
Supplemental Material:

Mean-field theory of graph neural networks in graph partitioning

Tatsuro Kawamoto, Masashi Tsubaki

Artificial Intelligence Research Center,
National Institute of Advanced Industrial Science and Technology,
2-3-26 Aomi, Koto-ku, Tokyo, Japan
{kawamoto.tatsuro, tsubaki.masashi}@aist.go.jp

Tomoyuki Obuchi

Department of Mathematical and Computing Science, Tokyo Institute of Technology,
2-12-1 Ookayama Meguro-ku Tokyo, Japan
obuchi@c.titech.ac.jp

A Derivation of the self-consistent equation

In this section, the detailed derivation of the self-consistent equation of the covariance matrix is derived. Here, we recast our starting-point equation:

$$1 = \int \mathcal{D}\hat{\mathbf{x}}^{t+1} \mathcal{D}\mathbf{x}^{t+1} e^{-\mathcal{L}_0} \langle e^{\mathcal{L}_1} \rangle_{A, W^t, X^t}, \quad \left\{ \begin{array}{l} \mathcal{L}_0 = \sum_{\sigma\mu} \gamma_{\sigma} \hat{x}_{\sigma\mu}^{t+1} x_{\sigma\mu}^{t+1} \\ \mathcal{L}_1 = \frac{1}{N} \sum_{\mu\nu} \sum_{ij} A_{ij} W_{\nu\mu}^t \hat{x}_{\sigma_i\mu}^{t+1} \phi(x_{j\nu}^t) \end{array} \right\}, \quad (1)$$

A.1 Random averages over W^t and A

We first take the average of $\exp(\mathcal{L}_1)$ over W^t . The Gaussian integral with respect to W^t yields

$$\begin{aligned} \langle e^{\mathcal{L}_1} \rangle_{W^t} &= \exp \left(\frac{1}{2DN^2} \sum_{\mu\nu} \sum_{ijkl} A_{ij} \hat{x}_{\sigma_i\mu}^{t+1} \phi(x_{j\nu}^t) A_{k\ell} \hat{x}_{\sigma_k\mu}^{t+1} \phi(x_{\ell\nu}^t) \right) \\ &= \exp \left(\frac{1}{2D} \sum_{\sigma_1\sigma_2\sigma_3\sigma_4} \sum_{\mu} \hat{x}_{\sigma_1\mu}^{t+1} \hat{x}_{\sigma_3\mu}^{t+1} \sum_{\nu} \psi_{\sigma_1\sigma_2}^{t,\nu} \psi_{\sigma_3\sigma_4}^{t,\nu} \right) = \exp \left(\frac{D}{2} \sum_{\sigma_1\sigma_3} u_{\sigma_1\sigma_3}^{t+1} v_{\sigma_1\sigma_3}^t \right), \end{aligned} \quad (2)$$

where we introduce the following quantities:

$$\psi_{\sigma\sigma'}^{t,\nu} \equiv \frac{1}{N} \sum_{i \in V_{\sigma}} \sum_{j \in V_{\sigma'}} A_{ij} \phi(x_{j\nu}^t), \quad (3)$$

$$u_{\sigma\sigma'}^{t+1} \equiv \frac{1}{D} \sum_{\mu} \hat{x}_{\sigma\mu}^{t+1} \hat{x}_{\sigma'\mu}^{t+1}, \quad (4)$$

$$v_{\sigma\sigma'}^t \equiv \frac{1}{D} \sum_{\mu} \sum_{\tilde{\sigma}\tilde{\sigma}'} \psi_{\sigma\tilde{\sigma}}^{t,\mu} \psi_{\tilde{\sigma}'\sigma'}^{t,\mu}. \quad (5)$$

Therefore, Eq. (1) can be written as

$$\begin{aligned}
1 &= \int \mathcal{D}\hat{\mathbf{x}}^{t+1} \mathcal{D}\mathbf{x}^{t+1} \left\langle \int \mathcal{D}\hat{\mathbf{u}}^{t+1} \mathcal{D}\mathbf{u}^{t+1} \int \mathcal{D}\hat{\mathbf{v}}^t \mathcal{D}\mathbf{v}^t \int \mathcal{D}\hat{\psi}^t \mathcal{D}\psi^t \right. \\
&\quad \times \left\langle \exp \left[-\mathcal{L}_0 + \frac{D}{2} \sum_{\sigma\sigma'} u_{\sigma\sigma'}^{t+1} v_{\sigma\sigma'}^t - \sum_{\sigma\sigma'} \hat{u}_{\sigma\sigma'}^{t+1} \left(D u_{\sigma\sigma'}^{t+1} - \sum_{\mu} \hat{x}_{\sigma\mu}^{t+1} \hat{x}_{\sigma'\mu}^{t+1} \right) \right. \right. \\
&\quad \left. \left. - \sum_{\sigma\sigma'} \hat{v}_{\sigma\sigma'}^t \left(D v_{\sigma\sigma'}^t - \sum_{\nu} \sum_{\tilde{\sigma}\tilde{\sigma}'} \psi_{\sigma\tilde{\sigma}}^{t,\nu} \psi_{\sigma'\tilde{\sigma}'}^{t,\nu} \right) \right. \right. \\
&\quad \left. \left. - \sum_{\sigma\sigma'} \sum_{\mu} \hat{\psi}_{\sigma\sigma'}^{t,\mu} \left(\psi_{\sigma\sigma'}^{t,\mu} - \frac{1}{N} \sum_{i \in V_{\sigma}} \sum_{j \in V_{\sigma'}} A_{ij} \phi(x_{j\mu}^t) \right) \right] \right\rangle_A \right\rangle_{\mathbf{X}^t}. \quad (6)
\end{aligned}$$

Note that as we will see below, \mathbf{u}^{t+1} , \mathbf{v}^t , ψ^t , and their conjugates are related to \mathbf{X}^t , and thus the average over \mathbf{X}^t is taken outside of their integral.

We next take the average over a random graph. In Eq. (6), only the final term in the exponent is relevant to \mathbf{A} . We denote this term as \mathcal{L}_2 . We also let $\Xi_{ij}^t = \sum_{\mu} \left(\hat{\psi}_{\sigma_i\sigma_j}^{t,\mu} \phi(x_{j\mu}^t) + \hat{\psi}_{\sigma_j\sigma_i}^{t,\mu} \phi(x_{i\mu}^t) \right) / N$. Because the graph is generated from the SBM, we have that

$$\langle e^{\mathcal{L}_2} \rangle_A = \prod_{i < j} \left[\sum_{A_{ij} \in \{0,1\}} \left(\rho_{\sigma_i\sigma_j} e^{\Xi_{ij}^t} \right)^{A_{ij}} (1 - \rho_{\sigma_i\sigma_j})^{1-A_{ij}} \right] \quad (7)$$

$$\approx \exp \left[\sum_{i < j} \rho_{\sigma_i\sigma_j} \left(e^{\Xi_{ij}^t} - 1 \right) \right] \quad (8)$$

$$\approx \exp \left[\sum_{i < j} \rho_{\sigma_i\sigma_j} \Xi_{ij}^t \right] \quad (9)$$

$$= \exp \left[\sum_{\mu} \sum_{\sigma\sigma'} \gamma_{\sigma} \rho_{\sigma\sigma'} \hat{\psi}_{\sigma\sigma'}^{t,\mu} \sum_{j \in V_{\sigma'}} \phi(x_{j\mu}^t) \right]. \quad (10)$$

At the second line, we used the fact that $\rho_{\sigma\sigma'} = O(N^{-1})$. Then, at the third line we used $\Xi_{ij}^t = O(D/N) \ll 1$. Finally, at the last line we used the symmetry of the undirected graph, $\rho_{\sigma\sigma'} = \rho_{\sigma'\sigma}$.

Note here that the degrees of freedom with respect to the feature dimension are factored out, and thus the dependence on μ can be omitted. Hereafter, the same notation will be employed for the variables without the μ -dependence. We also introduce the notation $\exp^D(f) \equiv \exp(Df)$. The factor inside of the average over \mathbf{X}^t in Eq. (6) can be written as follows:

$$\begin{aligned}
&\int \mathcal{D}\hat{\mathbf{u}}^{t+1} \mathcal{D}\mathbf{u}^{t+1} \int \mathcal{D}\hat{\mathbf{v}}^t \mathcal{D}\mathbf{v}^t \int \mathcal{D}\hat{\psi}^t \mathcal{D}\psi^t e^{\mathcal{L}^*(\hat{\mathbf{u}}^{t+1})} \times \exp^D \left[\sum_{\sigma\sigma'} B_{\sigma\sigma'} \hat{\psi}_{\sigma\sigma'}^t \frac{1}{\gamma_{\sigma'} N} \sum_{j \in V_{\sigma'}} \phi(x_j^t) \right. \\
&\quad \left. + \sum_{\sigma\sigma'} \left(\frac{1}{2} u_{\sigma\sigma'}^{t+1} v_{\sigma\sigma'}^t - \hat{u}_{\sigma\sigma'}^{t+1} u_{\sigma\sigma'}^{t+1} - \hat{v}_{\sigma\sigma'}^t v_{\sigma\sigma'}^t - \hat{\psi}_{\sigma\sigma'}^t \psi_{\sigma\sigma'}^t - \hat{v}_{\sigma\sigma'}^t \sum_{\tilde{\sigma}\tilde{\sigma}'} \psi_{\sigma\tilde{\sigma}}^t \psi_{\sigma'\tilde{\sigma}'}^t \right) \right], \quad (11)
\end{aligned}$$

where

$$\mathcal{L}^*(\hat{\mathbf{u}}^{t+1}) = \log \int \mathcal{D}\hat{\mathbf{x}}^{t+1} \mathcal{D}\mathbf{x}^{t+1} \exp^D \left(- \sum_{\sigma} \gamma_{\sigma} \hat{x}_{\sigma}^{t+1} x_{\sigma}^{t+1} + \sum_{\sigma\sigma'} \hat{u}_{\sigma\sigma'}^{t+1} \hat{x}_{\sigma}^{t+1} \hat{x}_{\sigma'}^{t+1} \right). \quad (12)$$

As in the main text, we have defined $B_{\sigma\sigma'} \equiv N \gamma_{\sigma} \rho_{\sigma\sigma'} \gamma_{\sigma'}$.

When $D \gg 1$, the saddle-point condition of the exponent in Eq. (11) yields $u_{\sigma\sigma'}^{t+1} = \langle \hat{x}_{\sigma}^{t+1} \hat{x}_{\sigma'}^{t+1} \rangle_{\mathcal{L}^*}$, $v_{\sigma\sigma'}^t = \psi_{\sigma}^t \psi_{\sigma'}^t$, $\psi_{\sigma\sigma'}^t = B_{\sigma\sigma'} (\gamma_{\sigma'} N)^{-1} \sum_{j \in V_{\sigma'}} \phi(x_j^t)$, $\hat{u}_{\sigma\sigma'}^{t+1} = v_{\sigma\sigma'}^t / 2$, $\hat{v}_{\sigma\sigma'}^t = u_{\sigma\sigma'}^{t+1} / 2$, and $\hat{\psi}_{\sigma\sigma'}^t =$

$\sum_{\tilde{\sigma}} (\hat{v}_{\sigma\tilde{\sigma}}^t + \hat{v}_{\tilde{\sigma}\sigma}^t) \psi_{\sigma\tilde{\sigma}}^t$, where $\psi_{\sigma}^t \equiv \sum_{\tilde{\sigma}} \psi_{\sigma\tilde{\sigma}}^t$, and $\langle \dots \rangle_{\mathcal{L}^*}$ is the average taken with the weight of the integrand of Eq. (12). Because the correlation between the auxiliary variables should be zero, owing to causality [1, 2], we finally arrive at

$$1 = \int \mathcal{D}\hat{\mathbf{x}}^{t+1} \mathcal{D}\mathbf{x}^{t+1} \exp \left(- \sum_{\sigma} \gamma_{\sigma} \hat{x}_{\sigma}^{t+1} x_{\sigma}^{t+1} \right) \left\langle \exp \left(\frac{1}{2} \sum_{\sigma\sigma'} \hat{x}_{\sigma}^{t+1} F_{\sigma\sigma'}(\mathbf{X}^t) \hat{x}_{\sigma'}^{t+1} \right) \right\rangle_{\mathbf{X}^t}, \quad (13)$$

where $F_{\sigma\sigma'}(\mathbf{X}^t)$ is defined as

$$F_{\sigma\sigma'}(\mathbf{X}^t) \equiv \sum_{\tilde{\sigma}\tilde{\sigma}'} B_{\sigma\tilde{\sigma}} B_{\sigma'\tilde{\sigma}'} \frac{1}{\gamma_{\tilde{\sigma}} N} \frac{1}{\gamma_{\tilde{\sigma}'} N} \sum_{k \in V_{\tilde{\sigma}}} \sum_{\ell \in V_{\tilde{\sigma}'}} \phi(x_k^t) \phi(x_{\ell}^t). \quad (14)$$

A.2 Stochastic process with a correlated noise

Here, we compare Eq. (13) with a Markovian discrete-time stochastic process $y_{\sigma}^{t+1} = \eta_{\sigma}^t$, in which each element is correlated via a random noise, i.e., $\langle \eta_{\sigma}^t \rangle_{\eta} = 0$, and $\langle \eta_{\sigma}^t \eta_{\sigma'}^t \rangle_{\eta} = C_{\sigma\sigma'}$ for any t . The corresponding normalization condition reads

$$\begin{aligned} 1 &= \int \prod_{\sigma} dy_{\sigma}^{t+1} \left\langle \prod_{\sigma} \delta(y_{\sigma}^{t+1} - \eta_{\sigma}^t) \right\rangle_{\eta} = \int \prod_{\sigma} \gamma_{\sigma} \frac{d\hat{y}_{\sigma}^{t+1} dy_{\sigma}^{t+1}}{2\pi i} \left\langle e^{-\sum_{\sigma} \gamma_{\sigma} \hat{y}_{\sigma}^{t+1} (y_{\sigma}^{t+1} - \eta_{\sigma}^t)} \right\rangle_{\eta} \\ &= \int \mathcal{D}\hat{\mathbf{y}}^{t+1} \mathcal{D}\mathbf{y}^{t+1} e^{-\sum_{\sigma} \gamma_{\sigma} \hat{y}_{\sigma}^{t+1} y_{\sigma}^{t+1}} \int \prod_{\sigma} \frac{d\eta_{\sigma}^t}{2\pi i} \exp \left(-\frac{1}{2} \sum_{\sigma\sigma'} \eta_{\sigma}^t C_{\sigma\sigma'}^{-1} \eta_{\sigma'}^t + \sum_{\sigma} \gamma_{\sigma} \hat{y}_{\sigma}^{t+1} \eta_{\sigma}^t \right) \\ &= \int \mathcal{D}\hat{\mathbf{y}}^{t+1} \mathcal{D}\mathbf{y}^{t+1} \exp \left(-\sum_{\sigma} \gamma_{\sigma} \hat{y}_{\sigma}^{t+1} y_{\sigma}^{t+1} + \frac{1}{2} \sum_{\sigma\sigma'} \hat{y}_{\sigma}^{t+1} \gamma_{\sigma} C_{\sigma\sigma'} \gamma_{\sigma'} \hat{y}_{\sigma'}^{t+1} \right). \end{aligned} \quad (15)$$

Analogously to the case of the GNN, we have defined $\mathcal{D}\hat{\mathbf{y}}^{t+1} \mathcal{D}\mathbf{y}^{t+1} \equiv \prod_{\sigma} \gamma_{\sigma} d\hat{y}_{\sigma}^{t+1} dy_{\sigma}^{t+1} / 2\pi i$.

A.3 Self-consistent equation

Finally, we compare Eqs. (13) and (15). However, note that these are not of exactly the same form, because the average over \mathbf{X}^t is taken outside of the exponential in Eq. (13). Two approximations are made in order to derive the self-consistent equation, and the assumptions that justify these approximations are discussed afterward.

First, if the approximation

$$\left\langle \exp \left(\frac{1}{2} \sum_{\sigma\sigma'} \hat{x}_{\sigma}^{t+1} F_{\sigma\sigma'}(\mathbf{X}^t) \hat{x}_{\sigma'}^{t+1} \right) \right\rangle_{\mathbf{X}^t} \approx \exp \left(\frac{1}{2} \sum_{\sigma\sigma'} \hat{x}_{\sigma}^{t+1} \langle F_{\sigma\sigma'}(\mathbf{X}^t) \rangle_{\mathbf{X}^t} \hat{x}_{\sigma'}^{t+1} \right) \quad (16)$$

holds in the stationary limit, then the group-wise state \mathbf{x}^t can be regarded as a Gaussian variable whose correlation matrix obeys

$$C_{\sigma\sigma'} = \frac{1}{\gamma_{\sigma} \gamma_{\sigma'}} \sum_{\sigma_1 \sigma_2} B_{\sigma\sigma_1} B_{\sigma'\sigma_2} \left\langle \frac{1}{\gamma_{\sigma_1} N \gamma_{\sigma_2} N} \sum_{i \in V_{\sigma_1}} \sum_{j \in V_{\sigma_2}} \phi(x_i) \phi(x_j) \right\rangle_{\mathbf{X}^t}. \quad (17)$$

This equation is still not closed, because the right-hand side of Eq. (17) depends on the statistic of \mathbf{X}^t , rather than \mathbf{x}^t . However, because the vertices within the group σ are statistically equivalent, $\{x_i\}_{i \in V_{\sigma}}$ are expected to obey the same distribution with mean \mathbf{x}_{σ} , which itself is a random variable. If $\sum_{i \in V_{\sigma}} \phi(x_i) / (\gamma_{\sigma} N) \approx \phi(\mathbf{x}_{\sigma})$ holds, then the right-hand side of Eq. (17) can be evaluated as the average with respect to the group-wise variable \mathbf{x}^t . Then, within this regime we arrive at the following self-consistent equation with respect to the covariance matrix $\mathbf{C} = [C_{\sigma\sigma'}]$:

$$C_{\sigma\sigma'} = \frac{1}{\gamma_{\sigma} \gamma_{\sigma'}} \sum_{\tilde{\sigma}\tilde{\sigma}'} B_{\sigma\tilde{\sigma}} B_{\sigma'\tilde{\sigma}'} \int \frac{d\mathbf{x} e^{-\frac{1}{2} \mathbf{x}^T \mathbf{C}^{-1} \mathbf{x}}}{(2\pi)^{\frac{N}{2}} \sqrt{\det \mathbf{C}}} \phi(\mathbf{x}_{\tilde{\sigma}}) \phi(\mathbf{x}_{\tilde{\sigma}'}). \quad (18)$$

Let us consider the first approximation that we adopted in Eq. (16). In the terminology of physics, this is the replacement of a free energy with an internal energy, or the neglect of the entropic contribution. It is difficult to evaluate this residual in general. However, note that this becomes closer to equality as every x_i approaches the same value. Therefore, this implies that the self-consistent equation is more accurate as we approach the detectability limit, and yields an adequate estimate of the critical value.

Let us next consider the second approximation we adopted in Eq. (17). Although the law of large numbers with respect to $\phi(x_i)$ (not x_i) ensures that $\sum_{i \in V_\sigma} \phi(x_i)/(\gamma_\sigma N)$ has a certain value characterized by the group, this may be different from $\phi(x_\sigma)$. In fact, the relation between these is in general an inequality (Jensen's inequality) when the activation function ϕ is a convex function. The (exact) equality holds only when $\{x_i\}$ is constant or the function ϕ is linear within the considered domain.

The second approximation can be justified in the following cases. The first case is when the fluctuation of $x_i - x_{\sigma_i}$ is negligible compared to the magnitude of x_σ . Note that this is the same assumption as we made in the first approximation. To see this precisely, let us express x_i as $x_i = x_\sigma + z_i$ for $i \in V_\sigma$. We can formally write the probability distribution $P(\{x_i\})$ of $\{x_i\}$ in a hierarchical fashion as follows:

$$P(\{x_i\}) = \int \prod_{\sigma} dx_{\sigma} \int \prod_i dz_i P_{\sigma}(\{x_{\sigma}\}) P(\{x_i\}) \prod_i \delta(z_i - x_i + x_{\sigma_i}), \quad (19)$$

where $P_{\sigma}(\{x_{\sigma}\})$ is the probability distribution with respect to \mathbf{x} . Thus, the expectation $\langle f(\mathbf{X}) \rangle_{\mathbf{X}}$ can be expressed as

$$\begin{aligned} \langle f(\mathbf{X}) \rangle_{\mathbf{X}} &\equiv \int \prod_i dx_i P(\{x_i\}) f(\{x_i\}) \\ &= \int \prod_{\sigma} dx_{\sigma} \int \prod_i dz_i P(\{z_i\}|\{x_{\sigma}\}) P_{\sigma}(\{x_{\sigma}\}) f(\{x_{\sigma_i} + z_i\}), \end{aligned} \quad (20)$$

where $P(\{z_i\}|\{x_{\sigma}\}) \equiv P(\{x_{\sigma_i} + z_i\})$, which can be a nontrivial function. However, whenever the contributions from the average with respect to z_i are negligible, Eq. (20) implies that the expectation in Eq. (17) can be evaluated using only the group-wise variables $\{x_{\sigma}\}$. Another case is when the activation function ϕ is almost linear within the domain over which z_i fluctuates. For example, in the case that $\phi = \tanh$, the present approximation does not deteriorate the accuracy even when $x_{\sigma} \approx 0$. When either of these assumption holds, the equality of Jensen's inequality is approximately satisfied, and our derivation of the self-consistent equation is justified.

B K-means classification using $\phi(\mathbf{X})$

Instead of \mathbf{X}^T , $\phi(\mathbf{X}^T)$ can be adopted to perform the k-means classification after the feedforward process. Again, we employ \tanh as the nonlinear activation function. The results of an untrained GNN and a trained GNN under the same experimental settings as in the main text are illustrated in Fig. 1 and Fig. 2, respectively. In Fig. 1a, the reader should note that the range of the color gradient is different from that in the phase diagram in the main text. For the untrained GNN, the obtained overlaps are clearly better than that using \mathbf{X}^T . It can be understood that the error is reduced because the nonlinear function drives each element of the state \mathbf{X}^T to either $+1$ or -1 , making the classification using the k-means method easier and more accurate. On the other hand, for the trained GNN, differences between the overlaps using \mathbf{X}^T and $\phi(\mathbf{X}^T)$ are hardly observable.

Particularly for the case of an untrained GNN in which $\phi(\mathbf{X}^T)$ is adopted for the readout classifier, the overlap gradually changes around the estimated detectability limit. This may be as result of the strong finite-size effect. Again, note that our estimate of the detectability limit is for the case that $N \rightarrow \infty$.

References

- [1] H. Sompolinsky and Annette Zippelius. Relaxational dynamics of the edwards-anderson model and the mean-field theory of spin-glasses. *Phys. Rev. B*, 25:6860–6875, Jun 1982.

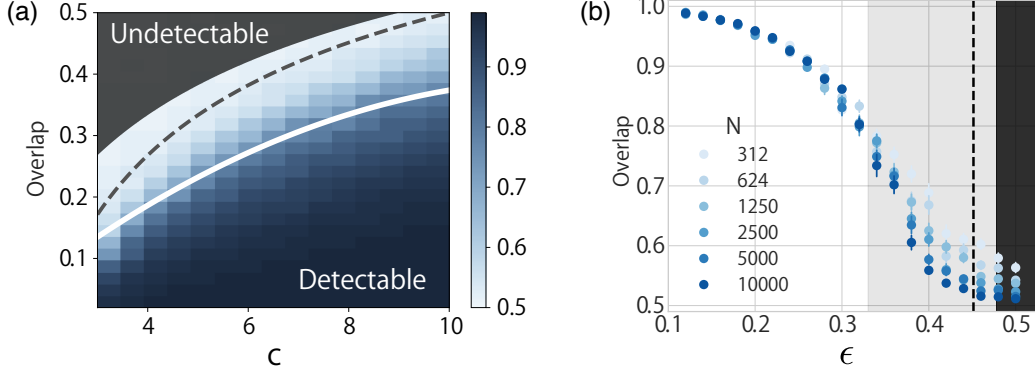


Figure 1: Performance of the untrained GNN using the k-means classifier with $\phi(\mathbf{X})$. (a) The detectability phase diagram and (b) the overlaps of the SBM with $c = 8$ are plotted in the same manner as in Fig. 3 in the main text. When the variation of the overlap is interpolated for each graph size, these curves are crossed at $\epsilon^* \approx 0.33$. It implies the presence of detectability phase transition around the value of ϵ predicted by our mean-field estimate.

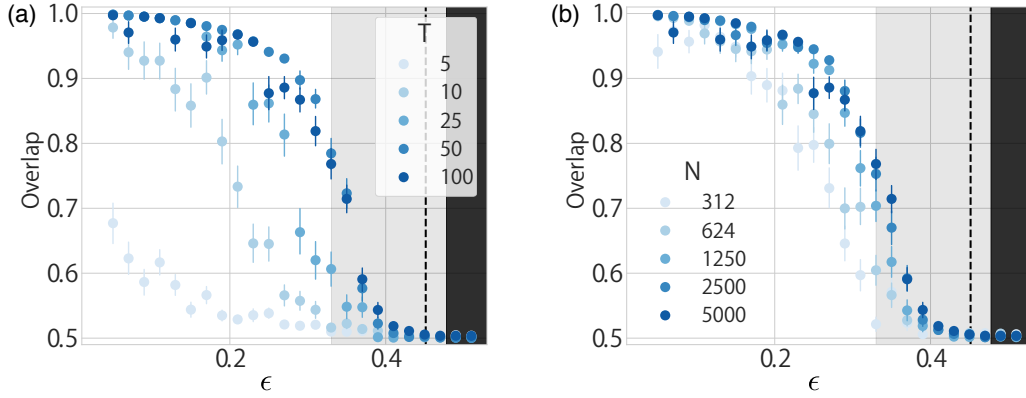


Figure 2: Performance of the trained GNN using the classifier with $\phi(\mathbf{X})$. For (a) and (b), the overlaps are plotted in the same manner as in Fig. 4 in the main text, respectively.

[2] A Crisanti, HJ Sommers, and H Sompolinsky. Chaos in neural networks: chaotic solutions. *preprint*, 1990.