text{cov}(z_i, z_j) = \begin{cases} c^2 - a^2, & i = j, \\ b^2 \sum_r \frac{F_{ir}F_{jr}}{D} & i \neq j. \end{cases}$$

- ✓ Transformation of the means is only linear
- ✓ RF cannot separate XOR inputs with Low SNR setting in feature space
- ✓ RF only separate when the centres of data are separated enough

Convergence of RF to kernel methods shows that kernel methods also fails for Low SNR setting

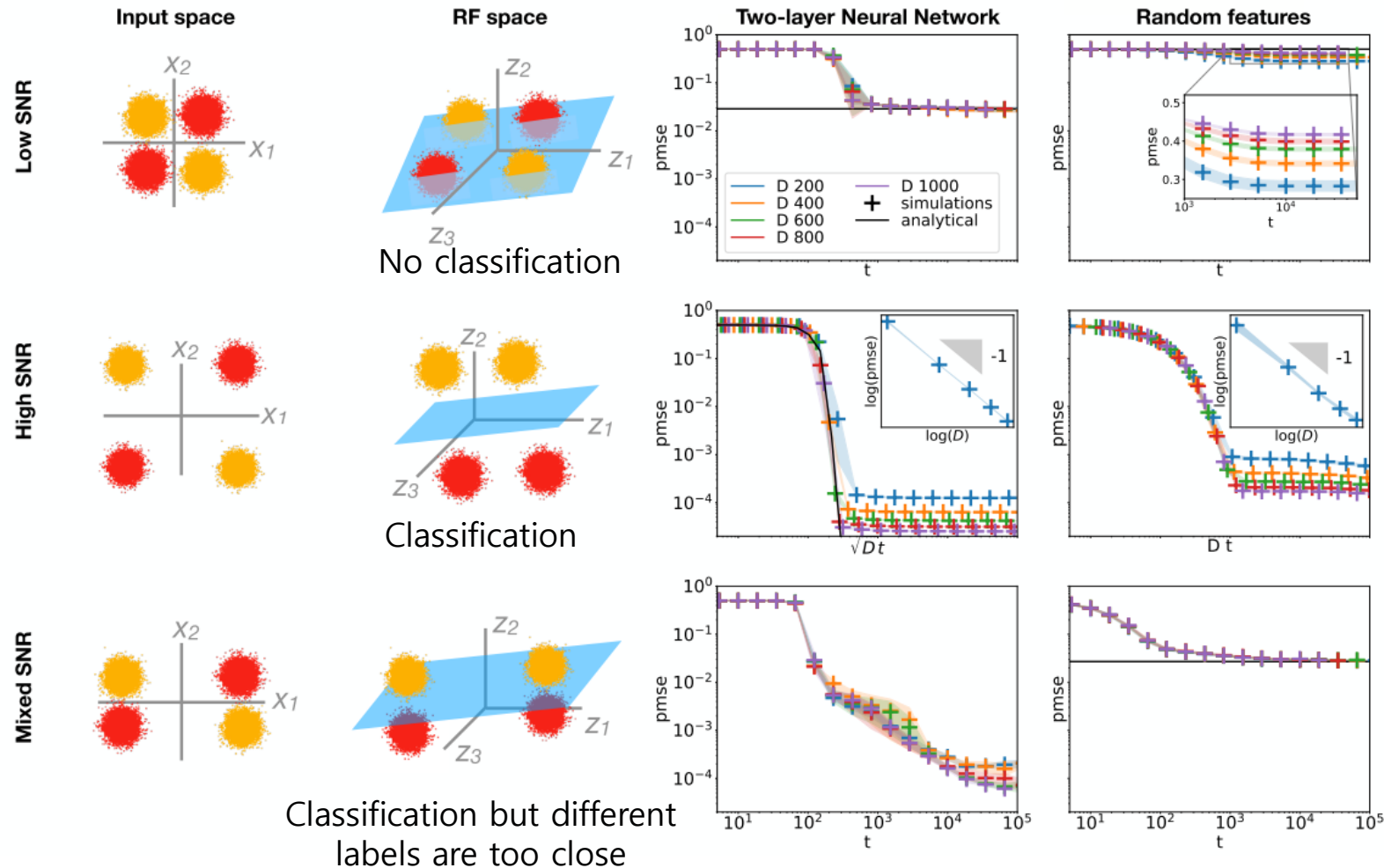
- Rahimi & Recht, 2008; 2009 ($D, P \rightarrow \infty$ and then $\gamma \rightarrow \infty$)

$$K(\mathbf{x}, \mathbf{y}) = \frac{1}{P} \sum_{i=1}^P \mathbb{E}_F \left[\psi \left(\sum_{r=1}^D \frac{x_r F_{ir}}{\sqrt{D}} \right) \psi \left(\sum_{s=1}^D \frac{y_s F_{is}}{\sqrt{D}} \right) \right]$$

- With Low SNR setting

$$\begin{aligned} c^2 &= \mathbb{E} K(\sigma \boldsymbol{\omega}_1, \sigma \boldsymbol{\omega}_1), \quad a^2 = \mathbb{E} K(\sigma \boldsymbol{\omega}_1, \sigma \boldsymbol{\omega}_2), \\ b^2 &= D \sigma^2 \left[-a^2 + \mathbb{E} K \left(\frac{\boldsymbol{\mu}}{\sqrt{D}} + \sigma \boldsymbol{\omega}_1, \frac{\boldsymbol{\mu}}{\sqrt{D}} + \sigma \boldsymbol{\omega}_2 \right) \right] \\ \mathbb{E} z_i &= a + b \sum_{r=1}^D \frac{F_{ir} \mu_r}{\sigma D} \\ \text{cov}(z_i, z_j) &= \begin{cases} c^2 - a^2, & i = j, \\ b^2 \sum_r \frac{F_{ir} F_{jr}}{D} & i \neq j. \end{cases} \end{aligned}$$

2LNN is better than RF for GM classification



Summary

