

Random Matrix Theory in Stochastic Optimisation

Chanju Park

With Gert Aarts and Biagio Lucini



Our result

- Stochastic optimisation of a matrix has universal features described by Random Matrix Theory.
- The stochasticity of the system scales as $\alpha/|\mathcal{B}|$.
 α : Step size $|\mathcal{B}|$: Batch size
- Linear Scaling Rule [Goyal et al, arXiv:1706.02677]

Outline

1. Stochastic Optimisation
2. Random Matrix Theory
3. Random Matrix Theory in Stochastic Gradient Descent
4. Experiments

Stochastic Optimisation

Robbins - Monro algorithm

Imagine we have a stochastic variable x and a function $f(x)$ with an extremum at $x = x^*$. We can find the value x^* by following iterative algorithm.

$$x_{n+1} = x_n - \alpha_n \left(\mathbb{E} [f(x)] - f(x^*) \right)$$

For α_n satisfying $\sum_{n=0}^{\infty} \alpha_n = \infty$, $\sum_{n=0}^{\infty} \alpha_n^2 < \infty$

Eg) Stochastic Gradient Descent

$$x = W_{ij}, \quad f(W_{ij}) = \frac{\partial \mathcal{L}}{\partial W_{ij}}, \quad f(W_{ij}^*) = 0$$

Langevin Dynamics

Stochastic optimisation

$$x_{n+1} = x_n - \alpha \left(\mathbb{E} [f(x_n)] - f(x^*) \right)$$

$$f(x^*) = 2, \quad P(f(x_n)) \sim \mathcal{N}(\mu_n, \sigma_n)$$

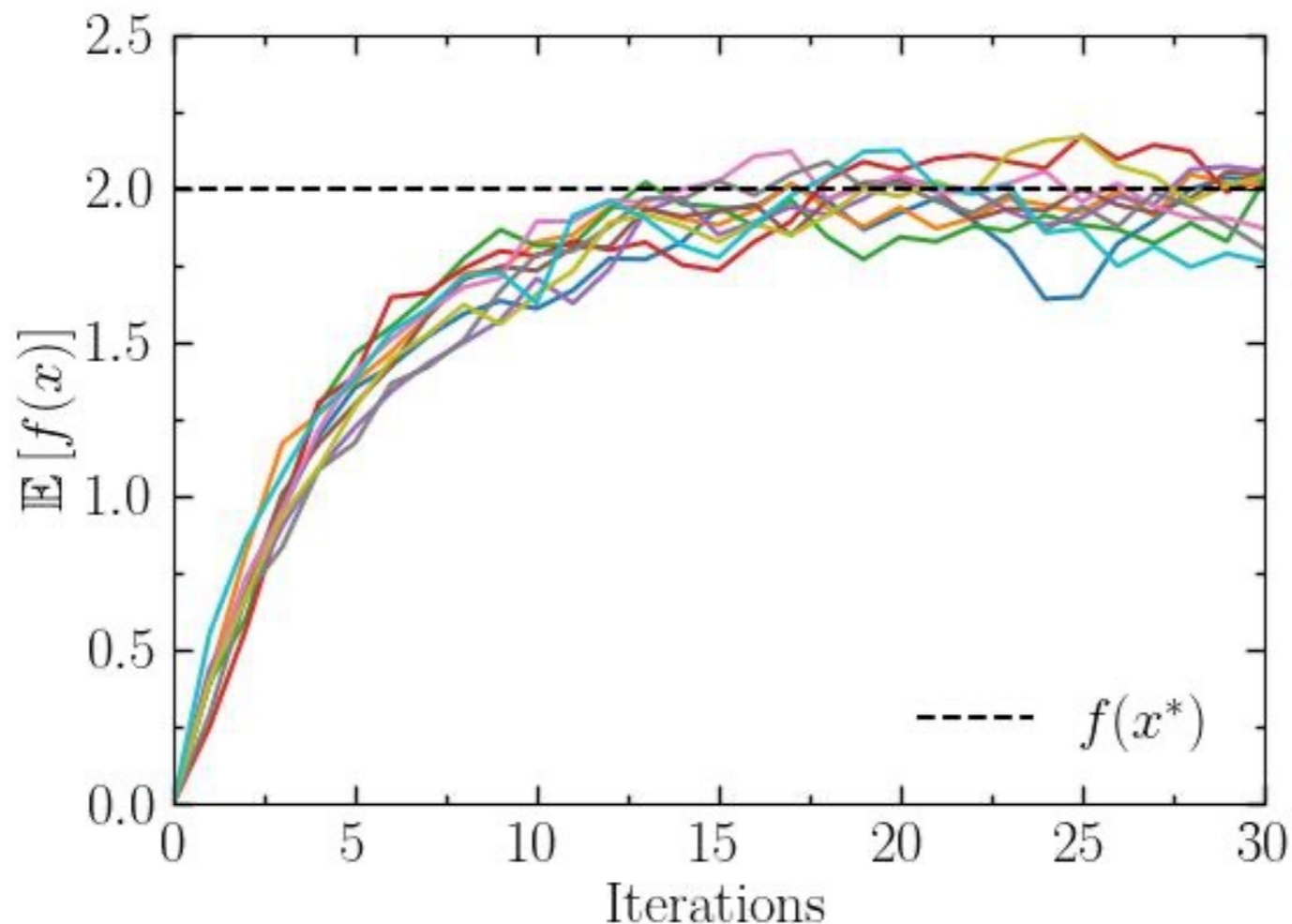
Langevin equation

$$\frac{dx}{dt} = -K(x; t) + \sqrt{2}g(x; t)\eta$$

$$K(x; t) = \alpha(\mu(t) - f(x^*)), \quad g(x; t) = \alpha \frac{\sigma(t)}{\sqrt{2}}$$



Simulation



Theory

Stochastic optimisation can be described by Langevin dynamics.

Random Matrix Theory

Gaussian Orthogonal Ensemble


Imagine we have an ensemble of symmetric matrices, where the elements are Gaussian random variables.

$$M_{ij} \sim \mathcal{N}(0,1)$$

$$P(M_{ij}) \propto e^{-\frac{1}{2}\text{Tr}M_{ij}^2} \quad \text{invariant under } M' = O^T M O, \quad O^T O = I$$

Because of the symmetry, there is degeneracy in the representation and it is convenient to choose eigenvalues as the basis.

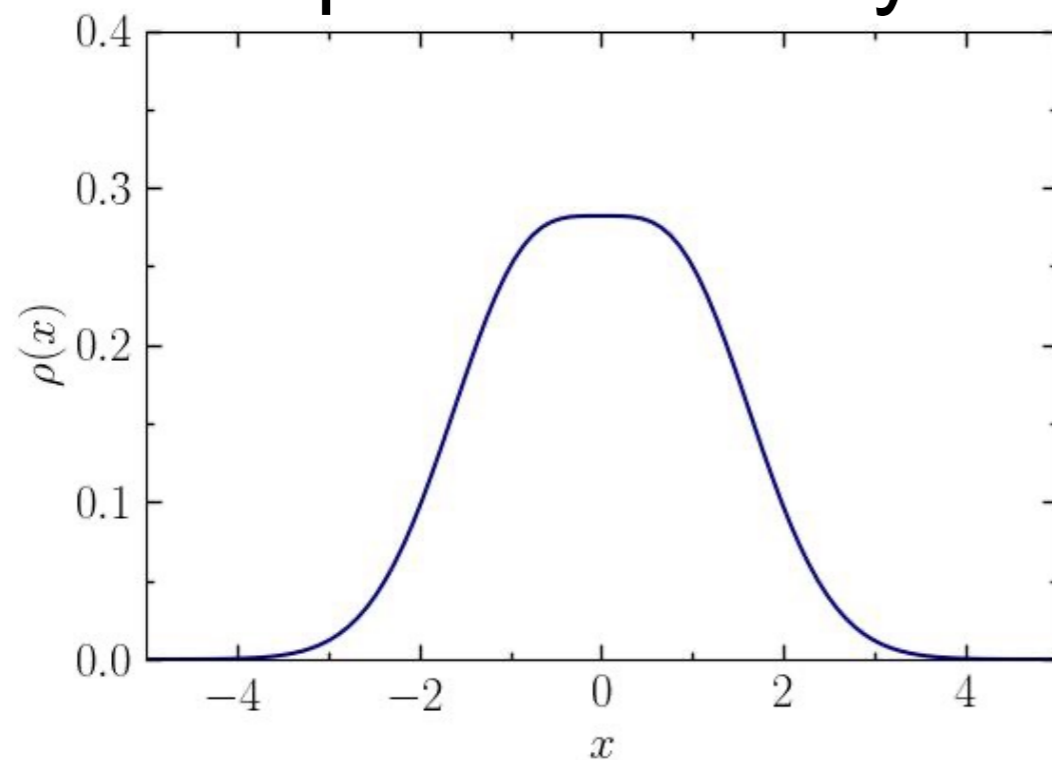
$$P(M_{ij}) \Rightarrow P(x_i) \propto \prod_{i < j} |x_i - x_j| e^{-\frac{1}{2} \sum_i x_i^2}$$

 Jacobian

Properties of Random Matrix eigenvalues

Some useful statistical properties of the eigenvalues are known for these types of matrices.

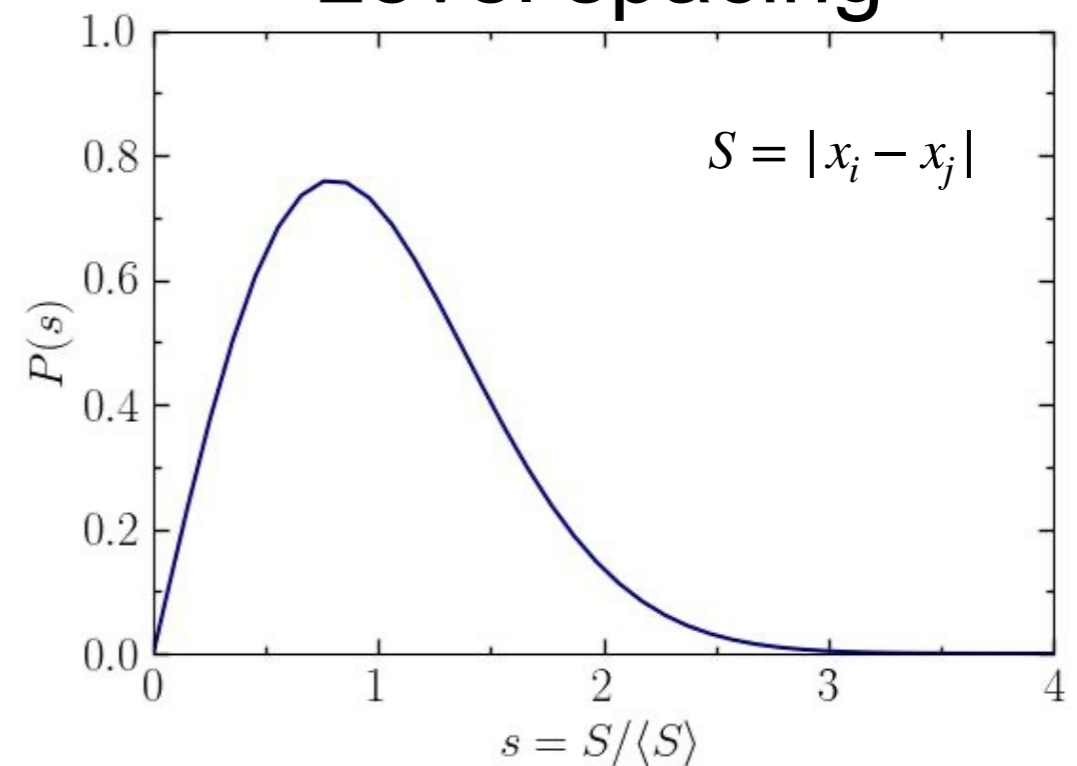
Spectral density



$$\rho(x) = \left\langle \frac{1}{N} \sum_{i=1}^N \delta(x - x_i) \right\rangle$$

“Wigner semi-circle”

Level spacing

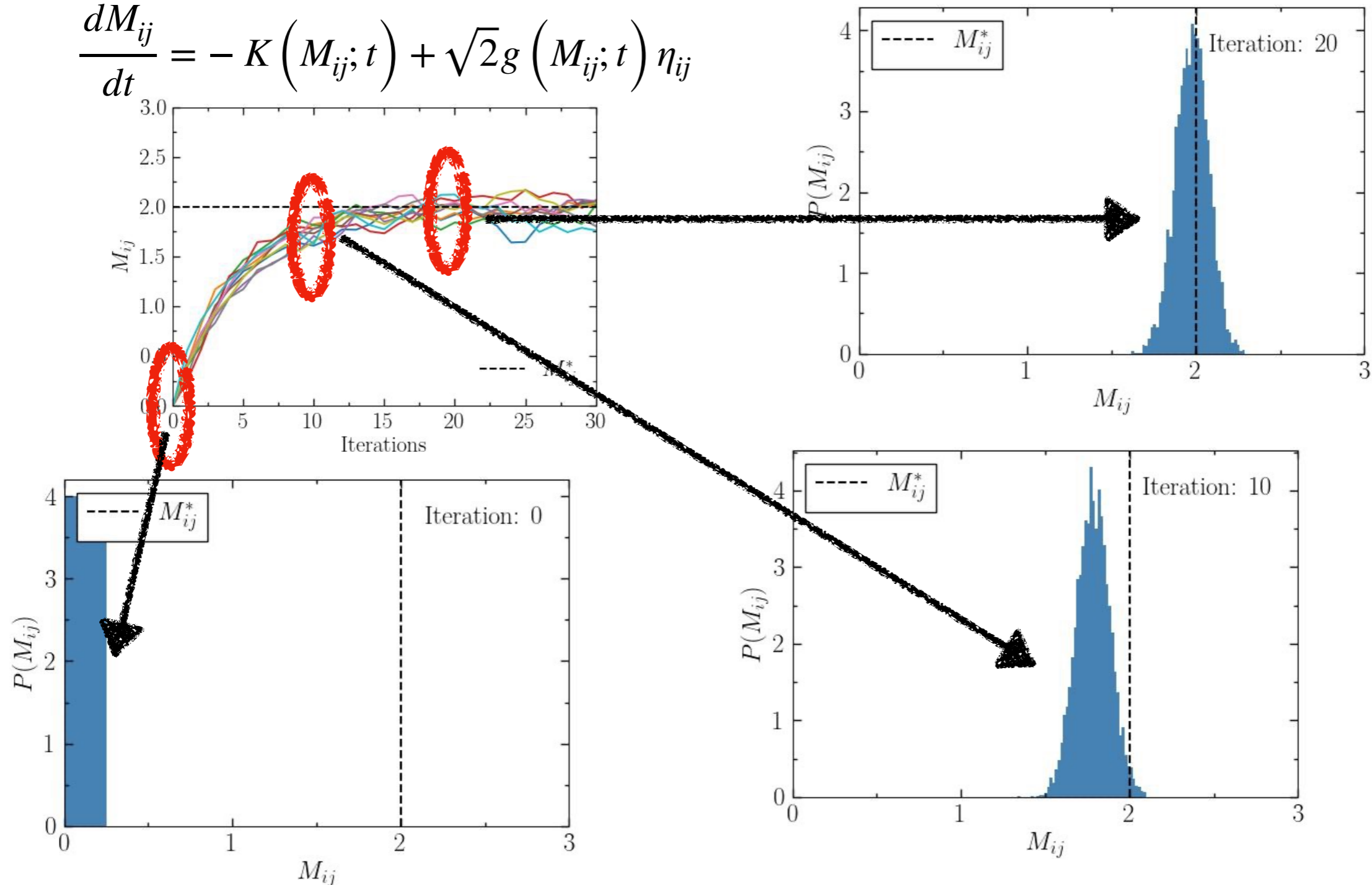


$$P(s) = \frac{\pi}{2} s e^{-\pi s^2/4}$$

“Wigner surmise”

Stochastic optimisation of Matrices

$$\frac{dM_{ij}}{dt} = -K(M_{ij}; t) + \sqrt{2}g(M_{ij}; t)\eta_{ij}$$



At each time slice, matrix elements are randomly distributed!
=> Random Matrix Theory!

Dynamics of eigenvalue

Dyson-Brownian motion

When a matrix is stochastically optimised,

$$\frac{dM_{ij}}{dt} = -K_{ij} + g_{ij}\eta_{ij}, \quad K_{ij} = \frac{\partial \mathcal{L}}{\partial M_{ij}}$$

Change of variable introduces Jacobian (Vandermonde determinant),

$$P(M_{ij}) \Rightarrow P(x_i) \propto \prod_{i < j} |x_i - x_j| e^{-\mathcal{L}} = e^{-\mathcal{L} + \sum_{i < j} \log(x_i - x_j)}$$

Dynamics of eigenvalues can be derived.

$$\frac{dx_i}{dt} = -K_{ii} + \sum_{j \neq i} \frac{g_{ij}^2}{x_i - x_j} + \sqrt{2} g_{ii} \eta_{ii}$$

Stochastic Gradient Descent (SGD)

Stochastic Gradient Descent is one of the Stochastic optimisation algorithms.

$$M_{n+1} = M_n + \alpha \left\langle \Delta_p \right\rangle_{p \in \mathcal{B}}$$

$$\left\langle \Delta_p \right\rangle_{p \in \mathcal{B}} \equiv \frac{1}{|\mathcal{B}|} \sum_{p \in \mathcal{B}} \Delta_p, \quad \Delta_p \equiv \left. \frac{\partial \mathcal{L}}{\partial M_n} \right|_p$$

Stochasticity is introduced from the finite batch size, and its strength can be extracted using the CLT.

$$\mathbb{V} \left[\Delta_p \right]_{p \in \mathcal{B}} \propto \frac{1}{|\mathcal{B}|} \tilde{g}^2, \quad \tilde{g}^2 \equiv \mathbb{V} \left[\left. \frac{\partial \mathcal{L}}{\partial M_n} \right] \right.$$

Where \tilde{g} is a variance of the gradient.

Langevin equation for SGD

Corresponding Langevin equation for stochastic gradient descent can be written as,

$$M'_{ij} = M_{ij} - \underbrace{\alpha \mathbb{E}_{ij} [\Delta_p]}_{= K_{ij}(M)} + \frac{\alpha}{\sqrt{|\mathcal{B}|}} \sqrt{\tilde{g}_{ij}^2} \eta_{ij}$$

$\eta_{ij} \sim \mathcal{N}(0,1)$

Dynamic equation for eigenvalue is given as,

$$x'_i = x_i + \underbrace{\alpha \mathbb{E}_{ii} [\Delta_p] + \frac{\alpha^2}{|\mathcal{B}|} \sum_{j \neq i} \frac{\tilde{g}_{ij}^2}{x_i - x_j}}_{= K_{ii}^{(\text{eff})}(x)} + \frac{\alpha}{\sqrt{|\mathcal{B}|}} \sqrt{2\tilde{g}_{ii}^2} \eta_i$$

$\eta_i \sim \mathcal{N}(0,1)$

Universal scaling factor

Stationary limit distribution

The distribution of the eigenvalues is obtained by solving the Fokker-Planck equation.

$$\partial_t P(\{x_i\}, t) = \sum_{i=1}^N \partial_{x_i} \left[\left(\frac{\alpha^2}{|\mathcal{B}|} \tilde{g}_i^2 \partial_{x_i} - K_{ii}^{(\text{eff})} \right) P(\{x_i\}, t) \right]$$

Stationary limit solution: Coulomb gas distribution

$$P(\{x_i\}) = \frac{1}{Z} \prod_{i < j} |x_i - x_j| e^{-\sum_i V_i(x_i)/\sigma_i^2}, \quad K_{ii}(x_i) = -\alpha \frac{dV_i(x_i)}{dx_i}$$

... details in Matteo's poster

and
$$\sigma_i^2 = \frac{\alpha}{|\mathcal{B}|} \frac{\tilde{g}_i^2}{\Omega_i}$$

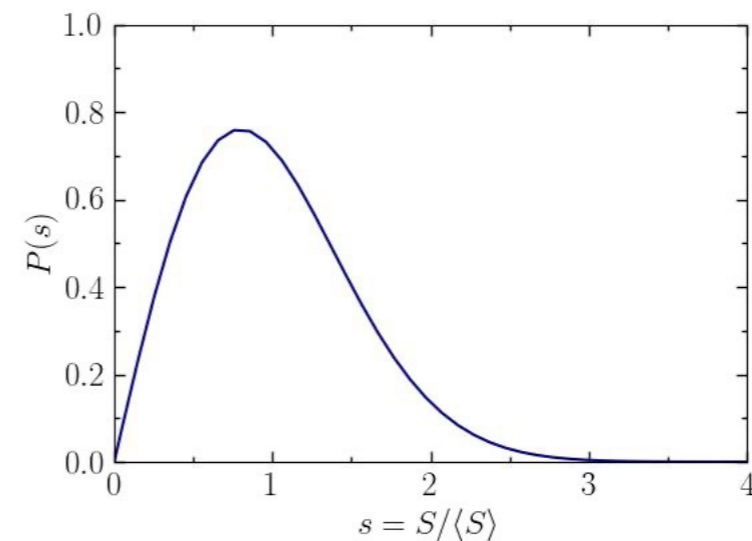
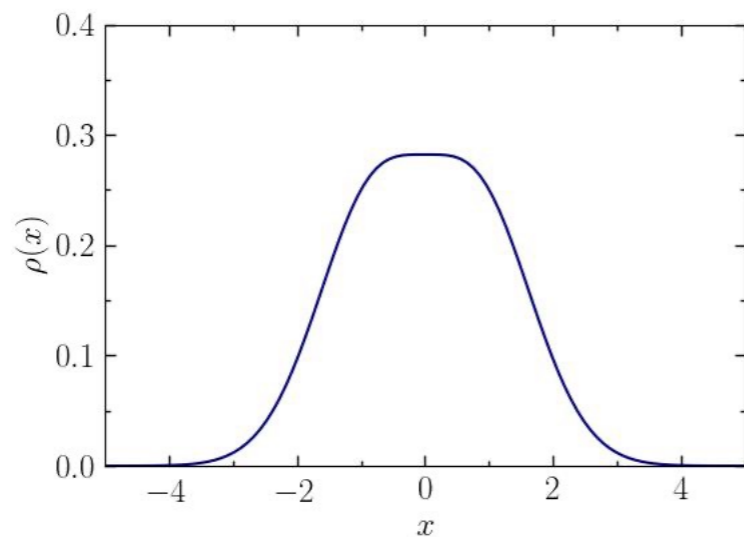
Universal scaling
(SGD)

Model-specific factor
(Loss function, architecture, etc.)

Experiment

$$\sigma_i^2 = \frac{\alpha}{|\mathcal{B}|} \frac{\tilde{g}_i^2}{\Omega_i}$$

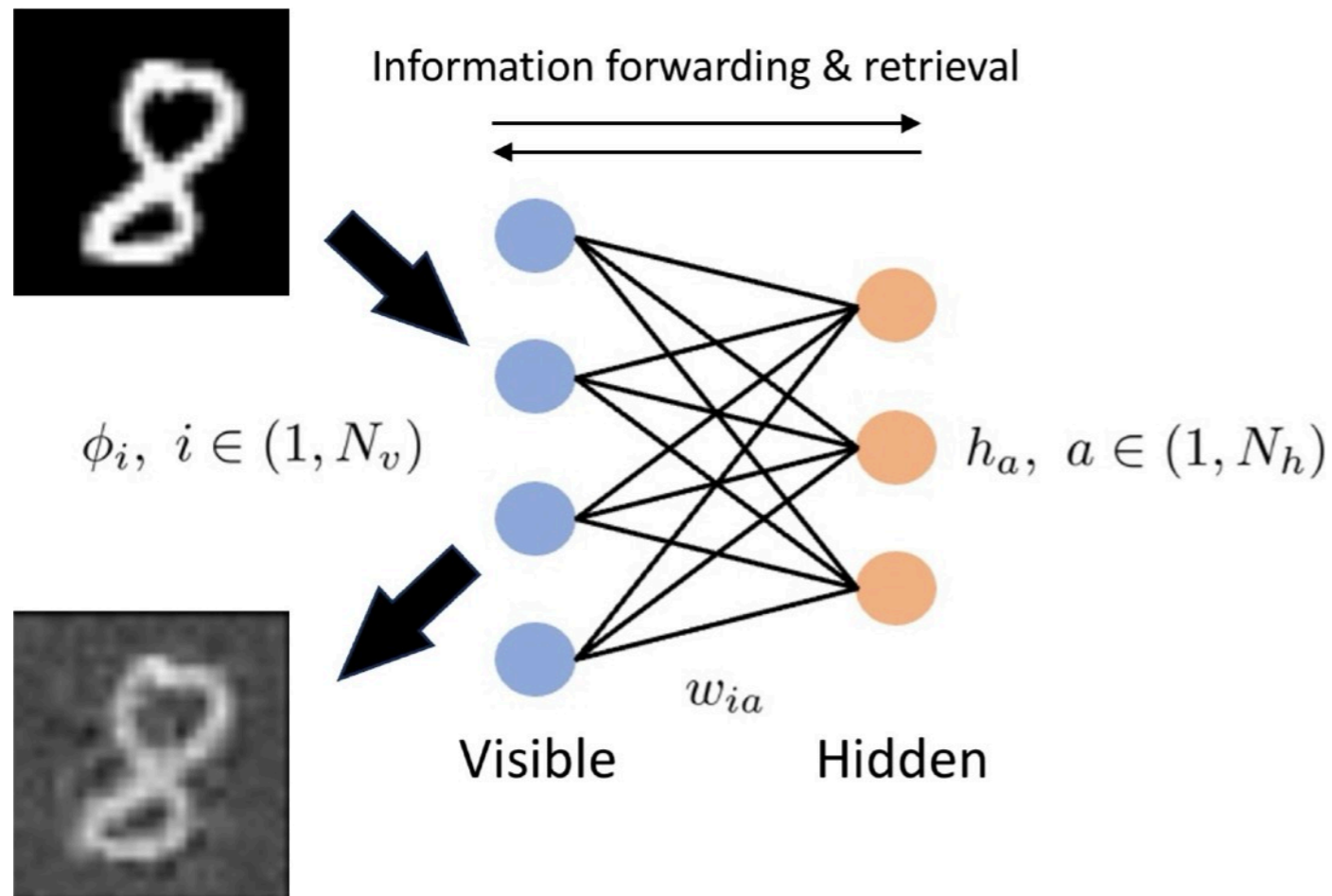
1. Ensembles of the model are trained with different values of α and $|\mathcal{B}|$.
2. Random matrix behaviour of the stochastic optimisation is checked by comparing eigenvalue distributions to Wigner surmise and semi-circle.



3. Universal scaling factor is checked by observing the scaling behaviour of Wigner surmise and semi-circle according to different values of α and $|\mathcal{B}|$.

Scalar field RBM

Gaussian Restricted Boltzmann Machine



- RBM is an energy-based generative model.
- Each layer is sampled based on the energy function.

$$p(\phi, h) = \frac{1}{Z} e^{-H(\phi, h)} \quad \text{with an energy function} \quad H(\phi, h).$$

Learning scalar field theory

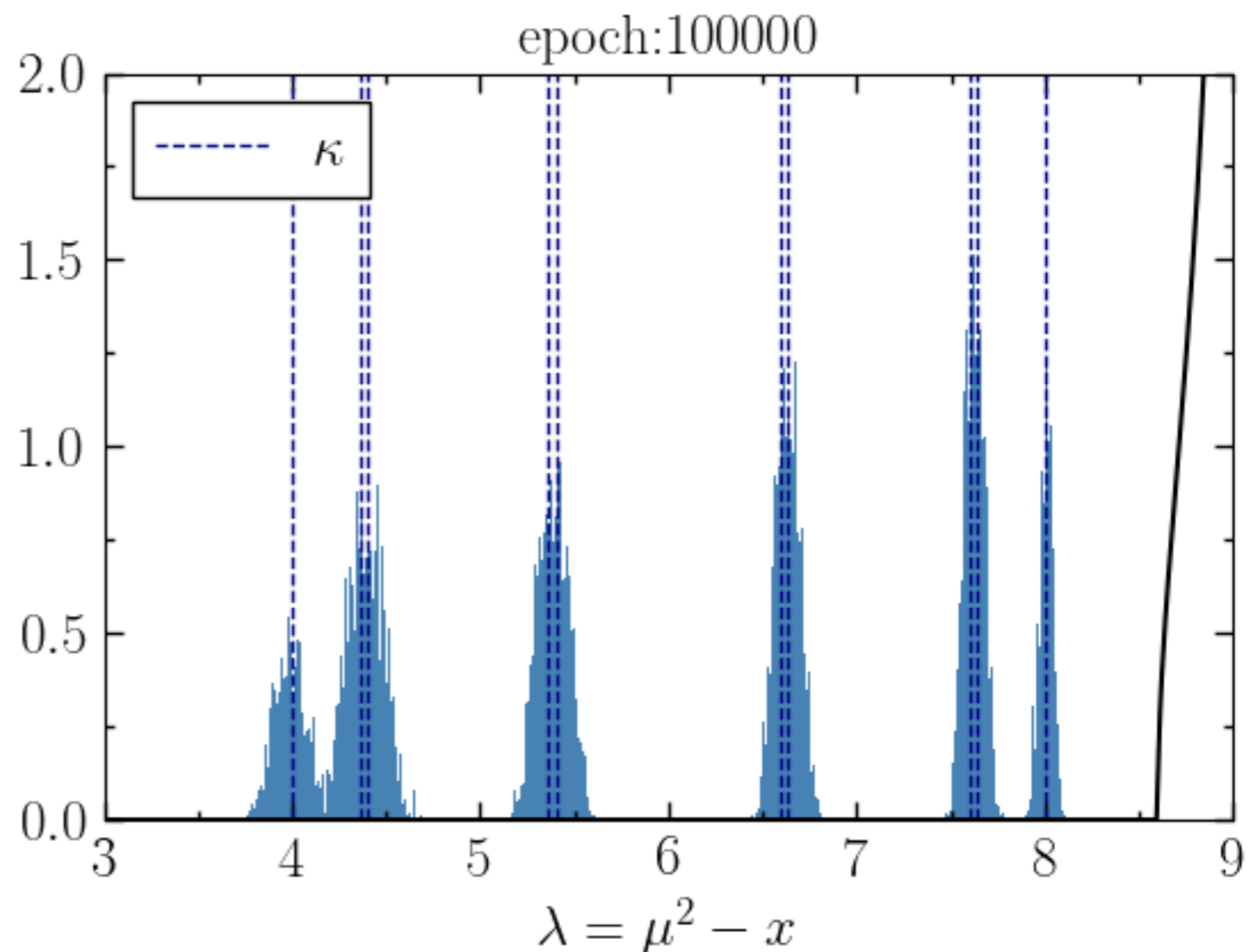
RBM is trained to learn a distribution representing a lattice scalar field theory.

Target distribution eigenvalue: $\kappa_i = m^2 + 2 - 2 \cos \left(\frac{2\pi i}{N} \right)$

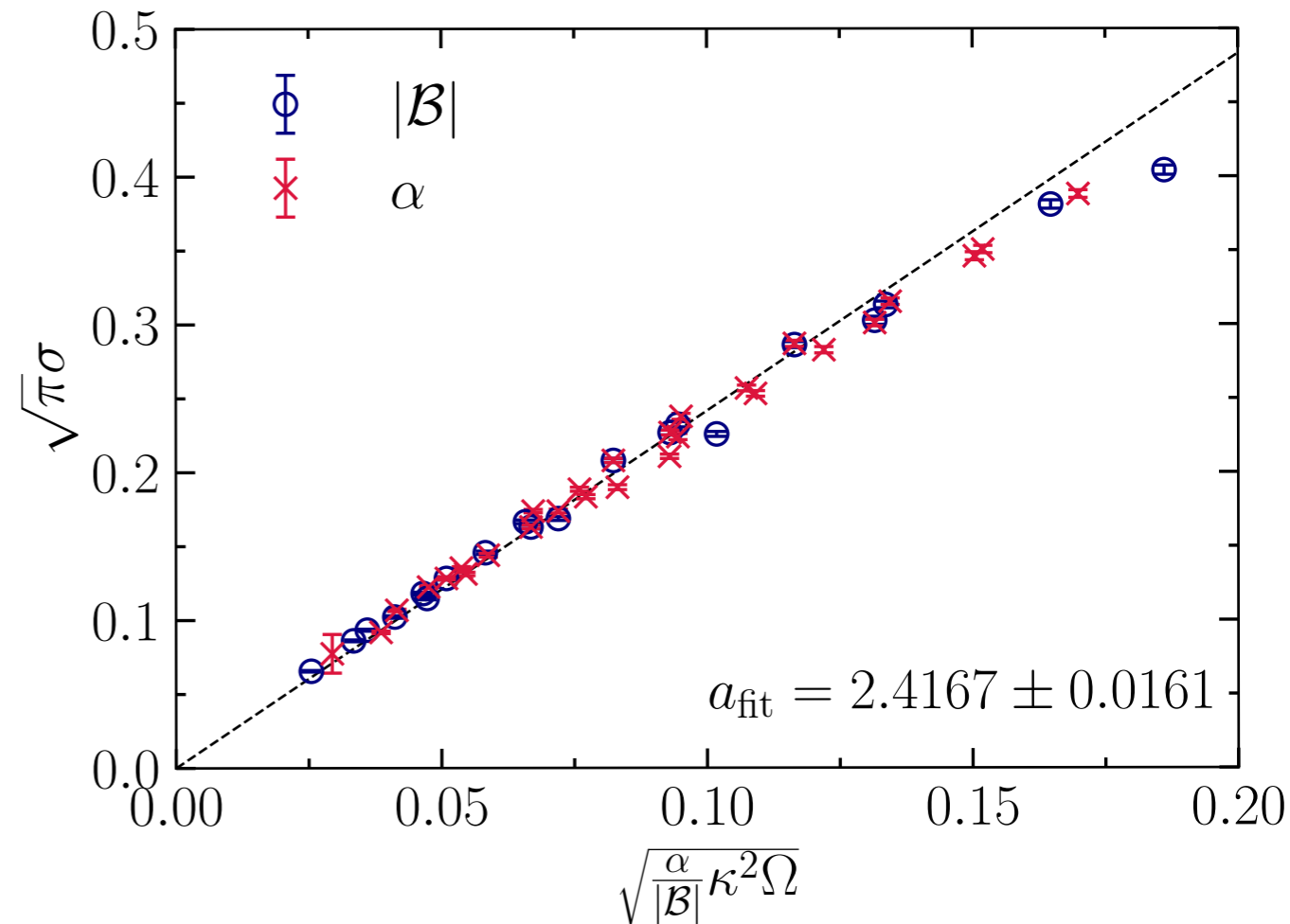
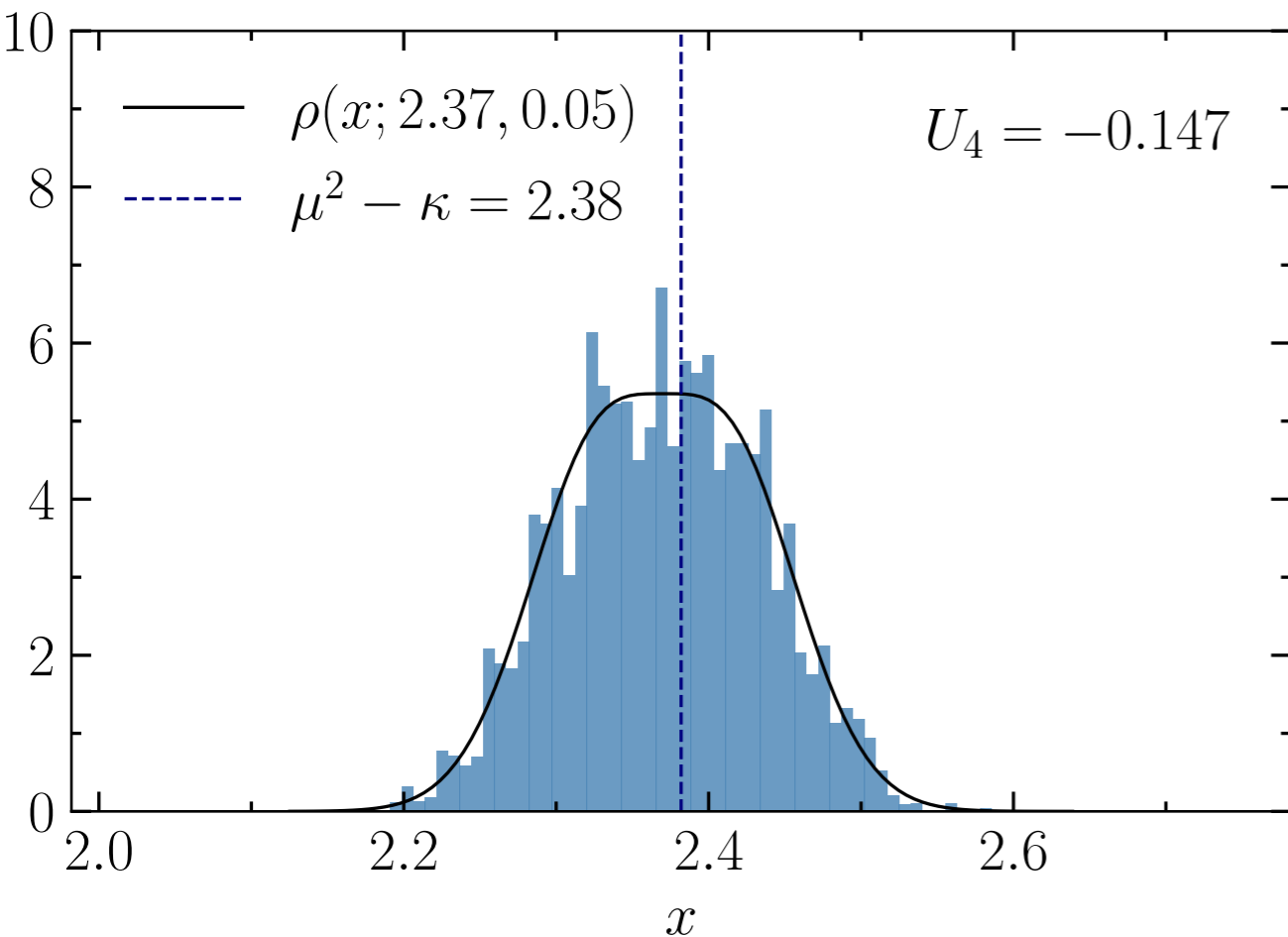
RBM eigenvalue distribution:

Gradient (drift) of Scalar field RBM:

$$\frac{\partial \mathcal{L}}{\partial M_{ii}} = K_i(x_i) = \left(\frac{1}{\kappa_i} - \frac{1}{\mu^2 - x_i} \right) x_i$$



Spectral density

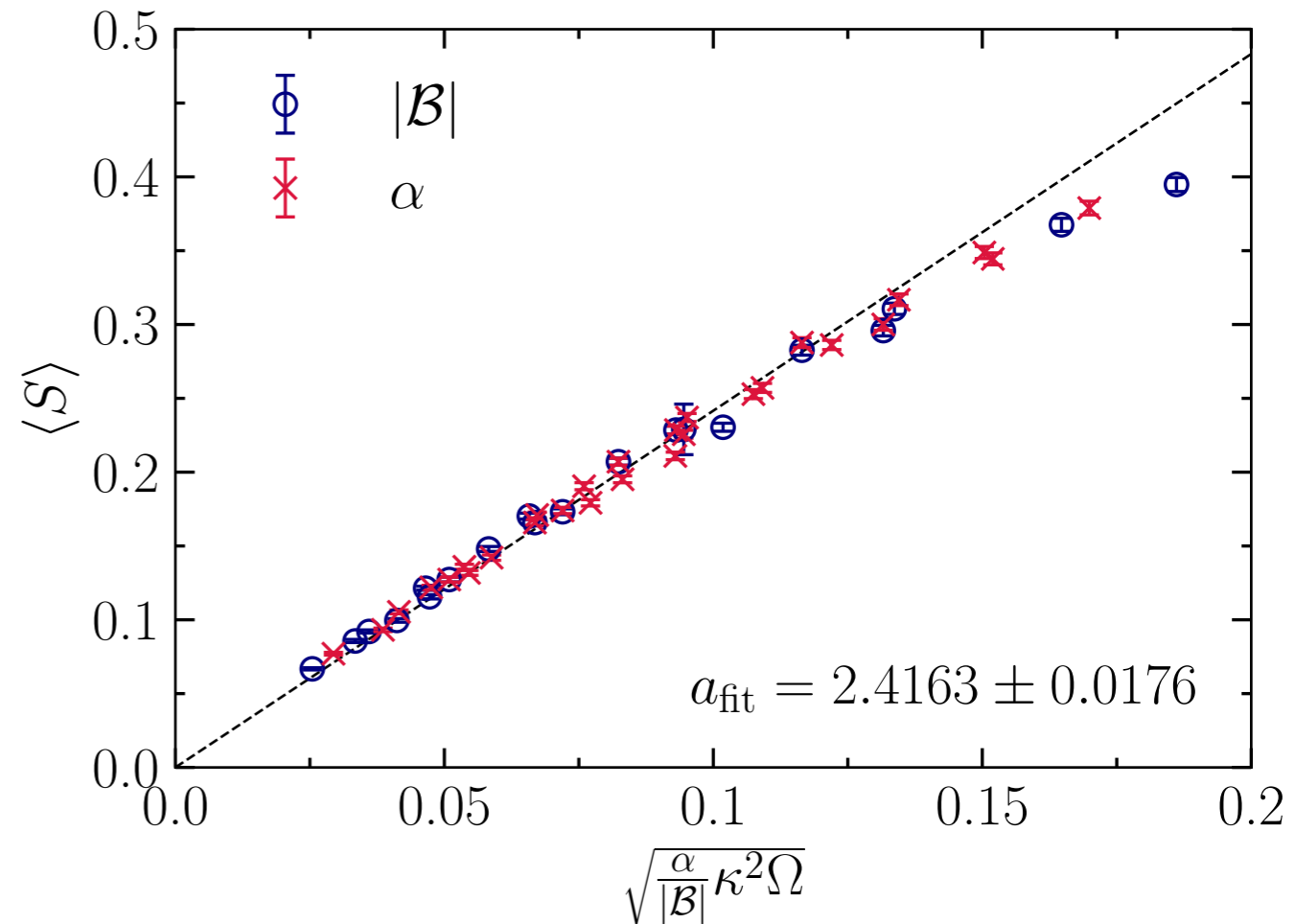
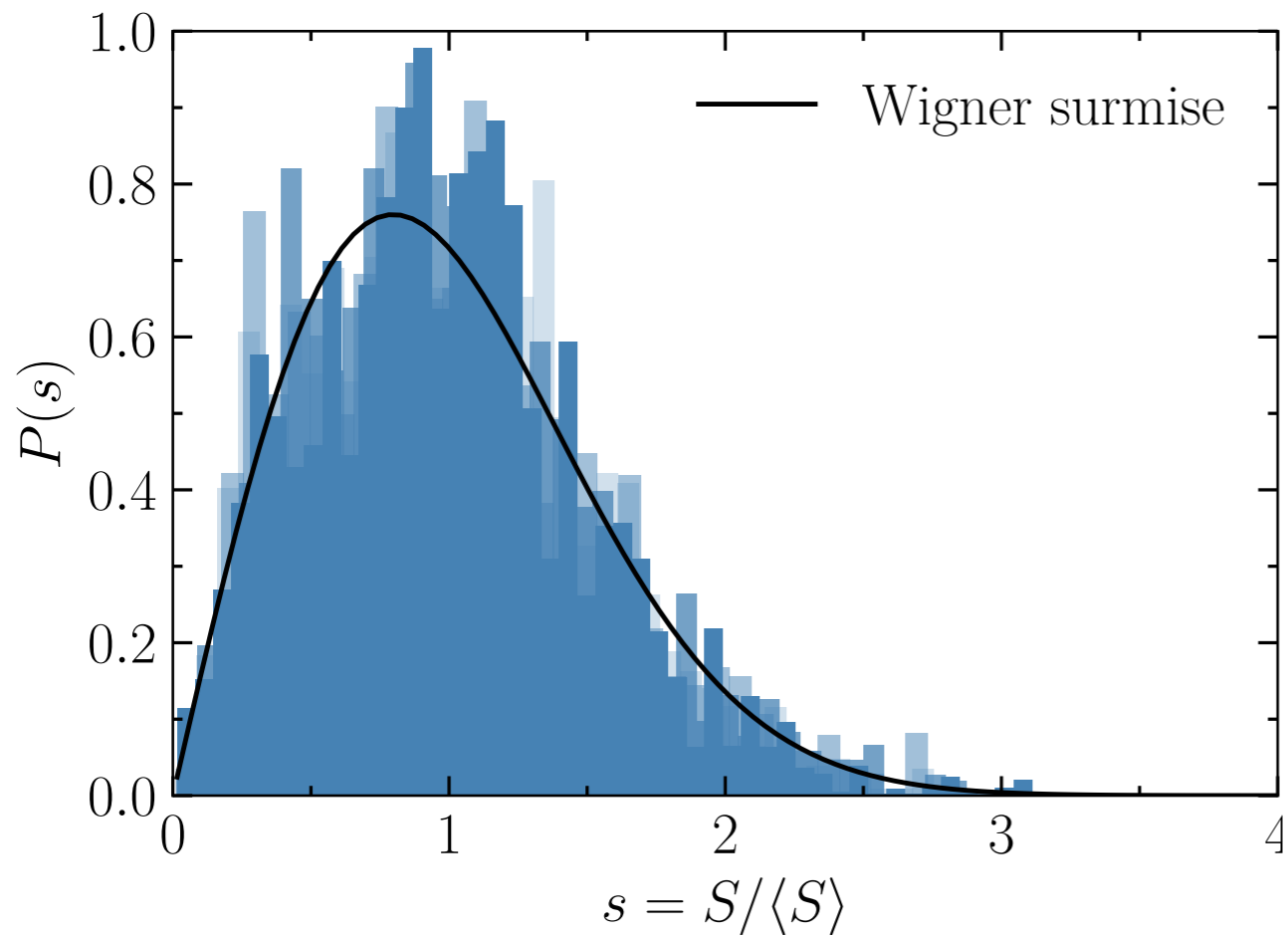


Eigenvalue distribution follows the Wigner semi-circle.

$$U_4 \equiv \frac{\langle \delta x^4 \rangle}{3 \langle \delta x^2 \rangle^2} - 1 = -\frac{4}{27} \approx -0.147\dots$$

The width of the distribution scales with the universal scaling factor $\alpha/|\mathcal{B}|$.

Level spacing

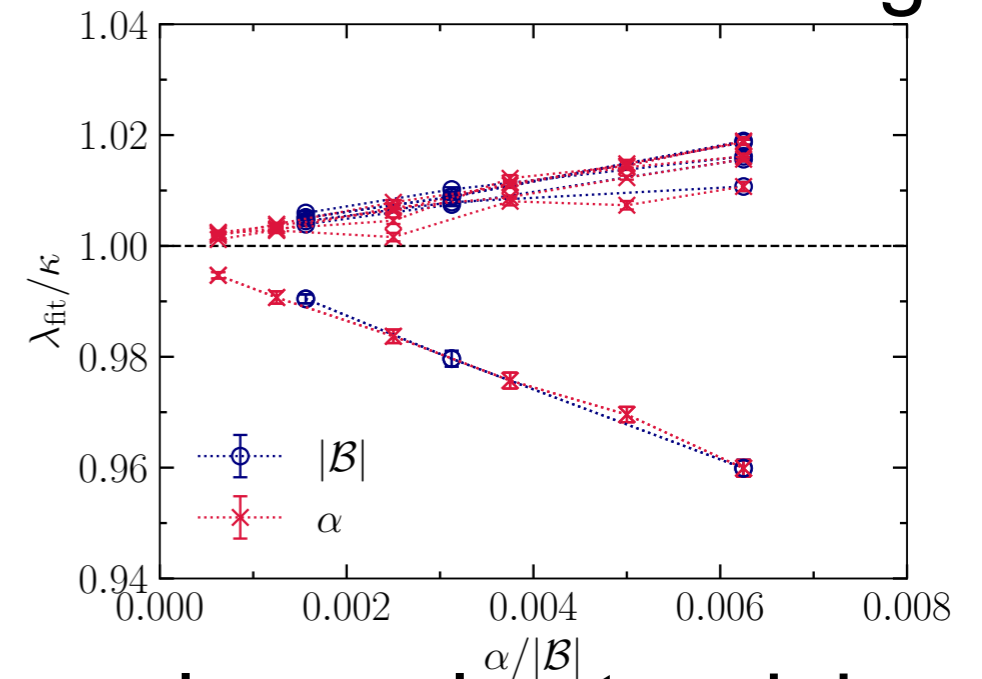


Mean level spacing collapses into the universal curve.

Level spacing scales with the universal scaling factor $\alpha/|\mathcal{B}|$.

Summary and Outlook

- The Linear Scaling Rule of the stochastic gradient descent can be derived from the random matrix theory.
- Stochasticity of the model scales with the universal scaling factor $\alpha/|\mathcal{B}|$.
- Training error (precision) is bounded below by the fluctuation.
- Extending the experiment to the general neural network is in process. ... see Matteo's poster



More nice plots

