

Tutorial 3: Identifying phase transitions using principal component analysis

May 29, 2019

The objective of this tutorial is to use the dimensional reduction technique known as principal component analysis (PCA) to identify phases without explicitly training with phase labels. You will reproduce the results in Figures 1 and 2 of Reference [1].

The goal of dimensional reduction is to generate a lower-dimensional representation $\mathcal{D}' = \{\mathbf{x}'\}$ of a high-dimensional dataset $\mathcal{D} = \{\mathbf{x}\}$, where $\mathbf{x}' \in \mathbb{R}^{N'}$, $\mathbf{x} \in \mathbb{R}^N$ and $N' < N$. The lower-dimensional dataset should still encode the important features of the original higher-dimensional data. The PCA method attempts to accomplish this goal by applying a linear transformation. In this tutorial we will apply PCA to N -dimensional spin configurations of the two-dimensional Ising model.

Our data is stored in an $M \times N$ matrix X , where each of the M rows stores a spin configuration for a system with N spins. For the two-dimensional Ising model we have $N = L^2$.

PCA can be performed on a matrix X^c where each column has mean 0. One can calculate X^c from X as

$$X_{ij}^c = X_{ij} - \frac{1}{M} \sum_{k=1}^M X_{kj}. \quad (1)$$

The principal components x'_1, x'_2, \dots are then stored in the columns of an $M \times N$ matrix

$$X' = X_c P, \quad (2)$$

where P is an $N \times N$ matrix. P is determined by solving the eigenvalue problem

$$\frac{1}{M-1} X_c^T X_c = P^T D P, \quad (3)$$

where D is a diagonal matrix with non-negative entries $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N \geq 0$.

Another important definition is the so-called *explained variance ratio* r_ℓ , which measures how much of the variance in the dataset X can be explained by the principal component x'_ℓ . This ratio is defined in terms of the eigenvalues λ_ℓ as

$$r_\ell = \frac{\lambda_\ell}{\sum_{i=1}^N \lambda_i}. \quad (4)$$

For this tutorial, you have been given a dataset containing rows of spin configurations $[s_1, s_2, \dots, s_N]$ for the two-dimensional Ising model on various sized lattices. Each spin ‘up’ is stored as 1 and each spin ‘down’ is stored as -1.

You have been given data for $L = 20, 40$ and 80 . Each spin configuration file contains 100 spin configurations at each of the 20 temperatures $T/J = 1.0, 1.1, 1.2, \dots, 2.9$ such that $M = 2000$ for each lattice size. For each L , there is a corresponding file storing the temperature at which each configuration was generated (using Monte Carlo simulation). The temperature data will not be used to determine the principal components and will only be used for data visualization purposes. If you wish, you could generate this data yourself using the code from Tutorial 1.

- a) Write code that reads in the spin configurations for the Ising model for a given lattice size and determines the principal components x'_1, x'_2, \dots . Make a scatter plot of x'_1 versus x'_2 for each of the lattice sizes. What do you notice about the behaviour of the resulting two-dimensional cluster(s) as L increases?

Hint: You may find it useful to use the function

```
(lamb, P) = np.linalg.eig(np.dot(Xc.T, Xc))
```

When `np.dot(Xc.T, Xc)` is an $N \times N$ matrix, this function will return the N eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_N$ in the array `lamb`. The eigenvector corresponding to `lamb[i]` will be returned in `P[:,i]`.

- b) Label the points in your plot such that they are coloured according to their temperature and compare with Figure 2 of Reference [1]. What does each cluster correspond to in terms of the phases of the two-dimensional Ising model?
- c) Consider now the explained variance ratios r_ℓ . Plot the largest 10 explained variance ratios for each lattice size and compare with Figure 1 of Reference [1]. How many principal components are needed to explain how the Ising spin configurations vary as a function of temperature?
- d) Let p_ℓ be the i^{th} column of the matrix P such that $x'_\ell = Xp_\ell$. Plot the elements of p_1 . What does your plot tell you about how x'_1 is computed from the data X ? Relate your plot to the magnetization order parameter for the Ising model, which is given by $\frac{1}{N} \sum_i s_i$.

References

- [1] L. Wang, Phys. Rev. B **94**, 195105 (2016), <https://arxiv.org/abs/1606.00318>.