

# Unsupervised machine learning of quenched gauge symmetries: A proof-of-concept demonstration

Daniel Lozano-Gómez \*, Darren Pereira \*, and Michel J. P. Gingras 

*Department of Physics and Astronomy, University of Waterloo, Waterloo, Ontario, Canada N2L 3G1*



(Received 6 December 2021; accepted 28 June 2022; published 18 November 2022)

One of the most prominent tasks of machine learning (ML) methods within the field of condensed matter physics has been to classify phases of matter. Given their many successes in performing this task, one may ask whether these methods—particularly unsupervised ones—can go beyond learning the thermodynamic behavior of a system. This question is especially intriguing when considering systems that have a “hidden order”. In this work we study two random spin systems with a hidden ferromagnetic order that can be exposed by applying a Mattis gauge transformation. We demonstrate that the principal component analysis, perhaps the simplest unsupervised ML method, can detect the hidden order, quantify the corresponding gauge variables, and map the original random models onto simpler gauge-transformed ferromagnetic ones, all *without any prior knowledge* of the underlying gauge transformation. Our work illustrates that ML algorithms can in principle identify not manifestly obvious symmetries of a system.

DOI: [10.1103/PhysRevResearch.4.043118](https://doi.org/10.1103/PhysRevResearch.4.043118)

## I. INTRODUCTION

Machine learning (ML) methods have in recent years proven to be a powerful pattern recognition tool with applications in numerous and varied branches of science. These have shown their ability to extract, identify, and even propose descriptive patterns found in the input data. In condensed matter physics, the application of ML techniques first rose to prominence with the use of the principal component analysis (PCA) method [1] and neural networks [2] by displaying their ability to identify and classify the ferromagnetic and paramagnetic phases of the Ising model. Since then, the use of ML in condensed matter physics has rapidly expanded to include a variety of techniques [3–5]. These and their applications can be broadly grouped into two categories: supervised ML (SML), in which the input data to train the machine is labeled [6–21]; and unsupervised ML (UML), in which the input data is unlabeled and the machine proposes its own classification scheme [10,19–31]. Across these two categories, the identification of thermodynamic phases in various models has remained a central theme. Evidence is accumulating that ML-based learning of phases can be guided by physical insights into a model or system, such as its symmetries. This has been most clearly demonstrated by exploiting properties such as locality and translational symmetry to expedite the learning of convolutional neural networks [2], or by taking advantage of expected symmetry breaking in spin models to extract underlying order parameters for hidden orders [7].

Given the benefits that these physically inspired insights provide, one may ask a question of foremost importance for the broad and growing usage of ML in physics: *can ML guide the identification of hidden or unknown properties of a model with minimal user assistance?* A fitting testing ground for such a question is found in topologically trivial physical models possessing gauge symmetries; in fact, the investigation of gauge-symmetric models with ML protocols is a topic of current interest [16,32,33]. Such models can be simplified by a suitable mathematical transformation, if “one” is *a priori* informed of the appropriate transformation. Our question then becomes a matter of determining whether ML can detect the hidden local gauge symmetry of such models without any prior knowledge. Achieving this would demonstrate that ML can in principle be used to learn the fundamental mathematical details of a model beyond its thermodynamic properties. Although this would be similar in flavor to finding the order parameter of a SU(2) lattice gauge theory [21], it would go beyond such context by demonstrating that ML can determine a gauge transformation to map one *entire model* to another, not simply expose its order parameter (even if hidden). In this vein, and perhaps more interestingly, the controlled mathematical nature of such gauge-symmetric models may suggest their use as a way to probe the inner mechanism of how a ML procedure does truly learn. To explore the aforementioned motivating question, we therefore require (i) a class of models that present themselves as seemingly complex, but which can be much simplified by a gauge transformation, and (ii) a ML method whose classification scheme can be exposed.

In light of (i), we study the Mattis Ising spin glass (MISG) [34,35] and the Mattis XY gauge glass (MXYGG) models [35]. To the “uninformed”, and at a first glance, these two models look very complex: their respective Hamiltonians possess random spin-spin bond interactions with many of the inherent intricacies of random systems. Crucially, a snapshot of their low-temperature ordered state configurations

\*These authors contributed equally to this work.

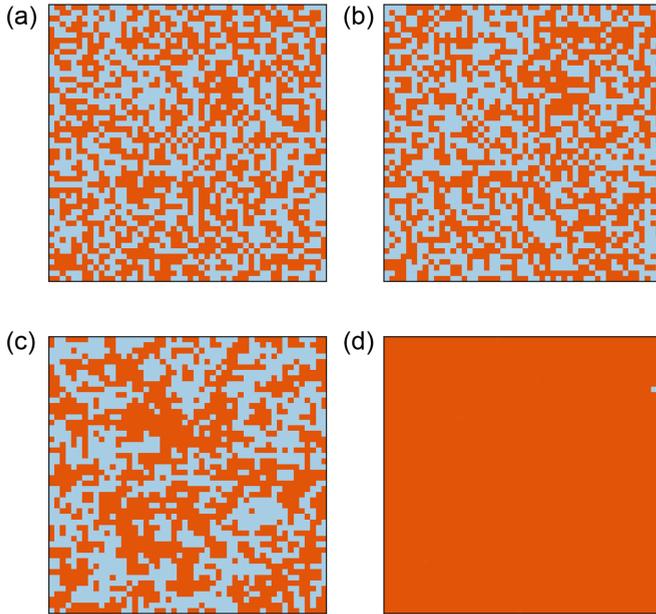


FIG. 1. (a) Disordered (paramagnetic) and (b) ordered spin configurations of the MISG model *before* a gauge transformation; (c) disordered (paramagnetic) and (d) ordered spin configurations of the MISG model *after* a Mattis gauge transformation has been applied, which renders the model equivalent to the ferromagnetic Ising model.

displays no immediately recognizable pattern but, instead, appears completely disordered (paramagnetic), as shown in Figs. 1(a) and 1(b) for the MISG model. However, the MISG and MXYGG models can be transformed into regular ferromagnetic Ising and XY models, respectively, under a (Mattis) gauge transformation [34,35], as shown in Figs. 1(c) and 1(d) for the MISG model. Our general interest is to determine whether ML can learn the mapping from Figs. 1(a) and 1(b) to Figs. 1(c) and 1(d), respectively [36].

Regarding (ii), an important consideration is the trade-off between interpretability and scalability [37,38]. Although a feedforward neural network (FFNN) is able to correctly classify the ordered and disordered configurations of the MISG model, as demonstrated in Appendix A, its lack of interpretability limits the access to arbitrary information it may have learned about the model’s gauge symmetry. With this in mind, we instead use PCA [39], which is highly interpretable, making it a suitable ML method for our exploration. It is also unsupervised, fulfilling our condition of “minimal user assistance”.

Performing PCA on the Mattis models, we provide an affirmative answer to our motivating question. We demonstrate that PCA is able to distinguish between the disordered and ordered phases, and show that the principal components contain excellent approximations of the site-dependent gauge variables that were hidden from PCA. Moreover, these principal components can be used to calculate various quantities that verify the equivalence of the MISG (MXYGG) model with the ferromagnetic Ising (XY) model, confirming that PCA has in essence learned how to map Figs. 1(a) and 1(b) onto Figs. 1(c) and 1(d), respectively.

## II. MODELS

The MISG model [34,35] on a square lattice is defined by the nearest-neighbor Hamiltonian

$$H_I = - \sum_{\langle i,j \rangle} J_{ij} \sigma_i \sigma_j, \quad (1)$$

where  $\sigma_i \equiv \pm 1$ . The couplings  $\{J_{ij}\}$  take values  $\pm J$  randomly ( $J > 0$ ), but with the *imposed constraint* that the product of  $J_{ij}$  couplings around a plaquette is positive:  $P \equiv \prod_{\langle i,j \rangle \in \square} J_{ij} > 0$ . This constraint enforces nonfrustrated plaquettes, and thus nonfrustrated ground states, allowing a Mattis gauge transformation to be applied [34,35]. This transformation reexpresses  $J_{ij}$  as  $J_{ij} \rightarrow \varepsilon_i \varepsilon_j J$ , where  $\{\varepsilon_i\}$  are random site (gauge) variables that take  $\varepsilon_i \equiv \pm 1$  values. Through this transformation, the Hamiltonian (1) becomes  $\tilde{H}_I = -J \sum_{\langle i,j \rangle} \tau_i \tau_j$ , with  $\tau_i \equiv \varepsilon_i \sigma_i \equiv \pm 1$  as the new Ising variables. It is now clear that this gauged system possesses an order parameter given by the Ising model “ $\tau$  magnetization”,  $M \equiv \langle \sum_i \tau_i \rangle = \langle \sum_i \varepsilon_i \sigma_i \rangle$ , illustrating that the MISG model is nothing but an Ising ferromagnetic model in disguise. Further details about this mapping are given in Appendix B.

Similarly, the MXYGG model is described by an XY model with random phase factors  $\{A_{ij}\}$  [40–42],

$$H_{XY} = -J \sum_{\langle i,j \rangle} \cos(\Delta\phi_{ij} - A_{ij}), \quad (2)$$

where  $J > 0$ , and  $\Delta\phi_{ij} = \phi_i - \phi_j$  is the difference between on-site angular variables  $\phi_i \in [0, 2\pi)$ .  $H_{XY}$  is unfrustrated as long as the sum of the random phase factors  $A_{ij}$  around a plaquette is a multiple of  $2\pi$ , i.e.,  $P_{XY} = (\sum_{\langle i,j \rangle \in \square} A_{ij}) \bmod 2\pi = 0$  [35]. Under this condition, a Mattis gauge transformation can be applied by defining random site (gauge) variables  $\{b_i\}$  such that  $A_{ij} = b_j - b_i$ , with  $b_i \in [0, 2\pi)$ .  $H_{XY}$  then becomes  $\tilde{H}_{XY} = -J \sum_{\langle i,j \rangle} \cos(\Delta\theta_{ij})$ , where  $\theta_i \equiv \phi_i + b_i$  are new XY variables. This gauge transformation maps the MXYGG model onto a ferromagnetic XY model which has a “ $\theta$  magnetization”  $\mathbf{M} \equiv \langle \sum_i (\cos \theta_i \hat{x} + \sin \theta_i \hat{y}) \rangle$ , or

$$\begin{aligned} M_x &= \left\langle \sum_i (\cos \phi_i \cos b_i - \sin \phi_i \sin b_i) \right\rangle, \\ M_y &= \left\langle \sum_i (\sin \phi_i \cos b_i + \cos \phi_i \sin b_i) \right\rangle. \end{aligned} \quad (3)$$

## III. METHODS

PCA is a dimensional reduction technique that identifies which linear combinations of the input data best characterize the data set. If the studied system has  $N$  sites, with a variable  $x_i$  (e.g.,  $\sigma_i$ ) associated with each site  $i$  ( $i = 1, \dots, N$ ), a single “state” of the system is a particular configuration of the variables  $\{x_i\}$ . The input data is defined as  $n$  such states  $\{x_i(T_t)\}$ , where each state is sampled at a temperature  $T_t$  ( $t = 1, \dots, n$ ). Note that multiple states can be sampled at a same temperature. The whole data set is formatted as a matrix

$X_{\text{data}}$ ,

$$X_{\text{data}} \equiv \begin{pmatrix} \{x_i(T_1)\} \\ \{x_i(T_2)\} \\ \vdots \\ \{x_i(T_n)\} \end{pmatrix}. \quad (4)$$

Each row is then centered by subtracting its mean value, producing a redefined matrix  $\bar{X}_{\text{data}}$ . The covariance matrix defined as  $\bar{X}_{\text{data}}^T \bar{X}_{\text{data}}$  is then diagonalized. The resulting normalized eigenvalues are the *explained variance ratios*  $\{\lambda_k\}$  with their corresponding normalized eigenvectors the *principal components*  $\{\bar{u}^{(k)}\}$ . The explained variance ratios  $\{\lambda_k\}$  (ordered from largest to smallest) quantify how correlated the data set is along the direction  $\bar{u}^{(k)}$  (most to least correlated). Furthermore, the *projection*  $\ell^{(k)}(T_i)$  of the  $i$ th state  $\{x_i(T_i)\}$  onto the  $k$ th principal component  $\bar{u}^{(k)}$  is defined as

$$\ell^{(k)}(T_i) \equiv \sum_i u_i^{(k)} x_i(T_i), \quad (5)$$

where  $u_i^{(k)}$  is the  $i$ th component of  $\bar{u}^{(k)}$ , implying that  $\bar{u}^{(k)}$  contain site-dependent information. Hence, for a fixed  $k$ ,  $\{\ell^{(k)}(T_i)\}$  is a set of  $n$  values; each state  $\{\{x_i(T_i)\}\}$  of the system is characterized by a *single* value via Eq. (5).  $\{\ell^{(k_1)}(T_i)\}$  and  $\{\ell^{(k_2)}(T_i)\}$  (for  $k_1 \neq k_2$ ) can be plotted against each other. As illustrated in Appendix C, this visually reveals how PCA “clusters” the  $n$  values from Eq. (5), along which projections the data is more or less correlated, and how important the corresponding projections are for characterizing the full data set.

PCA is applied to the MISG model by using spin configurations  $\{\sigma_i(T_i)\}$  as input data, i.e. the  $\{x_i(T_i)\}$  in Eq. (4). States are sampled for  $T_i \in [J, 4J]$ . A standard single spin flip Monte Carlo (MC) algorithm of a system of  $N = 2500$  spins ( $L = 50$ ) is used. The  $\{\varepsilon_i\}$  gauge variables at each site are randomly chosen as either  $\pm 1$  with equal probabilities. For the MXYGG model, Eq. (2), a MC simulation is performed on a system of 900 spins ( $L = 30$ ) for temperatures  $T_i \in [0.2J, 1.8J]$ . PCA is then applied to three different data sets:  $\{\cos[\phi_i(T_i)]\}$  (the “X data set”),  $\{\sin[\phi_i(T_i)]\}$  (the “Y data set”), or  $\{\{\cos[\phi_i(T_i)]\}, \{\sin[\phi_i(T_i)]\}\}$  (the “XY data set”). The gauge variables  $\{b_i\}$  are randomly drawn from a discrete distribution  $\{\frac{2\pi n}{5} \mid n = 1, \dots, 5\}$  [43].  $3 \times 10^4$  thermalization sweeps and  $5 \times 10^4$  measurement sweeps are used at every temperature for both models with 50 different temperatures in the above intervals. Sampling is done every 50 (100) measurement sweeps for the MISG (MXYGG) model, providing  $n = 5 \times 10^4$  ( $n = 2.5 \times 10^4$ ) states, respectively. For both the MISG and MXYGG cases, PCA has no prior information about the gauge variables  $\{\varepsilon_i\}$  or  $\{b_i\}$ , nor the gauged transformed Ising or XY variables  $\{\tau_i\}$  and  $\{\theta_i\}$ . The objective is thus to determine whether PCA can identify the underlying “gauge-hidden” long-range order in these systems.

## IV. RESULTS

### A. MISG model

PCA is applied to the states  $\{\sigma_i(T_i)\}$  sampled through MC simulations for the MISG model. The first principal component has a significantly larger explained variance ratio in

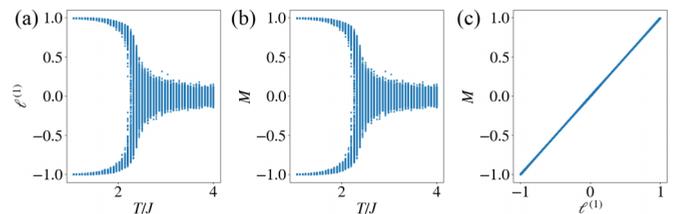


FIG. 2. (a) Projection  $\ell^{(1)}(T)$  of the MISG spin configurations  $\{\sigma_i\}$  as a function of temperature. (b)  $\tau$  magnetization  $M$  of the MISG model for the sampled spin configurations. Note that the input data contains states corresponding to both up and down  $\tau$  magnetization in order to consider all ground states allowed by the global  $\mathbb{Z}_2$  symmetry of the model. (c) Plot of  $\ell^{(1)}(T)$  against  $M$ .

comparison to all other principal components, similar to the regular Ising model as shown in Fig. 7(a) of Appendix C. Furthermore, the corresponding projection  $\ell^{(1)}(T)$  clusters the input data into a central high-temperature cluster and two adjacent low-temperature clusters, similar to Fig. 7(b) of Appendix C. This clustering pattern is similar to the pattern reported for the standard ferromagnetic Ising (FI) model, where the projection  $\ell^{(1)}(T)$  is identified as the total magnetization [24]. Indeed, the projection  $\ell^{(1)}(T)$  of the MISG model also behaves like an order parameter signaling a transition at the temperature  $T_c \approx 2.269J$  [see Fig. 2(a)]. For the MISG model, the order parameter is the  $\tau$ -magnetization  $M$ , which is calculated in the MC simulations and illustrated in Fig. 2(b). Moreover, when plotting  $\ell^{(1)}(T)$  against  $M$ , as shown in Fig. 2(c), essentially perfect agreement is found for all input states, revealing that  $\ell^{(1)}(T)$  is *indeed* the  $\tau$  magnetization.

This identification leads to the key observation: since the  $\tau$  magnetization is a function of the spin variables  $\{\sigma_i\}$  and the gauge variables  $\{\varepsilon_i\}$ , but only the former are provided to PCA, the first principal component  $\bar{u}^{(1)}$  must contain a set of gauge variables *learned* by PCA. More specifically, by directly comparing the expressions for  $\ell^{(1)}$  [Eq. (5) with  $x_i = \sigma_i$ ] and  $M = \langle \sum_i \varepsilon_i \sigma_i \rangle$ , the *learned* gauge variable  $\tilde{\varepsilon}_i$  is identified as the  $i$ th component of  $\bar{u}^{(1)}$ .

These  $\{\tilde{\varepsilon}_i\}$  gauge variables define a set of *learned* bond interactions  $\{\tilde{J}_{ij}\} \equiv \{\tilde{\varepsilon}_i \tilde{\varepsilon}_j\}$  and learned square plaquette values  $\{\tilde{P}\} \equiv \{\prod_{(i,j) \in \square} \tilde{J}_{ij}\}$ . A comparison between the distribution of the learned ( $\{\tilde{\varepsilon}_i\}$ ,  $\{\tilde{J}_{ij}\}$  and  $\{\tilde{P}\}$ ) variables identified by PCA and the original values ( $\{\varepsilon_i\}$ ,  $\{J_{ij}\}$  and  $\{P\}$ ) used in simulations is shown in Fig. 3. As can be seen in Figs. 3(a)

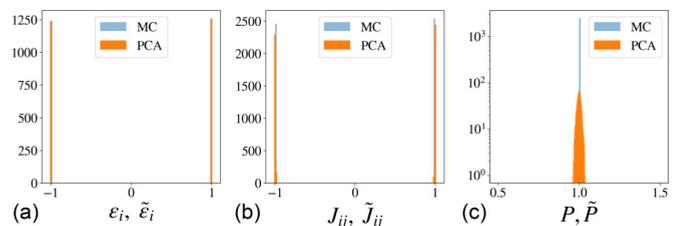


FIG. 3. Histograms of the (a) gauge variables  $\{\varepsilon_i\}$ , (b) bond interactions  $\{J_{ij}\}$ , and (c) square plaquette values  $\{P\}$  for the MISG model, comparing real values from MC (known) and values from PCA (learned).

and 3(b), the gauge variables  $\{\tilde{\varepsilon}_i\}$  and bond interactions  $\{\tilde{J}_{ij}\}$  are both described by bimodal distributions centered around  $\pm 1$ , agreeing with the distributions for the original random variables  $\{\varepsilon_i\}$  and  $\{J_{ij}\}$  introduced in the MC simulation. Moreover, a remarkable result comes from the distribution for the plaquette values  $\{\tilde{P}\}$ , shown in Fig. 3(c): this distribution is centered around the value  $P \equiv 1$ , which defines the plaquette constraint used in the construction of the MISG model. The origin of the minor spread in  $\tilde{P}$  can be traced back to an imprecision in the determination of the learned gauge variables  $\{\tilde{\varepsilon}_i\}$ . This imprecision arises from the matrix operations used in PCA. However, when comparing the learned gauge variables  $\{\tilde{\varepsilon}_i\}$  to the original gauge variables  $\{\varepsilon_i\}$  sitewise, this imprecision is only about 1% (see Appendix D). Hence, the learned gauge variables are thus a rather faithful reproduction of the original gauge variables. The faithfulness of the learned gauge variables is further confirmed by performing a MC simulation with the *learned* gauge variables  $\{\tilde{\varepsilon}_i\}$  and comparing the resulting thermodynamic behavior with that from the simulations using the *original* gauge variables  $\{\varepsilon_i\}$ , as shown in Fig. 11 of Appendix E.

Finally, as was seen when comparing  $\ell^{(1)}$  with  $M$ , sitewise multiplication of some quantity (e.g.,  $\{\sigma_i\}$ ) by  $\{\tilde{\varepsilon}_i\} \equiv \{u_i^{(1)}\}$  gauge transforms that quantity into its analog within the FI model (e.g.,  $\{\tau_i\}$ ). As shown in Appendix F, this implies that *any* principal component  $\{u_i^{(n)}\}$  of the MISG model can be mapped onto that of the FI model by this sitewise multiplication. To summarize, our results demonstrate that PCA is not only able to differentiate between the ordered and disordered configurations of the MISG model [i.e., between Figs. 1(a) and 1(b)], as a FFNN does (see Appendix A), but also produces a direct mapping of these configurations onto those of the regular FI model [i.e., Figs. 1(c) and 1(d)], thus exposing the hidden quenched gauge transformation.

## B. MXYGG model

From Eq. (5), one may wonder if PCA can only discover gauge transformations that are linearly applied to the microscopic variables. We now study the MXYGG model to demonstrate that this is not the case. For the MXYGG model, following Refs. [22,24] for the XY model, PCA is first applied to the aforementioned XY data set generated from MC simulations of Hamiltonian (2). As in the regular XY model [24], PCA finds *two* principal components with comparably high explained variance ratios. Moreover, the associated projections  $\ell^{(1)}(T)$  and  $\ell^{(2)}(T)$  corresponding to these two principal components reveal a similar clustering pattern as for the regular XY model [24] when plotted against each other, as in Fig. 4(a). For the regular XY model, these first two projections are related to the total magnetization via the quantity  $[(\ell^{(1)})^2 + (\ell^{(2)})^2]^{1/2}$ , see Fig. 12 of Ref. [24]. This relation also appears to hold for the MXYGG model when considering the same quantity, as shown in Fig. 4(b). In other words, the  $\theta$  magnetization of the MXYGG model is identified through PCA, and therefore the gauge variables must have been learned since only the “bare microscopic” angular variables  $\{\phi_i\}$  were provided to PCA. To determine the values of these learned gauge variables and motivated by the analysis of the MISG model, we turn to the expression of the  $M_x$  and

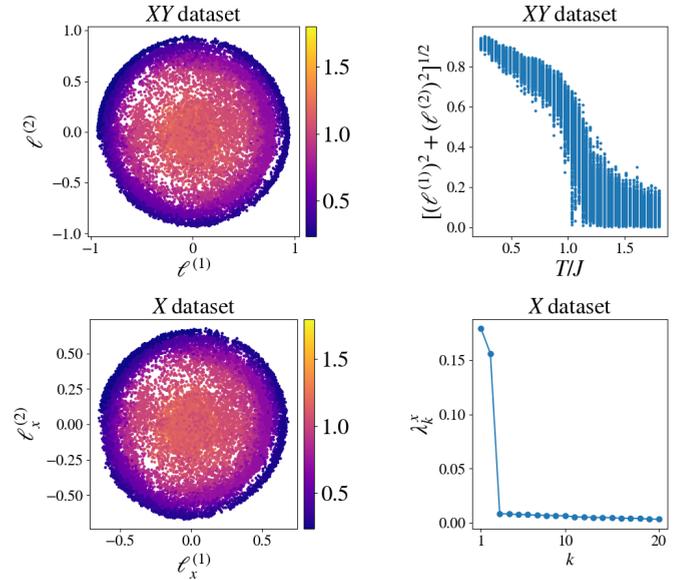


FIG. 4. Projections onto the first two principal components for (a) the XY data set and (b) the X data set. The color bar indicates the temperature at which each state is sampled. (c) Quadrature sum of  $\ell^{(1)}$  and  $\ell^{(2)}$  for the XY data set. (d) First 20 explained variance ratios for the MXYGG model, for the X data set. A similar result was obtained for  $\lambda_k^y$  for the Y data set (not shown).

$M_y$  components of the  $\theta$  magnetization in Eq. (3). These two expressions reveal how the gauge transformation decomposes the on-site spin components into two terms depending on the spin variables  $\{\phi_i\}$  and the gauge variables  $\{b_i\}$ . Motivated by this, we apply PCA onto the X and Y data sets separately.

Focusing on the X data set  $\{\cos[\phi_i(T_i)]\}$ , Fig. 4(c), two principal components are found to have comparably high explained variance ratios, as shown in Fig. 4(d). In contrast, when performing the same analysis on the regular XY model, *only one* relevant principal component is observed for the X data set, as shown in Fig. 8 of Appendix C. The projections  $\ell_x^{(1)}(T)$  and  $\ell_x^{(2)}(T)$  of the X data set of the MXYGG model are illustrated in Fig. 4(c). The reason that two principal components of similar magnitude are found for the MXYGG model’s X data set is traced back to the presence of  $\cos(\phi_i)$  in *both*  $M_x$  and  $M_y$  in Eq. (3) [44], which originates from the Mattis gauge transformation. In other words, finding that there are *two* most relevant principal components with similar explained variance ratios for the X data set indicates that PCA has identified the presence of a gauge transformation. A quantitative estimate of the learned gauge variables  $\{b_i\}$  is extracted as detailed in Appendix G. Similar to the PCA-learned  $\{\varepsilon_i\}$  of the MISG model in Fig. 3(a), PCA has learned the fivefold equally spaced distributed  $\{b_i\}$  of the MXYGG model illustrated in Fig. 16 of Appendix G. The extraction of the learned gauge variables, in addition to the finding of two most relevant principal components in the X data set in Fig. 4(c) and the  $\theta$  magnetization in Fig. 4(b), demonstrates that PCA is able to expose the gauge transformation of the MXYGG model.

V. CONCLUSION

We have applied the principal component analysis (PCA) to two spin models with random interactions, the Mattis Ising spin glass and Mattis XY gauge glass models on a square lattice, which can be respectively simplified into the ferromagnetic Ising and XY models under local gauge transformations. We have demonstrated that PCA is able to learn these gauge transformations without any prior information. Our work indicates that unsupervised machine learning (UML) protocols are indeed capable of more than just classifying data, discriminating between ordered and disordered phases, or learning other thermodynamic properties of a model. Interpretable UML methods may additionally learn hidden features of an underlying model, such as symmetries and gauge transformations. It might be interesting, as a future extension of this idea, to use more sophisticated machine learning techniques to study *quantum* models displaying a gauge symmetry. In this context, the use of PCA, neural networks [32,33], autoencoders [16], or other ML techniques as gauge-identifying protocols may become an invaluable tool. Incidentally, this line of reasoning could help elucidate how PCA and a neural network are able, when applied together, to learn the SU(2) gauge theory order parameter [21]. Taken more broadly, our results suggest that since gauge-symmetric models represent a class of models with underlying mathematical simplifications, applying UML methods to such models may help provide a deeper understanding of how these methods work and what exactly they learn.

*Note added.* In the process of finalizing this manuscript for submission, we became aware of a very recent study [45] reporting complementary results. In this study, the authors examined the MISG model using a *supervised* machine learning method. They approximately identified the gauge transformation with the weights within an intermediate layer of the neural network.

ACKNOWLEDGMENTS

We thank C. Cerkaskas, K. Chung, A. Golubeva, W. Jin, R. Melko, and S. Wetzal for useful discussions. This work was supported by the Canada Research Chair program (M.J.P.G., Tier 1) and by the NSERC of Canada CGS-M program (D.P. and M.J.P.G.).

APPENDIX A: FEEDFORWARD NEURAL NETWORK ON THE MISG MODEL

A simple 100-neuron, two-layer feedforward neural network (FFNN) was implemented with the TensorFlow Python module [46]. We use the cross entropy in addition to an L2 correction as the cost function, and an Adam optimizer to train the network. The neural network was trained for 200 epochs on configurations of the MISG model above and below the ordering transition. The prediction (output) and accuracy of the network are shown in Fig. 5. The prediction is the probability of the FFNN labelling a configuration as “ordered”. The accuracy is the percentage of configurations correctly predicted at a given temperature. Even though the FFNN correctly classifies the ordered and disordered configurations, the complexity of the network does not allow us to expose

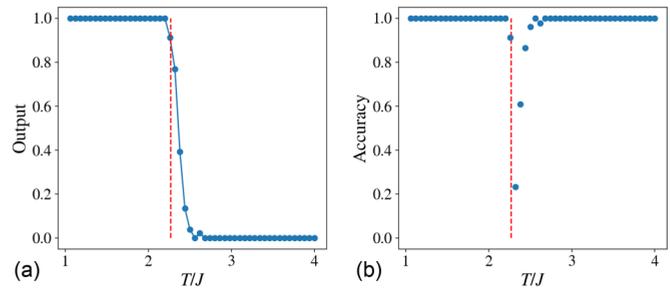


FIG. 5. (a) Output layer of the FFNN applied, to the MISG model. (b) Accuracy of the output layer of the FFNN. In both figures, the red dashed line corresponds to the analytical value of the transition temperature of the ferromagnetic Ising model.

how the network has learned to classify the configurations, or what quantities it is internally using to do so. In particular, it is not clear if the network has learned the Mattis gauge transformation.

APPENDIX B: DEFINITION OF PLAQUETTES FOR THE MISG MODEL

A plaquette in a lattice is defined as the smallest region contained within a closed loop of neighboring sites. On the square lattice, the resulting plaquettes are composed of four sites. For the Mattis transformation, we introduce gauge variables  $\epsilon_i$  for every site to define the coupling constant  $J_{ij} = \epsilon_i \epsilon_j J$  on every nearest neighbor bond. This procedure is sketched in Fig. 6 below.

APPENDIX C: PCA CLUSTERS FOR THE REGULAR ISING AND XY MODELS

For completeness and comparison, and following Refs. [1,22–24], PCA was applied to the regular (disorder-free, or pure) Ising and XY models on a square lattice. The MC simulation parameters are the same as those used for the MISG and MXYGG models as detailed in the main text, unless otherwise indicated. For the regular Ising model, PCA identifies a single principal component with a high explained variance as illustrated in Fig. 7(a). The corresponding projection  $\ell^{(1)}$  shown in Fig. 7(b) classifies the states  $\{\sigma_i\}$  into a central high temperature cluster and

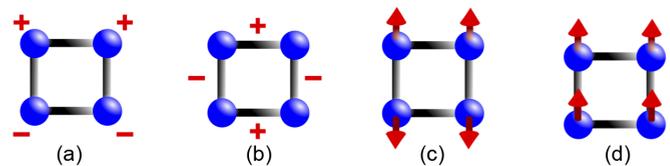


FIG. 6. Plaquette in the MISG model. (a) Example of gauge variables  $\{\epsilon_i\}$  for the four sites. (b) Resulting signs of the bond interactions  $\{J_{ij}\}$  for the four bonds. (c) Example of ground state spin configuration of  $\sigma_i$  variables for these random bond interactions. Note that the coupling in the Hamiltonian is  $-J_{ij}$  between a pair of sites  $i$  and  $j$ . (d) Ground state configuration of  $\tau_i = \epsilon_i \sigma_i$  resulting from the spin configuration illustrated in (c) with the gauge variables in (a).

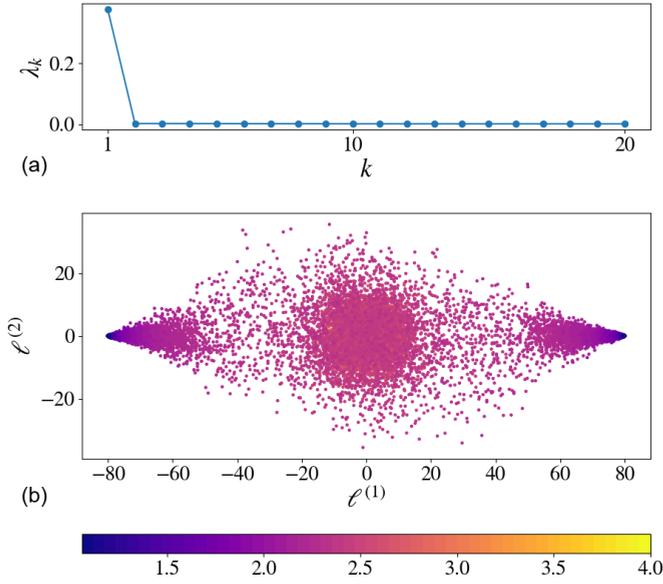


FIG. 7. (a) First 20 explained variance ratios for the regular Ising model on a  $L = 80$  square lattice and (b) principal component projection  $\ell^{(1)}$  versus  $\ell^{(2)}$ .

two adjacent low temperature clusters according to the magnetization of the state [1].

For the regular  $XY$  model, PCA is applied to the  $X$  data set ( $\{\cos[\phi_i(T_i)]\}$ ) and the  $Y$  data set ( $\{\sin[\phi_i(T_i)]\}$ ) separately. The projections onto the first two principal components for each of the two data sets are shown in Figs. 8(a) and 8(b), respectively. These projections should be compared to Fig. 4(b) of the main text, illustrating a clear difference between the clustering pattern observed here for the regular  $XY$  and the MXYGG model in the main text. This difference is indicative of the decomposition produced by the gauge transformation to the components of the magnetization Eq. (3), as detailed in the main text.

Additionally, each of the two projections obtained for the regular  $XY$  model, Figs. 8(a) and 8(b), are related to the two components of the magnetization. To support this, consider the projections  $\ell_x^{(1)}$  and  $\ell_y^{(1)}$  obtained from the  $X$  and  $Y$  data sets respectively. Consider also  $\ell^{(1)}$  and  $\ell^{(2)}$ , the

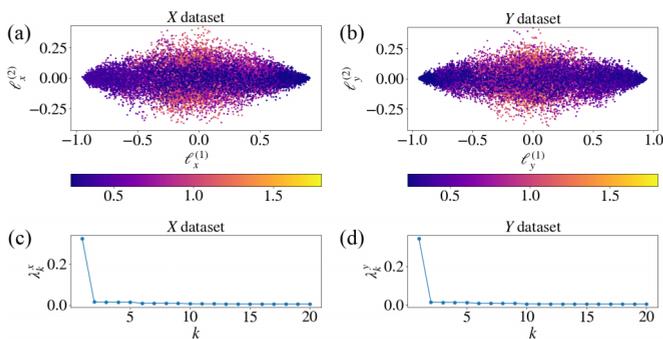


FIG. 8. (a) and (b) Principal component projections  $\ell_\alpha^{(1)}$  versus  $\ell_\alpha^{(2)}$  for the  $X$  data set ( $\alpha = x$ ) and  $Y$  data set ( $\alpha = y$ ), respectively. (c) and (d) First 20 explained variance ratios  $\lambda_k^\alpha$  for the regular  $XY$  model, for the  $X$  data set ( $\alpha = x$ ) and  $Y$  data set ( $\alpha = y$ ), respectively. The configurations were sampled from an  $L = 30$  square lattice.

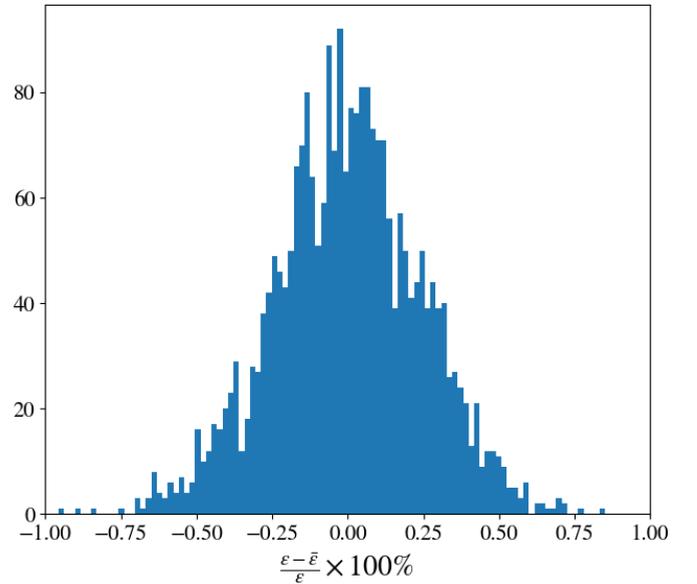


FIG. 9. Histogram of the percentage difference between the original and predicted gauge variables,  $\Delta \varepsilon_i = \varepsilon_i - \tilde{\varepsilon}_i$ , for the MISG model on a  $50 \times 50$  lattice.

first two principal components obtained for the full data set. The total magnetization is given by  $[(\ell^{(1)})^2 + (\ell^{(2)})^2]^{1/2}$ , as shown in Fig. 12 of Ref. [24]. This can be compared with  $[(\ell_x^{(1)})^2 + (\ell_y^{(1)})^2]^{1/2}$  (not shown). Good agreement between both expressions as a function of temperature is found. This supports the identification of the quantity  $[(\ell_x^{(1)})^2 + (\ell_y^{(1)})^2]^{1/2}$  with the total magnetization, which allows the projections  $\ell_x^{(1)}$  and  $\ell_y^{(1)}$  to be identified as the components of the total magnetization vector. The total magnetization  $[(\ell^{(1)})^2 + (\ell^{(2)})^2]^{1/2}$  is also studied for the full data set of the MXYGG model in the main text and illustrated in Fig. 4(b).

In the regular ferromagnetic Ising model, the points in Fig. 7(b) corresponding to high temperatures and the paramagnetic state are clustered near the middle ( $\ell^{(1)} \sim 0$ ). This is similarly seen for the regular  $XY$  model, in the center of Figs. 8(a) and 8(b) ( $\ell_x^{(1)} \sim 0$  and  $\ell_y^{(1)} \sim 0$ ). In contrast, the points in Fig. 7(b) corresponding to low temperatures form two clusters on the left and right sides. For the pure  $XY$  model, the points corresponding to low-temperature states instead form an oval-shaped cluster that cuts through the middle of Figs. 8(a) and 8(b). This difference between the low temperature projections in the Ising and  $XY$  models is easily understood, as the spin variables in the regular  $XY$  model are *continuous*. Therefore, the low-temperature projection forms one continuous horizontal overall cluster rather than two separate clusters.

#### APPENDIX D: NUMERICAL IMPRECISION OF THE LEARNED GAUGE VARIABLES FOR THE MISG MODEL

By comparing the learned and known values of the gauge variables in a sitewise manner, i.e.,  $\Delta \varepsilon_i = \varepsilon_i - \tilde{\varepsilon}_i$ , a histogram of the imprecision of PCA's results is obtained (see Fig. 9). This comparison reveals an imprecision in the learned values that is not greater than  $\sim 1\%$ ; in other words,  $\Delta \varepsilon_i$  is, site-by-

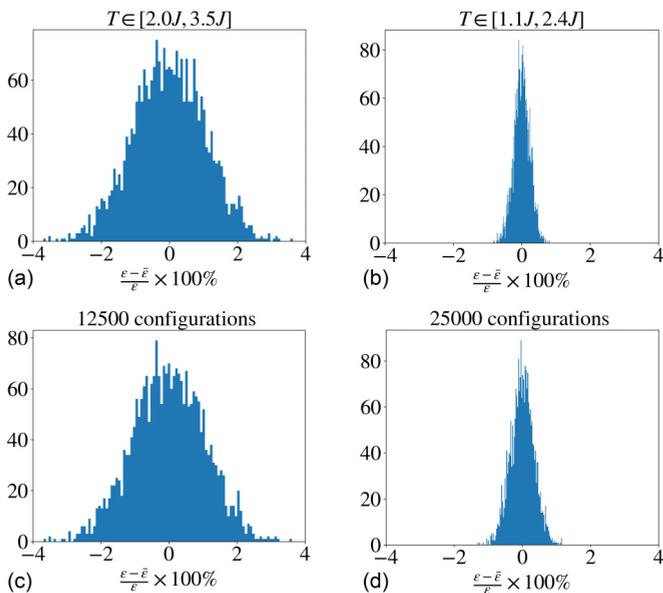


FIG. 10. Histogram of the percentage difference between the original and predicted gauge variables,  $\Delta\varepsilon_i = \varepsilon_i - \tilde{\varepsilon}_i$ , for the MISG model on a  $50 \times 50$  lattice, for (a) 24 000 configurations sampled from  $T \in [2.0J, 3.5J]$ , (b) 24 000 configurations sampled from  $T \in [1.1J, 2.4J]$ , (c) 12 500 configurations sampled from  $T \in [1.1J, 4.0J]$ , and (d) 25 000 configurations sampled from  $T \in [1.1J, 4.0J]$ .

site, of the order 1%. This minor numerical imprecision in the predicted gauge variables produces the small spread observed in the distribution of the learned plaquettes  $\{\tilde{P}\}$  in Fig. 3(c).

There are a number of ways one could proceed to reduce this imprecision. One strategy would be to apply PCA to a data set containing a greater fraction of ordered spin configurations. This can be seen by comparing Fig. 10(a), which samples configurations from mostly above  $T_c$ , and Fig. 10(b), which samples configurations from mostly below  $T_c$ , recalling that  $T_c \approx 2.269J$  for the FI model. Clearly the data set containing more ordered spin configurations results in a greater precision in the learned gauge variables. However, if the configurations provided are sampled from *just above* the critical temperature, where a finite-size effect parameter begins to acquire a nonzero value, PCA is still capable of learning the gauge variables, albeit with significantly greater imprecision. This can be thought of another way: in a temperature window in the well-ordered phase, thermal fluctuations are less prominent than near the critical temperature or in the paramagnetic phase. The reduced thermal fluctuations then result in less noisy data sets for PCA to manipulate through matrix operations and learn from, resulting in more precise learning. A second strategy for improving this precision pertains to the number of samples used in learning. From a statistics standpoint, it is evident that a larger number of samples in general should improve the precision. This is clearly seen by comparing Fig. 10(c) (with 12 500 configurations) and Fig. 10(d) (with 25 000 configurations), where the latter is more precise.

To summarize: the best way to reduce numerical imprecision in the learned gauge variables would thus be to apply

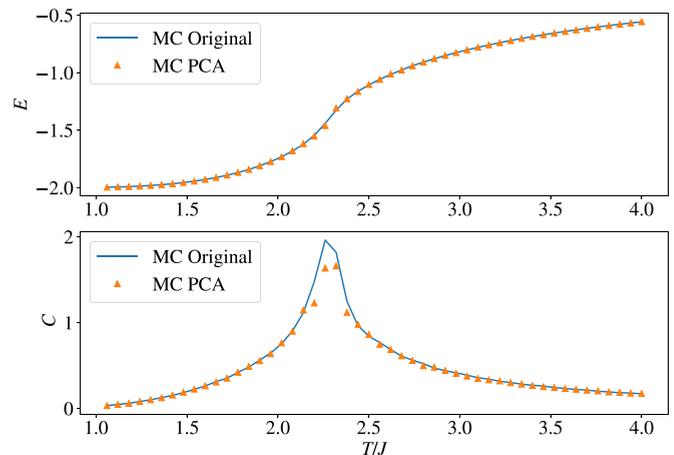


FIG. 11. Comparison of thermodynamic quantities (energy per spin  $E$  and specific heat  $C$ ) calculated within MC simulations, using the original known gauge variables and the learned gauge variables from PCA.

PCA to a large sample taken from the ordered phase. This may also prove to be a useful approach for the application of PCA to other models and further generalizations of our work; the learning precision could be improved by using data sets from the ordered phase of a given model. This can be done by taking data sets from the ordered phase from the start (i.e., if the ordered phase is already known), or by taking data sets from the ordered phase after the fact (i.e., once the “location”—that is, the critical temperature—of the ordered phase has been learned by PCA).

#### APPENDIX E: MC SIMULATION WITH THE LEARNED GAUGE VARIABLES FOR THE MISG MODEL

After applying PCA to the MISG model, we study the faithfulness of the learned gauge variables in their ability to produce thermodynamic behavior consistent with a pure ferromagnetic Ising (FI) model. We perform a Monte Carlo (MC) simulation on the MISG model using the learned gauge variables  $\{\tilde{\varepsilon}_i\}$  and compute the internal energy  $E$  and specific heat  $C$ . In Fig. 11, we compare these thermodynamic quantities with those measured for the MC simulation with the *original* set of gauge variables  $\{\varepsilon_i\}$ . As can be seen, both MC simulations possess a closely similar behavior for  $E$  and  $C$  over the studied range of temperature. Note that the small discrepancy in the specific heat  $C$  near the critical temperature is likely caused by the numerical uncertainty found in the learned  $\{\tilde{\varepsilon}_i\}$  variables that is discussed in the main text and in Appendix D.

#### APPENDIX F: GAUGE TRANSFORMATION OF HIGHER PRINCIPAL COMPONENTS FOR THE MISG MODEL

As stated in the main text, the gauge transformation identified by PCA can be used to transform higher principal components into the principal components expected for the regular FI model. This observation originates from the similarity between the projections of the MISG model and those obtained for the FI model. In other words, the gauge

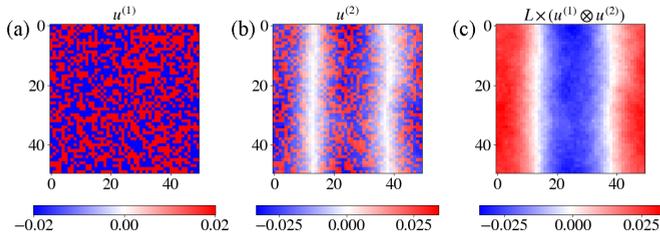


FIG. 12. Values of the (a) first and (b) second principal components of the MISG model, plotted on the lattice sites that they are associated with. (c) Product of the values of the first and second principal components on each site. Note that this product has been rescaled by the lattice dimension  $L$ .

transformation must also be present in the higher principal components. Recall that  $u_i^{(1)}\sigma_i \approx \tau_i$  maps the  $\sigma_i$  variables of the MISG model onto the  $\tau_i$  variables of the FI model. Therefore, the projection onto any principal component can be written in the following way, starting from Eq. (5) in the main text:

$$\ell^{(k)} = \sum_i u_i^{(k)} \sigma_i \quad (\text{F1})$$

$$\begin{aligned} &\approx \sum_i u_i^{(k)} (u_i^{(1)})^2 \sigma_i \\ &= \sum_i (u_i^{(k)} u_i^{(1)}) (u_i^{(1)} \sigma_i) \\ &\approx \sum_i (u_i^{(k)} u_i^{(1)}) \tau_i \\ &\equiv \sum_i \tilde{u}_i^{(k)} \tau_i, \end{aligned} \quad (\text{F2})$$

where, in the second line, we used  $u_i^{(1)} = \tilde{\varepsilon}_i \approx \varepsilon_i = \pm 1$ , so  $(u_i^{(1)})^2 \approx 1$ . This approximation comes from the finding in the main text that the learned gauge variable  $\tilde{\varepsilon}_i$  (which is  $u_i^{(1)}$ ) is almost identical to the original gauge variable  $\varepsilon_i$ . The final result in Eq. (F2) has been introduced in order to draw an analogy with Eq. (F1), but now considering the  $\tau$  magnetization of the FI model. This expression further suggests that the  $k$ th principal component of the FI model can be obtained by considering the sitewise multiplication with the first principal component [i.e.,  $u_i^{(1)} u_i^{(k)}$  in the line above Eq. (F2) should yield the  $k$ th principal component of the FI model]. To explore this further, consider for example the second principal component  $\tilde{u}^{(2)}$  of the MISG model [Fig. 12(b)] and its transformation  $\tilde{u}^{(1)} \otimes \tilde{u}^{(2)}$  [Fig. 12(c)], defined as the sitewise product of the two principal components (i.e.,  $[\tilde{u}^{(1)} \otimes \tilde{u}^{(2)}]_i = u_i^{(1)} u_i^{(2)}$ ). The transformed second principal component illustrated in Fig. 12(c) is highly similar to the second principal component obtained for the FI model reported in earlier work [24] (see Fig. 4(b) in Ref. [24]).

#### APPENDIX G: GLOBAL ROTATION OF THE REGULAR XY MODEL MAGNETIZATION

For the regular XY model, the  $x$  and  $y$  components of the magnetization vector for a state  $\{\cos[\phi_i], \sin[\phi_i]\}$  is given by

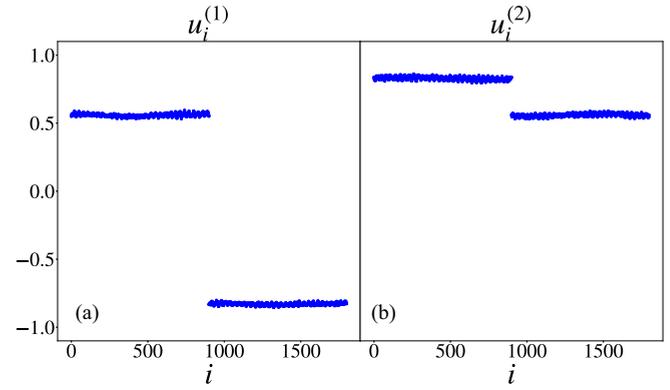


FIG. 13. Components of the (a) first and (b) second principal component eigenvectors  $\tilde{u}^{(1)}$  and  $\tilde{u}^{(2)}$  for the XY data set (defined in the main text) of the regular XY model, before applying a global rotation by an angle  $\alpha$ .

the expressions

$$M_x = \sum_i \cos(\phi_i), \quad M_y = \sum_i \sin(\phi_i). \quad (\text{G1})$$

Here we show that the projections onto the first two principal components  $\ell^{(1)}$  and  $\ell^{(2)}$  are identified as the two components of the magnetization vector in Eq. (G1) up to a global rotation of the angles  $\{\phi_i\}$ . The origin of this rotation can be understood using the global  $U(1)$  symmetry of the XY model. The XY model is invariant under global rotations of the  $x$  and  $y$  axes. Therefore, PCA can learn the magnetization components along a new set of axes,  $\bar{x}$  and  $\bar{y}$ , that are related to the original  $x$  and  $y$  axes by a global rotation by an arbitrary angle  $\alpha$ . Under such a global rotation  $\alpha$ , the magnetization vector becomes

$$\begin{aligned} M_{\bar{x}} &= \sum_i \cos(\phi_i + \alpha) \\ &= \sum_i (\cos \phi_i \cos \alpha - \sin \phi_i \sin \alpha), \\ M_{\bar{y}} &= \sum_i \sin(\phi_i + \alpha) \\ &= \sum_i (\sin \phi_i \cos \alpha + \cos \phi_i \sin \alpha). \end{aligned} \quad (\text{G2})$$

Comparing the expressions for  $M_{\bar{x}}$  and  $M_{\bar{y}}$  with Eq. (5) of the main text implies that, when the XY data set is provided to

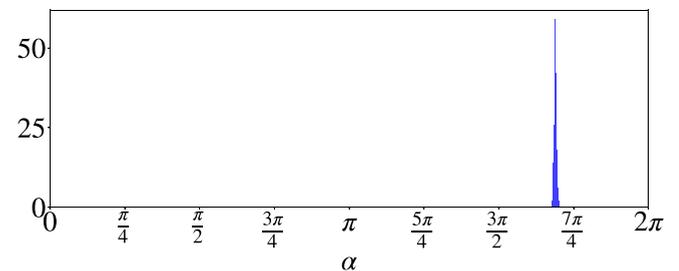


FIG. 14. Histogram of extracted global rotation angle  $\alpha$  for the regular XY model.

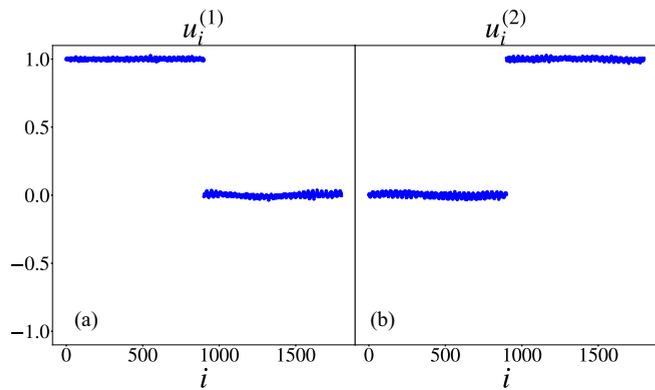


FIG. 15. Components of the (a) first and (b) second principal component eigenvectors  $\bar{u}^{(1)}$  and  $\bar{u}^{(2)}$  for the XY data set (defined in the main text) of the regular XY model, *after* having applied a global rotation by an angle  $\alpha$ .

PCA (namely,  $\{\{\cos[\phi_i]\}, \{\sin[\phi_i]\}\}$ ), the coefficients  $u_i^{(k)}$  in Eq. (5) are simply  $\cos(\alpha)$  and  $\sin(\alpha)$  (up to a minus sign). The principal components therefore give a direct measurement of the global rotation angle  $\alpha$  along which the magnetization vector is learned by PCA.

In fact, since the  $u_i^{(k)}$  coefficients correspond to  $\sin(\alpha)$  and  $\cos(\alpha)$ , the value of  $\alpha$  is determined from the  $u_i^{(k)}$  values (e.g., by taking the arctangent of the ratio of pairs of  $u_i^{(k)}$  values, corresponding to the two components of the magnetization for the same site  $i$ ). Extracting  $\alpha$  in this way (using the principal components  $u_i^{(k)}$  shown in Fig. 13) produces the histogram shown in Fig. 14. Furthermore, if the angles  $\{\phi_i\}$  are rotated by  $\alpha$  *prior* to applying PCA (i.e., giving PCA  $\{\{\cos[\phi_i + \alpha]\}, \{\sin[\phi_i + \alpha]\}\}$ ), the principal components only take values of only 1s or 0s, as shown in Fig. 15. Having only values of

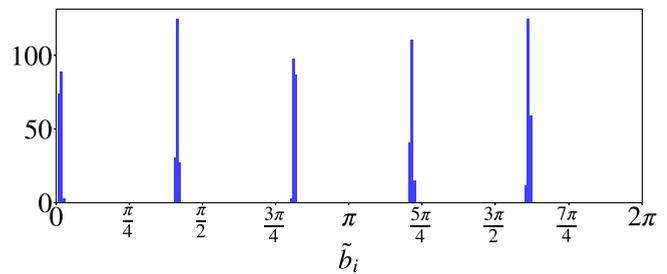


FIG. 16. Histogram of learned gauge variables  $\{\tilde{b}_i\}$  for the MXYGG model, revealing five equally spaced peaks corresponding to the five discrete choices of  $b_i$  in the MC simulation.

1s and 0s can be easily interpreted by comparing Eq. (5) of the main text and Eq. (G2). It means that PCA is directly summing  $\cos(\phi_i + \alpha)$  across all sites to calculate  $M_{\bar{x}}$  (and similarly for  $M_{\bar{y}}$ ), and is therefore learning the magnetization vector along these new axes.

This same analysis can be applied to the principal components of the MXYGG model to extract the *local* rotations produced by the gauge variables  $\{b_i\}$ . The resulting distribution for the *learned* gauge variables  $\{\tilde{b}_i\}$  is shown in Fig. 16, revealing five equally spaced peaks as expected for the five equally spaced choices of the original gauge variables defining the chosen (realized) MXYGG model. Note that a global rotation  $\alpha$  is also possible for the MXYGG model in addition to the local rotations  $\{b_i\}$  on each site, but the two rotations cannot be separately identified with the approach discussed here. Due to this possibility of a global rotation by an angle  $\alpha$ , there is an overall shift in the peaks of the histogram in Fig. 16 relative to the expected values  $\{\frac{2\pi n}{5} \mid n = 1, \dots, 5\}$  that are used in the MC simulations of the MXYGG model.

- [1] L. Wang, Discovering phase transitions with unsupervised learning, *Phys. Rev. B* **94**, 195105 (2016).
- [2] J. Carrasquilla and R. G. Melko, Machine learning phases of matter, *Nat. Phys.* **13**, 431 (2017).
- [3] P. Mehta, M. Bukov, C.-H. Wang, A. G. R. Day, C. Richardson, C. K. Fisher, and D. J. Schwab, A high-bias, low-variance introduction to machine learning for physicists, *Phys. Rep.* **810**, 1 (2019).
- [4] V. Dunjko and H. J. Briegel, Machine learning & artificial intelligence in the quantum domain: A review of recent progress, *Rep. Prog. Phys.* **81**, 074001 (2018).
- [5] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, Machine learning and the physical sciences, *Rev. Mod. Phys.* **91**, 045002 (2019).
- [6] K.-W. Zhao, W.-H. Kao, K.-H. Wu, and Y.-J. Kao, Generation of ice states through deep reinforcement learning, *Phys. Rev. E* **99**, 062106 (2019).
- [7] J. Greitemann, K. Liu, and L. Pollet, Probing hidden spin order with interpretable machine learning, *Phys. Rev. B* **99**, 060404(R) (2019).
- [8] M. J. S. Beach, A. Golubeva, and R. G. Melko, Machine learning vortices at the Kosterlitz-Thouless transition, *Phys. Rev. B* **97**, 045207 (2018).
- [9] J. Greitemann, K. Liu, L. D. C. Jaubert, H. Yan, N. Shannon, and L. Pollet, Identification of emergent constraints and hidden order in frustrated magnets using tensorial kernel methods of machine learning, *Phys. Rev. B* **100**, 174408 (2019).
- [10] P. Ponte and R. G. Melko, Kernel methods for interpretable machine learning of order parameters, *Phys. Rev. B* **96**, 205146 (2017).
- [11] K. Liu, J. Greitemann, and L. Pollet, Learning multiple order parameters with interpretable machines, *Phys. Rev. B* **99**, 104410 (2019).
- [12] H. Théveniaut and F. Alet, Neural network setups for a precise detection of the many-body localization transition: Finite-size scaling and limitations, *Phys. Rev. B* **100**, 224202 (2019).
- [13] A. Canabarro, S. Brito, and R. Chaves, Machine Learning Non-local Correlations, *Phys. Rev. Lett.* **122**, 200401 (2019).
- [14] X. Liang, W.-Y. Liu, P.-Z. Lin, G.-C. Guo, Y.-S. Zhang, and L. He, Solving frustrated quantum many-particle models with convolutional neural networks, *Phys. Rev. B* **98**, 104426 (2018).
- [15] M. August and X. Ni, Using recurrent neural networks to optimize dynamical decoupling for quantum memory, *Phys. Rev. A* **95**, 012335 (2017).
- [16] A. Decelle, V. Martin-Mayor, and B. Seoane, Learning a local symmetry with neural networks, *Phys. Rev. E* **100**, 050102(R) (2019).

- [17] H. Xu, J. Li, L. Liu, Y. Wang, H. Yuan, and X. Wang, Generalizable control for quantum parameter estimation through reinforcement learning, *npj Quantum Inf.* **5**, 82 (2019).
- [18] A. Bohrdt, C. S. Chiu, G. Ji, M. Xu, D. Greif, M. Greiner, E. Demler, F. Grusdt, and M. Knap, Classifying snapshots of the doped Hubbard model with machine learning, *Nat. Phys.* **15**, 921 (2019).
- [19] C. Casert, T. Viejira, J. Nys, and J. Ryckebusch, Interpretable machine learning for inferring the phase boundaries in a nonequilibrium system, *Phys. Rev. E* **99**, 023304 (2019).
- [20] A. Canabarro, F. F. Fanchini, A. L. Malvezzi, R. Pereira, and R. Chaves, Unveiling phase transitions with machine learning, *Phys. Rev. B* **100**, 045129 (2019).
- [21] S. J. Wetzel and M. Scherzer, Machine learning of explicit order parameters: From the Ising model to  $SU(2)$  lattice gauge theory, *Phys. Rev. B* **96**, 184410 (2017).
- [22] C. Wang and H. Zhai, Machine learning of frustrated classical spin models (I): Principal component analysis, *Phys. Rev. B* **96**, 144432 (2017).
- [23] C. Wang and H. Zhai, Machine learning of frustrated classical spin models (II): Kernel principal component analysis, *Front. Phys.* **13**, 130507 (2018).
- [24] W. Hu, R. R. P. Singh, and R. T. Scalettar, Discovering phases, phase transitions, and crossovers through unsupervised machine learning: A critical examination, *Phys. Rev. E* **95**, 062122 (2017).
- [25] S. J. Wetzel, Unsupervised learning of phase transitions: From principal component analysis to variational autoencoders, *Phys. Rev. E* **96**, 022140 (2017).
- [26] W. Zhang, J. Liu, and T.-C. Wei, Machine learning of phase transitions in the percolation and  $XY$  models, *Phys. Rev. E* **99**, 032142 (2019).
- [27] Y. Iwasaki, R. Sawada, V. Stanev, M. Ishida, A. Kirihaara, Y. Omori, H. Someya, I. Takeuchi, E. Saitoh, and S. Yorozu, Identification of advanced spin-driven thermoelectric materials via interpretable machine learning, *npj Comput. Mater.* **5**, 103 (2019).
- [28] Y. Wu and H. Zhai, Generalized independent component analysis for extracting eigen-modes of a quantum system, *Machine Learn. Sci. Technology* **1**, 025010 (2020).
- [29] T. Hou, K. Y. M. Wong, and H. Huang, Minimal model of permutation symmetry in unsupervised learning, *J. Phys. A: Math. Theor.* **52**, 414001 (2019).
- [30] R. B. Jadrich, B. A. Lindquist, and T. M. Truskett, Unsupervised machine learning for detection of phase transitions in off-lattice systems. I. Foundations, *J. Chem. Phys.* **149**, 194109 (2018).
- [31] S. N. Shah, Variational approach to unsupervised learning, *J. Phys. Commun.* **3**, 075006 (2019).
- [32] D. Luo, G. Carleo, B. K. Clark, and J. Stokes, Gauge Equivariant Neural Networks for Quantum Lattice Gauge Theories, *Phys. Rev. Lett.* **127**, 276402 (2021).
- [33] D. L. Boyda, M. N. Chernodub, N. V. Gerasimeniuk, V. A. Goy, S. D. Liubimov, and A. V. Molochkov, Finding the deconfinement temperature in lattice Yang-Mills theories from outside the scaling window with machine learning, *Phys. Rev. D* **103**, 014509 (2021).
- [34] D. Mattis, Solvable spin systems with random interactions, *Phys. Lett. A* **56**, 421 (1976).
- [35] K. H. Fischer and J. A. Hertz, *Spin Glasses* (Cambridge University Press, Cambridge, 1991).
- [36] Previous work [16] used an autoencoder to identify whether two instances of a local  $\mathbb{Z}_2$  Mattis gauge transformation (though the authors did not use this terminology) were mutually related under the application of Wilson and/or Polyakov loops, culminating in a binary classification (i.e., equivalent or not equivalent). By contrast, our work applies a ML protocol directly to the thermodynamic degrees of freedom, the spin configuration, and obtains a quantitative estimate of the underlying  $\mathbb{Z}_2$  Mattis transformation without prior knowledge of its existence.
- [37] D. Kim and D.-H. Kim, Smallest neural network to learn the Ising criticality, *Phys. Rev. E* **98**, 022138 (2018).
- [38] Y. Zhang, P. Ginsparg, and E.-A. Kim, Interpreting machine learning of topological quantum phase transitions, *Phys. Rev. Res.* **2**, 023283 (2020).
- [39] K. Pearson, On lines and planes of closest fit to systems of points in space, *London, Edinburgh, Dublin Philos. Mag. J. Sci.* **2**, 559 (1901).
- [40] M. J. P. Gingras, Real-space renormalization-group study of random-superconductor models, *Phys. Rev. B* **43**, 13747 (1991).
- [41] M. J. P. Gingras, Universality class of  $XY$ -like spin glasses lacking time-reversal symmetry, *Phys. Rev. B* **44**, 7139 (1991).
- [42] M. J. P. Gingras, Numerical study of vortex-glass order in random-superconductor and related spin-glass models, *Phys. Rev. B* **45**, 7547 (1992).
- [43] It is easier to determine PCA's success with determining the gauge variables if they are taken from a discrete distribution, which can be clearly identified in a histogram such as Fig. 16 of Appendix G, as opposed to a continuous distribution.
- [44] The contribution of  $\cos(\phi_i)$  to both components of the magnetization in Eq. (3) explains the similarity between Fig. 4(a) of the full  $XY$  data set and Fig. 4(b) of the  $X$  data set.
- [45] T. Morishita and S. Todo, Randomized-gauge test for machine learning of Ising model order parameter, *J. Phys. Soc. Jpn.* **91**, 044001 (2022).
- [46] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg *et al.*, TensorFlow: Large-scale machine learning on heterogeneous systems (2015), software available from [tensorflow.org](https://www.tensorflow.org).