



HAL
open science

Evidence-based data mining method to reveal similarities between materials based on physical mechanisms

Minh-Quyet Ha, Duong-Nguyen Nguyen, Viet-Cuong Nguyen, Hiori Kino, Yasunobu Ando, Takashi Miyake, Thierry Dencœux, Van-Nam Huynh, Hieu-Chi Dam

► **To cite this version:**

Minh-Quyet Ha, Duong-Nguyen Nguyen, Viet-Cuong Nguyen, Hiori Kino, Yasunobu Ando, et al.. Evidence-based data mining method to reveal similarities between materials based on physical mechanisms. *Journal of Applied Physics*, 2023, 133 (5), pp.053904. 10.1063/5.0134999 . hal-03976318

HAL Id: hal-03976318

<https://hal.science/hal-03976318>

Submitted on 7 Feb 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Evidence-based data mining method to reveal similarities between 2 materials based on physical mechanisms

3 Minh-Quyet Ha,¹ Duong-Nguyen Nguyen,¹ Viet-Cuong Nguyen,² Hiori Kino,³ Yasunobu Ando,⁴ Takashi Miyake,⁴
4 Thierry Dencœux,⁵ Van-Nam Huynh,¹ and Hieu-Chi Dam¹

5 ¹*Japan Advanced Institute of Science and Technology, 1-1 Asahidai, Nomi, Ishikawa 923-1292,*
6 *Japan*

7 ²*HPC SYSTEMS Inc., 3-9-15 Kaigan, Minato, Tokyo 108-0022, Japan*

8 ³*Research and Services Division of Materials Data and Integrated System, National Institute for Materials Science,*
9 *1-2-1 Sengen, Tsukuba, Ibaraki 305-0044, Japan*

10 ⁴*Research Center for Computational Design of Advanced Functional Materials,*
11 *National Institute of Advanced Industrial Science and Technology, 1-1-1 Umezono, Tsukuba, Ibaraki 305-8568,*
12 *Japan*

13 ⁵*Université de Technologie de Compiègne, CNRS, UMR 7253 Heudiasyc, Compiègne,*
14 *France*

15 (*Electronic mail: dam@jaist.ac.jp)

16 (Dated: 16 January 2023)

17 Measuring the similarity between materials is essential for estimating their properties and revealing the asso-
18 ciated physical mechanisms. However, current methods for measuring the similarity between materials rely
19 on theoretically derived descriptors and parameters fitted from experimental or computational data, which
20 are often insufficient and biased. Further, outliers and data generated by multiple mechanisms are usually
21 included in the dataset, making the data-driven approach challenging and mathematically complicated. To
22 overcome such issues, we apply the Dempster–Shafer theory to develop an evidential regression-based similar-
23 ity measurement (eRSM) method, which can rationally transform data into evidence. It then combines such
24 evidence to conclude the similarities between materials, considering their physical properties. To evaluate the
25 eRSM, we used two materials datasets, including 3d transition metal–4f rare-earth binary and quaternary
26 high-entropy alloys with target properties, Curie temperature and magnetization. Based on the informa-
27 tion obtained on the similarities between the materials, a clustering technique is applied to learn the cluster
28 structures of the materials that facilitate the interpretation of the mechanism. The unsupervised learning
29 experiments demonstrate that the obtained similarities are applicable to detect anomalies and appropriately
30 identify groups of materials whose properties correlate differently with their compositions. Furthermore,
31 significant improvements in the accuracies of the predictions for the Curie temperature and magnetization
32 of the quaternary alloys are obtained by introducing the similarities, with the reduction in mean absolute
33 errors (MAE) of 36% and 18%, respectively. The results show that the eRSM can adequately measure the
34 similarities and dissimilarities between materials in these datasets with respect to mechanisms of the target
35 properties.

36 I. INTRODUCTION

37 The concept of machine learning has great potential
38 for application in several areas of materials science, espe-
39 cially for discovering new materials. In materials science,
40 a number of the problems addressed by data-driven ap-
41 proaches require the effective utilization of existing ma-
42 terial data for predicting the properties of new mate-
43 rials and understanding the underlying physicochemical
44 mechanisms¹.

45 From an engineering point of view, developing a data-
46 driven model that quickly and accurately predicts the
47 physical properties of possible materials from accumu-
48 lated data can reduce the time required for material de-
49 velopment. By applying a data-driven model to screen
50 materials *in-silico*, we narrow down the candidates that
51 require expensive calculations and experiments to verify.
52 If there are sufficient independent supervised data from
53 the distribution of the target material data, a model with
54 high prediction accuracy can be built using state-of-the-

55 art data-driven techniques. However, because materi-
56 als research and development aim to develop materials
57 that are superior to existing ones, the distribution of the
58 target prediction data may be completely different from
59 the distribution of the original training data. Therefore,
60 there are concerns about whether data-driven models can
61 accurately predict the physical properties of new materi-
62 als.

63 On the contrary, considering the history of materi-
64 als science, researchers have discovered various materi-
65 als through a loop of hypothesis and verification based
66 on their knowledge, experience, and serendipity. Partic-
67 ularly, hypothesizing relies heavily on describing, inter-
68 preting, and understanding the underlying physicochem-
69 ical mechanisms of the observed physical phenomena of
70 materials. Scientifically, applying a data-driven approach
71 to extracting knowledge from existing complicated mate-
72 rial data can accelerate the process of describing, inter-
73 preting, and understanding the physicochemical mech-
74 anisms underlying the observed physical phenomena of

1 materials. This reduces the time required for material
2 development. Hence, to be effectively applied to materi-
3 als science, data-driven approaches that are interpretable
4 and understandable to humans must be developed.

5 One of the most intuitive and interpretable data-driven
6 approaches for humans is analogy-based inductive reason-
7 ing, which infers the properties of a new instance using
8 the information of the observed instances that are
9 most similar to it²⁻⁵. By applying analogy-based mod-
10 els, we can easily explain the reasoning process behind
11 the predictions and reveal the physicochemical mecha-
12 nisms rationalizing the observations^{6,7}. Materials scien-
13 tists have resolved different problems in materials science
14 by systematizing information about analogies in compo-
15 sition or structure between materials that exhibit similar
16 physicochemical properties⁸⁻¹¹.

17 Especially, in a discipline based on fundamental prin-
18 ciples, such as condensed matter physics, it is essential
19 to elucidate the physical mechanisms and which materi-
20 als are manifested through each of these physical mecha-
21 nisms. However, despite several new materials and super-
22 rior properties having been discovered, it is still difficult
23 to appropriately quantify the similarities between mate-
24 rials to elucidate the underlying physicochemical mech-
25 anisms of these properties. Furthermore, this difficulty
26 arises from the fact that the mechanisms of materials'
27 properties are typically interpreted in terms of physico-
28 chemical concepts based on relative criteria.

29 The phenomenon of superconductivity in materials,
30 which originates from the instability of metals, is a well-
31 known example of the above difficulty. One of the most
32 successful theories that describe the microscopic mech-
33 anisms is the Bardeen-Cooper-Schrieffer (BCS) theory
34 for superconductivity¹², the origin of which is electron-
35 phonon interactions. However, there also exist other
36 mechanisms. For example, one of the most plausible
37 origins of superconductivity in the high- T_C cuprates is
38 electron-electron interactions. Nevertheless, it is not easy
39 to achieve a consensus of classifying the superconducting
40 mechanism of materials among researchers as the ori-
41 gins. Although the emergence of superconductivity is
42 basically due to the instability in the metallic phase, it
43 is not easy to achieve the consensus because both the
44 mentioned and other mechanisms can contribute cooper-
45 atively in increasing the T_C value, for example. Although
46 it is challenging to classify individual materials when con-
47 sidering phenomena that cause such a situation, it is ex-
48 pected that the underlying physical mechanisms can be
49 discovered if we can inductively quantify the similarities
50 between the materials of interest and group similar ma-
51 terials using all observation data.

52 Incidentally, inductive reasoning with inefficient sim-
53 ilarity assessment can lead to misidentification of
54 outliers¹³ and difficulty in explaining the underlying
55 physicochemical mechanisms of datasets using single
56 models. Therefore, regarding predefined material de-
57 scriptors, an exhaustive examination of all possible hy-
58 potheses about the unknown physicochemical mecha-

59 nisms is necessary to assess the similarity between the
60 materials. Furthermore, similarity measures are usually
61 context-dependent. Because the context changes, the
62 similarity measure must be modified to adequately cap-
63 ture the phenomena under study^{14,15}. Thus, a quanti-
64 tative measure of similarity needs to consider the uncer-
65 tainty arising from the context or the measurement itself,
66 especially in situations where material data are often in-
67 sufficient and heavily biased. Moreover, similarities from
68 different contexts may not be directly comparable in the
69 integration to draw conclusions about the similarity be-
70 tween materials. These reasons make it challenging to
71 apply data-driven approaches to materials science.

72 To overcome such issues and efficiently extract knowl-
73 edge from the data, we propose a new approach that
74 shifts from measuring the similarity between materials
75 to quantitatively measure the confidence in their simi-
76 larities. We adopt the Dempster-Shafer theory¹⁶⁻¹⁸, re-
77 ferred to as the evidence theory, to develop an eviden-
78 tial regression-based similarity measurement (eRSM) for
79 detecting subgroups of materials such that learned mod-
80 els from the subgroups show high correlations between
81 descriptors and the target property of the constituent
82 materials. Further analysis of models describing the sub-
83 groups provide valuable information to extract, interpret,
84 and understand physical mechanisms. The Dempster-
85 Shafer theory can be regarded as a generalization of the
86 Bayesian approach for solving the problem of incomplete
87 and insufficient information. Moreover, it is suitable for
88 solving material data problems^{19,20}. The measure of sim-
89 ilarity here refers to whether the observed physical prop-
90 erties of the materials under study are described using
91 the same hidden mechanism that has not yet been re-
92 vealed. In other words, we consider any pair of materials
93 (in the dataset) as similar if their physical properties can
94 be described by the same hidden mechanism; otherwise,
95 the pair of materials is considered dissimilar. We then
96 first generate numerous hypothetical mechanisms by ran-
97 domly choosing subsets of data instances and construct-
98 ing regression models for each subset. Each regression
99 model is considered a source of evidence of the similar-
100 ities between materials. Thereafter, the Dempster-Shafer
101 theory¹⁶⁻¹⁸, which has a foundation for modeling and
102 combining the uncertainty of evidence, is applied to inte-
103 grate the collected pieces of evidence to draw conclusions
104 about the similarities between materials. The eRSM con-
105 sists of three main steps as follows:

- 106 1. *Collect sources of evidence:* Hypothetical mecha-
107 nisms are collected from a dataset by applying re-
108 gression analysis with single or mixture models and
109 are used as sources of evidence to rationalize the
110 similarity states of materials.
- 111 2. *Model similarity evidence:* An appropriate mass
112 function is designed to model the obtained evidence
113 within the framework of the evidence theory.
- 114 3. *Combine pieces of evidence:* Dempster's rule of

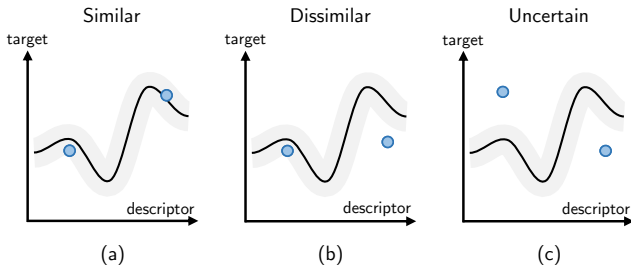


FIG. 1. Illustrative figures of the three possible similarity states between two data instances (blue circles), including similar (a), dissimilar (b), and uncertain (c), considering a referential regression model f_r (black line). The gray region is the interval that determines whether a data instance can be considered to have been generated by regression model f_r .

1 combination is used to integrate the pieces of the
2 evidence.

3 The steps of the eRSM are explained in detail in Sec-
4 tion II. Regarding the framework of the evidence theory,
5 the essential contributions of the eRSM are collecting
6 sources of evidence about the similarities between mate-
7 rials from datasets and designing suitable mass functions
8 to model the pieces of evidence rationally. The effective-
9 ness of obtained similarities using the eRSM for subdiv-
10 iding alloys from datasets into homogenous subgroups
11 is supported by experiments on 1) a dataset of binary
12 alloys with their Curie temperature as a target property
13 (Section III B); and 2) two dataset of quaternary alloys
14 with their magnetization (Section III C) and Curie tem-
15 perature (Section III D) as the target properties. Further
16 analysis of the detected subgroups to interpret the under-
17 lying physical mechanisms is shown in Section III E

18 II. METHODOLOGY

19 We consider a dataset \mathcal{D} consisting of p data in-
20 stances. We assume that a data instance with index
21 i in \mathcal{D} is described by n predefined descriptors and
22 is represented by an n -dimensional numerical vector,
23 $\mathbf{x}_i = (x_i^1, x_i^2, \dots, x_i^n) \in \mathbb{R}^n$. The target property of
24 the data instance \mathbf{x}_i is $y_i \in \mathbb{R}$. Thereafter, the dataset
25 $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_p, y_p)\}$ is represented using
26 a $(p \times (n + 1))$ matrix. In this study, we consider that
27 \mathcal{D} may contain pairs of data instances \mathbf{x}_i and \mathbf{x}_j , where
28 $\mathbf{x}_i \approx \mathbf{x}_j$; however, the value of y_i is far from y_j .

29 A. Collecting sources of similarity evidence

30 We perform random subset sampling of the data in-
31 stances without replacement to collect a large amount of
32 evidence of the similarity between pairs of data instances
33 in \mathcal{D} . Considering each sample, we obtain two datasets:
34 the reference dataset, \mathcal{D}_{ref} , and the evaluation dataset,

35 \mathcal{D}_{eval} ($\mathcal{D}_{ref} \cap \mathcal{D}_{eval} = \emptyset$ and $\mathcal{D}_{ref} \cup \mathcal{D}_{eval} = \mathcal{D}$). Con-
36 sidering \mathcal{D}_{ref} , we can generate a single or multiple ref-
37 erence functions $f_r : \mathbb{R}^n \rightarrow \mathbb{R}$ using a Gaussian process
38 (GP)²¹ or a mixture of Gaussian processes (MGP)²², re-
39 spectively. This study applies GP- or MGP-based models
40 instead of other nonlinear regression models such as ker-
41 nel ridge regression²³, random forest regression²⁴, or arti-
42 ficial neural networks²⁵ because GP or MGP can quantify
43 the uncertainty of its prediction without introducing any
44 other statistical validation. The sampling ratios of \mathcal{D}_{ref}
45 from \mathcal{D} are fixed at 0.3 and 0.7 for the experiments with
46 GP and MGP, respectively. Each reference function f_r
47 is considered as a source to provide pieces of evidence
48 for the similarity between (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) in \mathcal{D}_{eval} .
49 The function f_r is not used to provide any information
50 about the similarities between the data instances in \mathcal{D}_{ref}
51 or between a data instance in \mathcal{D}_{ref} and a data instance
52 in \mathcal{D}_{eval} . This is to exclude self-evaluation to ensure the
53 objectivity of the evidence. Regarding a reference func-
54 tion f_r , we consider the state of the similarity between
55 (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) as:

- 56 • Similar: Both data instances can be considered to
57 have been generated by the function f_r (Fig. 1 a).
- 58 • Dissimilar: Only one of the data instances can be
59 considered to have been generated by the function
60 f_r (Fig. 1 b).
- 61 • Uncertain: Neither of the data instances can be
62 considered to have been generated by the function
63 f_r (Fig. 1 c). The uncertain state indicates that f_r
64 does not provide any information about the simi-
65 larity between (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) .

66 To quantitatively evaluate whether (\mathbf{x}_i, y_i) can be con-
67 sidered to have been generated by the regression function
68 f_r , we use the likelihood $p(O_i|f_r)$, the probability of event
69 O_i that a data instance (\mathbf{x}_i, y_i) is observed, considering
70 f_r . The likelihood $p(O_i|f_r)$ is modeled using a normal
71 distribution with mean and standard deviation depend-
72 ing on the predicted target value $\hat{y}_i = f_r(\mathbf{x}_i)$ and the cor-
73 responding standard error $\sigma_{\mathbf{x}_i}$ by f_r , respectively. This
74 is expressed as:

$$p(O_i|f_r) = \begin{cases} 1 & \text{if } \Delta_i \leq 3\bar{\sigma} \\ 2 \times \int_{\Delta_i - 3\bar{\sigma}}^{+\infty} \mathcal{N}(u|0, \alpha\sigma_{\mathbf{x}_i}) du & \text{otherwise} \end{cases}, \quad (1)$$

75 where $\Delta_i = |y_i - \hat{y}_i| = |y_i - f_r(\mathbf{x}_i)|$ is the deviation
76 from the true to the predicted target values of data in-
77 stance i using f_r , and $\bar{\sigma}$ is the average of the predictive
78 standard error of all the data instances in \mathcal{D}_{ref} . α is
79 the hyperparameter used to adjust the condition that re-
80 stricts the data instances belonging to the function f_r .
81 In other words, the interval that determines the proba-
82 bility that a data instance (\mathbf{x}_i, y_i) belongs to f_r is $\alpha\sigma_{\mathbf{x}_i}$,
83 and if the data instance falls outside this interval, it is
84 determined that it does not belong to f_r . By increasing
85 or decreasing the value of the parameter α , the condition

1 for determining whether a data instance (\mathbf{x}_i, y_i) belongs
2 to f_r is relaxed or tightened, making $p(O_i|f_r)$ larger or
3 smaller, respectively. Optimal values of α can be chosen
4 using statistical criteria and appropriate validation meth-
5 ods; however, we set $\alpha = 2$ for all experiments in this
6 work to reduce model complexity. We consider $p(O_i|f_r)$
7 as the probability that (\mathbf{x}_i, y_i) is generated by f_r , and
8 $p(\bar{O}_i|f_r) = 1 - p(O_i|f_r)$ is the probability that (\mathbf{x}_i, y_i) is
9 not generated by f_r . Supplementary Figure 1 illustrates
10 the process of modeling the probability $p(O_i|f_r)$.

11 Events where (\mathbf{x}_i, y_i) or (\mathbf{x}_j, y_j) is generated by the
12 function f_r are independent events. Therefore, consider-
13 ing the function f_r , we can evaluate the joint probabilities
14 of observing:

- 15 • Both data instances:

$$p(O_i, O_j|f_r) = p(O_i|f_r) \times p(O_j|f_r); \quad (2)$$

- 16 • Only one of the data instances:

$$\begin{aligned} & p(O_i, \bar{O}_j|f_r) + p(\bar{O}_i, O_j|f_r) \\ &= p(O_i|f_r) \times p(\bar{O}_j|f_r) + p(\bar{O}_i|f_r) \times p(O_j|f_r); \end{aligned} \quad (3)$$

- 17 • Neither of the data instances:

$$\begin{aligned} & p(\bar{O}_i, \bar{O}_j|f_r) = p(\bar{O}_i|f_r) \times p(\bar{O}_j|f_r) \\ &= 1 - p(O_i, O_j|f_r) - p(O_i, \bar{O}_j|f_r) - p(\bar{O}_i, O_j|f_r). \end{aligned} \quad (4)$$

18 B. Modeling evidence by mass functions

19 Considering the Dempster–Shafer theory framework¹⁶,
20 we begin by defining the frame of discernment Ω . Let
21 $\Omega = \{s, ds\}$ be the universal set representing the similar-
22 ity states of any two data instances (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) .
23 s and ds denote the similarity and dissimilarity states
24 between the two data instances, respectively.

25 According to the Dempster–Shafer theory, the evidence
26 of the similarity states between these two data instances
27 is represented by a mass function $m^{i,j}$ (or a basic proba-
28 bility assignment)¹⁶. This assigns probability masses to
29 all the nonempty subsets of Ω ($\mathcal{X} = \{\{s\}, \{ds\}, \{s, ds\}\}$).
30 It is defined as follows:

$$m^{i,j} : \mathcal{X} \rightarrow [0, 1] \text{ with } \sum_{E \in \mathcal{X}} m(E) = 1. \quad (5)$$

31 The masses assigned to $\{s\}$ and $\{ds\}$ reflect the degrees of
32 belief exactly committed to the evidence to support the
33 similarity and dissimilarity between (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) ,
34 respectively. The weight assigned to $\{s, ds\}$ expresses the
35 degree of belief that the evidence provides no information
36 about the similarity (or dissimilarity) between (\mathbf{x}_i, y_i)
37 and (\mathbf{x}_j, y_j) .

38 Therefore, the mass function $m^{i,j}_{f_r}$, which models a
39 piece of evidence of the similarity between (\mathbf{x}_i, y_i) and

40 (\mathbf{x}_j, y_j) collected from f_r , is defined as follows:

$$m^{i,j}_{f_r}(\{s\}) = \frac{p(O_i, O_j|f_r)}{\gamma_{i,j}} \quad (6)$$

$$m^{i,j}_{f_r}(\{ds\}) = \frac{p(O_i, \bar{O}_j|f_r) + p(\bar{O}_i, O_j|f_r)}{\gamma_{i,j}} \quad (7)$$

$$m^{i,j}_{f_r}(\{s, ds\}) = 1 - \frac{1}{\gamma_{i,j}} + \frac{p(\bar{O}_i, \bar{O}_j|f_r)}{\gamma_{i,j}}, \quad (8)$$

41 where $\gamma_{i,j} = (e^{\frac{\Delta_y}{\bar{\sigma}}} + 1) \times (\frac{\sigma_{\mathbf{x}_i}}{\bar{\sigma}} + 1) \times (\frac{\sigma_{\mathbf{x}_j}}{\bar{\sigma}} + 1)$ is a discount-
42 ing factor^{16,26}, which describes the unreliability of evi-
43 dence about the similarity between (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j)
44 collected from a source of evidence f_r . Δ_y is the varia-
45 tion range of the target variable y in the dataset \mathcal{D} . The
46 smaller $\bar{\sigma}$ is relative to Δ_y , the more reliable the learned
47 regression function f_r is. Also, when $\sigma_{\mathbf{x}_i}$ and $\sigma_{\mathbf{x}_j}$ are
48 smaller than $\bar{\sigma}$, f_r can provide reliable evidence for the
49 relationship between (\mathbf{x}_i, y_i) and (\mathbf{x}_j, y_j) . By contrast,
50 when $\sigma_{\mathbf{x}_i}$ and $\sigma_{\mathbf{x}_j}$ are large compared to $\bar{\sigma}$, f_r cannot pro-
51 vide reliable evidence for the relationship between (\mathbf{x}_i, y_i)
52 and (\mathbf{x}_j, y_j) . A detailed explanation of each component
53 in $\gamma_{i,j}$ is provided in Supplementary Section I.

54 C. Dempster’s rule in combining evidence

55 Assuming that we can collect q pieces of evidence from
56 $\mathcal{F}_r = \{f_r^1, \dots, f_r^q\}$, a set of q reference functions is gen-
57 erated from \mathcal{D} to evaluate the similarity between a pair
58 of data instances with indices i and j . According to the
59 Dempster–Shafer theory framework, any two pieces of evi-
60 dence collected from the reference functions f_r^l and f_r^k ,
61 which are modeled by the corresponding mass functions
62 $m^{i,j}_{f_r^l}$ and $m^{i,j}_{f_r^k}$, respectively, can be combined using the
63 Dempster rule of combination to assign the joint mass
64 $m^{i,j}_{\{f_r^l, f_r^k\}}$ to each nonempty subset E of Ω as follows:

$$\begin{aligned} m^{i,j}_{\{f_r^l, f_r^k\}}(E) &= (m^{i,j}_{f_r^l} \oplus m^{i,j}_{f_r^k})(E) \\ &= \frac{\sum_{E_t \cap E_v = E} m^{i,j}_{f_r^l}(E_t) \times m^{i,j}_{f_r^k}(E_v)}{1 - \sum_{E_t \cap E_v = \emptyset} m^{i,j}_{f_r^l}(E_t) \times m^{i,j}_{f_r^k}(E_v)}, \end{aligned} \quad (9)$$

65 where E , E_t , and E_v are nonempty subsets of Ω . Demp-
66 ster’s rule is commutative and associative.

67 Based on Dempster’s rule, the obtained mass functions
68 corresponding to the q pieces of evidence are combined
69 to assign the final mass $m^{i,j}_{\mathcal{F}_r}$ as follows:

$$m^{i,j}_{\mathcal{F}_r}(E) = \left(m^{i,j}_{f_r^1} \oplus m^{i,j}_{f_r^2} \oplus \dots \oplus m^{i,j}_{f_r^q} \right)(E). \quad (10)$$

70 We perform similar analyses for all pairs of data in-
71 stances in \mathcal{D} to construct symmetric matrices M com-
72 prising the similarities ($M[i, j] = M[j, i] = m^{i,j}_{\mathcal{F}_r}(\{s\})$)

1 between them. Thereafter, the obtained matrix is ap-
 2 plied for further unsupervised data mining analysis, such
 3 as clustering or data visualization.

4 III. EXPERIMENTS AND RESULTS

5 In this section, we perform three experiments to
 6 demonstrate the application of our similarity measure-
 7 ment in dealing with outliers and data generated by mul-
 8 tiple mechanisms when designing materials descriptors.
 9 We apply the eRSM to measure similarities between mag-
 10 netic of three datasets for detecting subgroups of ma-
 11 terials: 1) The experimentally observed Curie tempera-
 12 ture dataset (\mathcal{D}_{binary}) of binary alloys for transitioning
 13 rare earth metals, 2) Dataset of calculated magnetiza-
 14 tion of quaternary high-entropy alloys ($\mathcal{D}_{quaternary}^{Mag}$), and
 15 3) Dataset of calculated Curie temperature of quaternary
 16 high-entropy alloys ($\mathcal{D}_{quaternary}^{TC}$). Note that the datasets
 17 $\mathcal{D}_{quaternary}^{Mag}$ and $\mathcal{D}_{quaternary}^{TC}$ contain similar alloys and
 18 differ only in the target properties.

19 A. Datasets

20 The details of the datasets investigated in this study
 21 are as follows.

- 22 • Binary alloys dataset \mathcal{D}_{binary} ²⁷: A material dataset
 23 containing 100 transition-rare earth metal binary
 24 alloys, comprising nickel (Ni), manganese (Mn),
 25 cobalt (Co), or iron (Fe), and the corresponding
 26 Curie temperatures (T_C). This dataset was col-
 27 lected from the Atomwork database of the National
 28 Institute of Materials Science^{28,29}. Each binary al-
 29 loy in \mathcal{D}_{binary} is represented using seven descrip-
 30 tors: (1,2) the atomic number of transition metal
 31 (Z_T) and rare-earth (Z_R) constituents; (3) projec-
 32 tion of the spin magnetic moment onto the total
 33 angular momentum of the $4f$ electrons ($J_{4f}(1 - g_j)$);
 34 (4, 5) covalent radius (r_{covT}) and first ionization
 35 (IP_T) of the transition metal; (6, 7) concentration
 36 of the transition metal (C_T) and rare-earth metal
 37 (C_R). The selection of these seven descriptors has
 38 been discussed in detail in previous studies^{10,30}.
- 39 • Quaternary high-entropy alloys datasets
 40 $\mathcal{D}_{quaternary}$ ²⁷: A material dataset contains
 41 990 equiatomic quaternary high-entropy alloys,
 42 which comprise 14 transition metals $\{Ag, Cd, Co,$
 43 $Cr, Cu, Fe, Mn, Mo, Ni, Pd, Rh, Ru, Tc, Zn\}$,
 44 and the corresponding calculated magnetizations
 45 and Curie temperatures in the BCC phase. The
 46 dataset was collected from an original dataset
 47 of 147,630 equiatomic quaternary high-entropy
 48 alloys calculated using Korringa-Kohn-Rostoker
 49 coherent approximation method³¹. Each alloy in
 50 $\mathcal{D}_{quaternary}$ is represented using 135 compositional

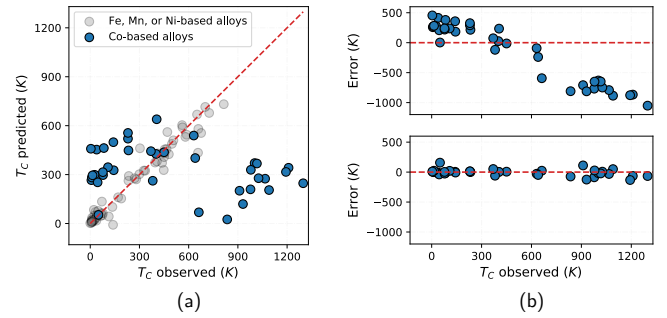


FIG. 2. (a) Observed and predicted Curie temperature of alloys in the dataset \mathcal{D}_{binary} using model generated for nickel (Ni), iron (Fe), and manganese (Mn)-based alloys. The blue and gray points indicate cobalt (Co)-based alloys and alloys of other transition metals (Ni, Fe, Mn), respectively. (b) Prediction error of Co-based alloys when excluding (top) or including (bottom) data of other Co-based alloys to the training dataset.

51 descriptors, including the means, standard devia-
 52 tions, and covariance of the atomic representations
 53 of their constituent elements¹³ and four categorical
 54 features indicating the elements comprising the
 55 quaternary alloy. The feature selection process
 56 applied to this dataset has been discussed in detail
 57 in Supplementary Section III.

58 B. Assessment of the similarity between transition-rare earth 59 metal binary alloys based on mechanisms of Curie temperature

60 In the first experiment, we show the versatility of the
 61 eRSM for detecting outliers and identifying a mixture of
 62 mechanisms. We apply the eRSM to assess the similar-
 63 ities between 100 transition rare earth metal binary alloys
 64 comprising nickel (Ni), manganese (Mn), cobalt (Co), or
 65 iron (Fe) in the dataset \mathcal{D}_{binary} based on their Curie tem-
 66 peratures. We can construct a regression model using a
 67 Gaussian process by considering the data instances in
 68 \mathcal{D}_{binary} . This shows a high prediction accuracy with an
 69 R^2 score of 0.963 and an MAE of 40 (K) in *ten*-fold
 70 cross-validation. However, such a nonparametric regres-
 71 sion model does not guarantee the reliability of the model
 72 in the subsequent exploratory predictions. This is be-
 73 cause the number of observable alloys is relatively small
 74 compared to the number of possible alloys.

75 Figure 2 (a) shows the results of the exploratory predic-
 76 tion of the Curie temperature of the Co-based binary al-
 77 loys in \mathcal{D}_{binary} using a Gaussian process regression model
 78 constructed from the data of binary alloys of Ni, Mn, and
 79 Fe. The regression model constructed from the data of
 80 binary alloys of Ni, Mn, and Fe shows a high predic-
 81 tion accuracy in *ten*-fold cross-validation ($R^2 = 0.946$
 82 and MAE= 35 (K)). Although the Co-based alloys with
 83 high Curie temperature tend to be underestimated by
 84 the model, the other Co-based alloys are often overesti-

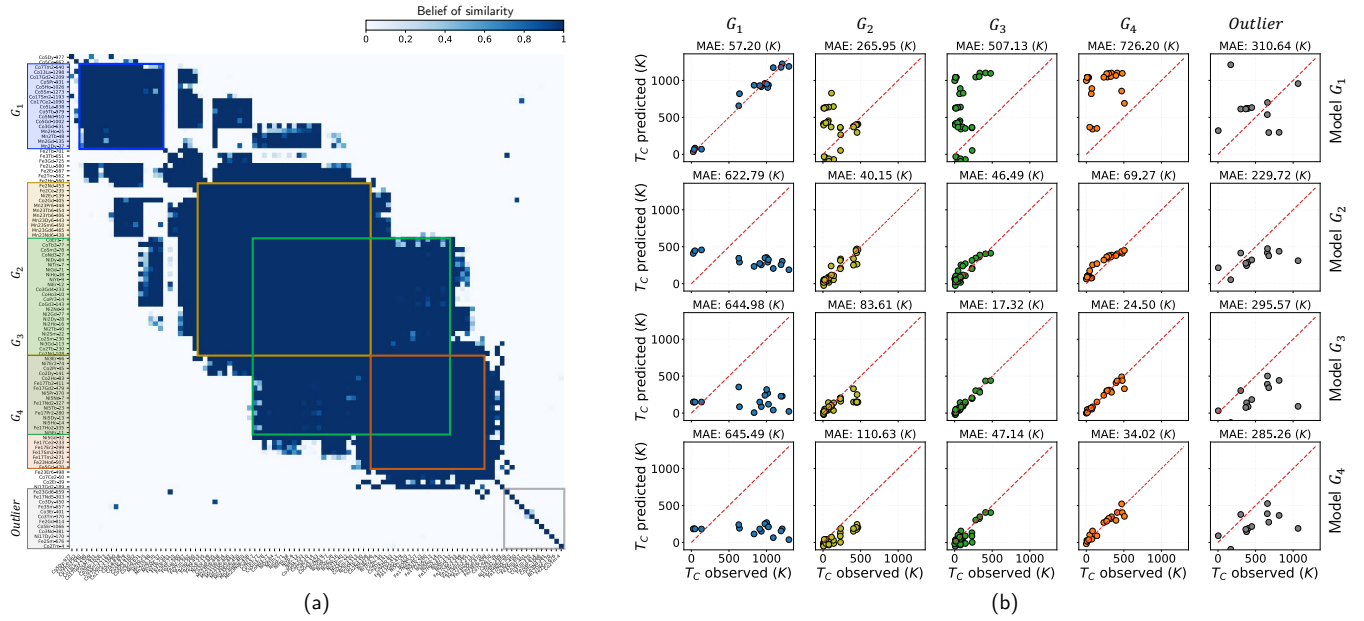


FIG. 3. (a) Heatmap illustrating the similarity matrix M_{binary} extracted for all the data instances in the \mathcal{D}_{binary} . (b) Confusion matrices measuring the regression-based similarities between alloys in four groups G_1 - G_4 and the dissimilarities between the models generated for alloys in different groups.

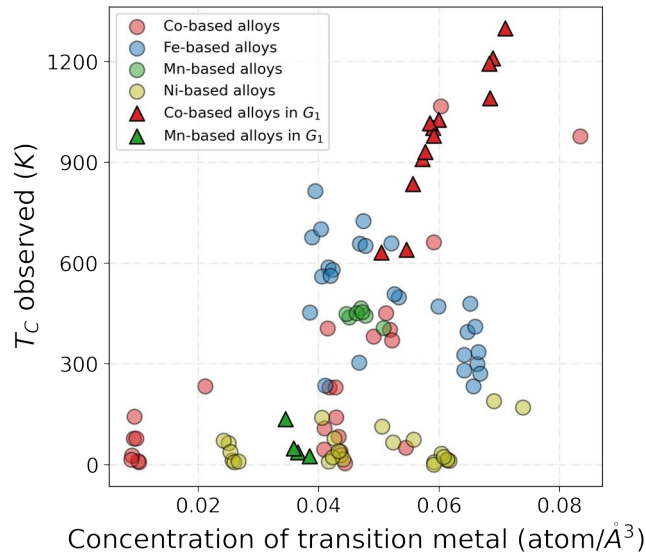


FIG. 4. Dependence of T_C on the concentration of the transition metal (C_T) in alloys. Red, blue, green and yellow scatters indicate alloys containing cobalt (Co), iron (Fe), manganese (Mn), and nickel (Ni). Alloys in G_1 are highlighted by triangles.

1 mated. The prediction error for the Co-based alloys is
 2 critically reduced when some data of the other Co-based
 3 alloys is included (Fig. 2 b). This observation supports
 4 the hypothesis that the underlying mechanisms are dif-
 5 ferent between the Co-based alloys and alloys of other
 6 transition metals. This facilitates the use of the eRSM

7 to clarify the mixture mechanism from this dataset.

8 By applying the eRSM on the dataset \mathcal{D}_{binary} , we
 9 obtain a similarity matrix M_{binary} with moderately high
 10 similarity values among the data instances (Fig. 3 a).
 11 Thus, approximately all the data instances can be re-
 12 gressed by a relatively smooth function. This is consist-
 13 ent with the high prediction accuracy of *ten-fold* cross-
 14 validation for all the alloys in the dataset. Considering
 15 the exploratory data analysis, to avoid false intuition or
 16 misunderstanding, the grouping of alloys in \mathcal{D}_{binary} is
 17 done such that the similarities between the alloys in each
 18 group are high. Moreover, one alloy can belong to more
 19 than one group simultaneously, or it can be in none of the
 20 groups. We apply a graph-based clustering method³² to
 21 the extracted similarity matrix to detect overlapping sub-
 22 groups of materials. As a result, we observe four groups
 23 of alloys, denoted as G_1 , G_2 , G_3 , and G_4 , which show
 24 high intragroup similarities, exceeding 0.7 (Fig. 3 a).
 25 Nevertheless, the similarity between the alloys in group
 26 G_1 and those in G_2 , G_3 , and G_4 is significantly dissimil-
 27 ar. In addition, a small group of alloys (Fig. 3 a, gray
 28 region) is approximately different from all the others and
 29 can be considered as outliers. The remaining alloys are
 30 not assigned to any group to have confidence in the clus-
 31 tering analysis results.

32 To evaluate the validity of the analysis process quan-
 33 titatively, we trained the regression models for T_C us-
 34 ing data from each of the four groups G_1 , G_2 , G_3 , and
 35 G_4 . Moreover, we monitored their prediction accuracy
 36 on these groups. The confusion matrix summarizing the
 37 correlation between the observed and predicted T_C by
 38 the four learned regression models is shown in Fig. 4.

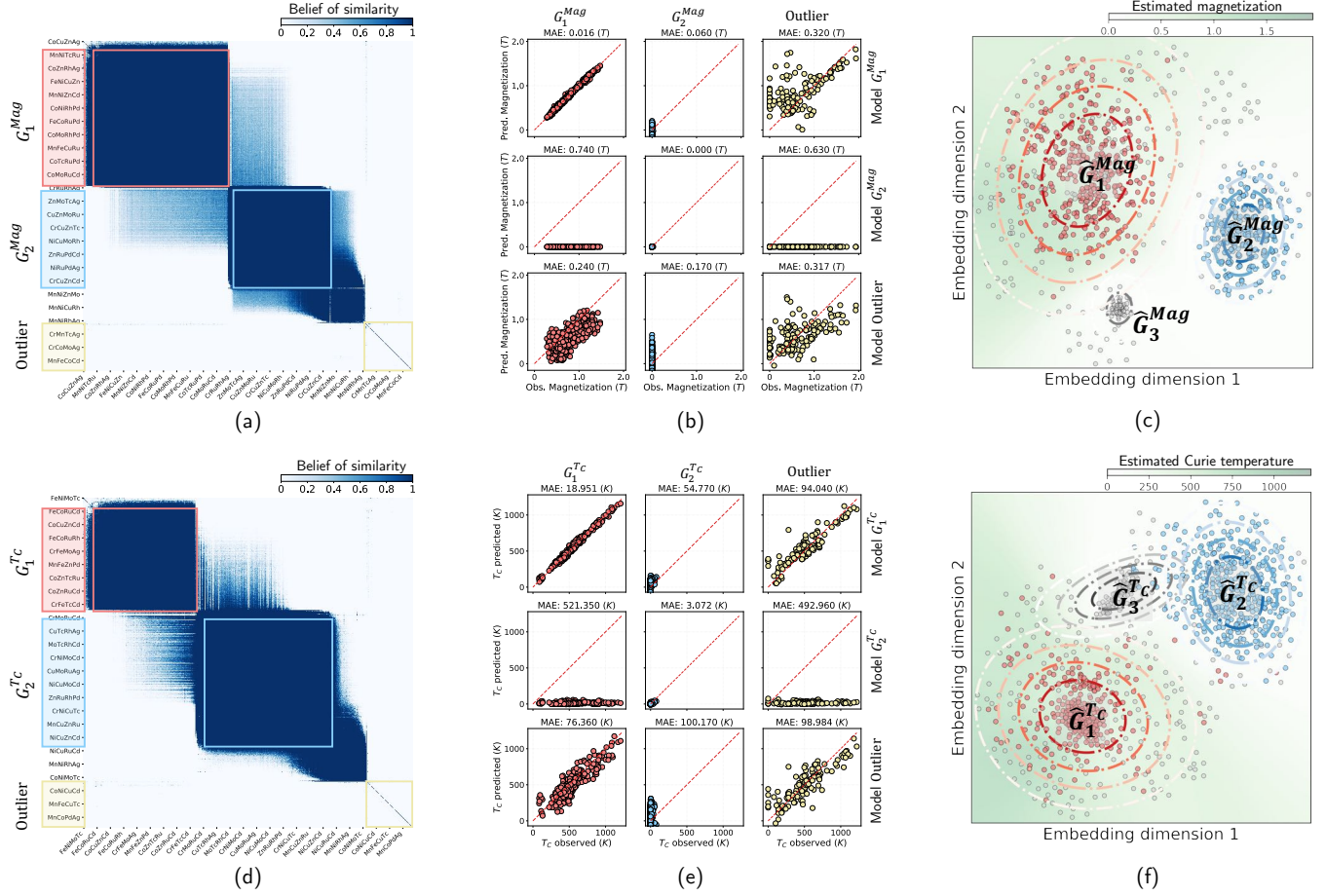


FIG. 5. (a,d) Heatmaps illustrating the similarity matrices $M_{quaternary}^{Mag}$ (a) and $M_{quaternary}^{TC}$ (d) extracted from datasets $\mathcal{D}_{quaternary}^{Mag}$ and $\mathcal{D}_{quaternary}^{TC}$, focusing on mechanisms of magnetization and T_C , respectively. (b,e) The confusion matrix summarizes the differences between the magnetization (b) or T_C (e) mechanisms of alloys in extracted groups. (c,f) Visualization of quaternary alloys in the two-dimensional embedding spaces constructed by applying the T-distributed Stochastic Neighbor Embedding (t-SNE) to $M_{quaternary}^{Mag}$ (c) and $M_{quaternary}^{TC}$ (f). Red, blue, and gray contours indicate Gaussian models \hat{G}_1^{Mag} (\hat{G}_1^{TC}), \hat{G}_2^{Mag} (\hat{G}_2^{TC}), and \hat{G}_3^{Mag} (\hat{G}_3^{TC}), respectively, learned by using the Gaussian Mixture Models³³ in the embedding space focusing on mechanisms of magnetization (T_C). In addition, red and blue points in sub-figures b and c (e and f) indicate the alloys in G_1^{Mag} (G_1^{TC}) and G_2^{Mag} (G_2^{TC}), respectively.

1 The diagonal plots illustrate the cross-validation results
 2 of the models learned from the four groups of alloys. The
 3 off-diagonal plot shows the correlation between the ob-
 4 served T_C and the predictions made by the model learned
 5 from the alloys of the other groups. The obtained results
 6 confirm the intragroup similarity of the alloys in groups
 7 G_1 , G_2 , G_3 , and G_4 , respectively, dissimilarity between
 8 the five groups, and intra-group dissimilarity of the alloys
 9 considered as outliers. This indicates that the obtained
 10 results suggest that the physical mechanisms of alloys in
 11 G_1 may be different from those of the alloys in G_2 , G_3 ,
 12 and G_4 . Nonetheless, it is difficult to determine the dif-
 13 ferences between the mechanisms of the T_C of alloys in
 14 G_2 , G_3 , and G_4 .

15 Moreover, considering the alloys in G_1 , there is a strong
 16 linear correlation between T_C and the concentration of
 17 transition metals in the alloys with a Pearson correla-

18 tion coefficient of 0.95 (Fig. 4, triangle scatters). This
 19 result is consistent with the observation of the previous
 20 research³⁰, when considering all binary alloys of transi-
 21 tion metals and rare earth metals in \mathcal{D}_{binary} ; the range
 22 of T_C is found to be correlated with the composition ra-
 23 tio of the transition metals. Furthermore, 13 of the 17
 24 alloys in G_1 are Co-based alloys with high Curie tem-
 25 peratures ($T_C > 600K$). By contrast, most of the other
 26 Co-based alloys in \mathcal{D}_{binary} have lower Curie temperatures
 27 ($T_C < 500K$) and are assigned to G_2 , G_3 , and G_4 . These
 28 results are consistent with the observation that the re-
 29 gression model for Fe-, Mn-, and Ni-based alloys tends to
 30 underestimate the T_C of the Co-based alloys with high
 31 T_C and overestimates the T_C of the remaining Co-based
 32 alloys (Fig. 2 a).

33 In addition, we examine the behavior of eRSM on
 34 toy datasets synthesized with outliers or multiple mech-

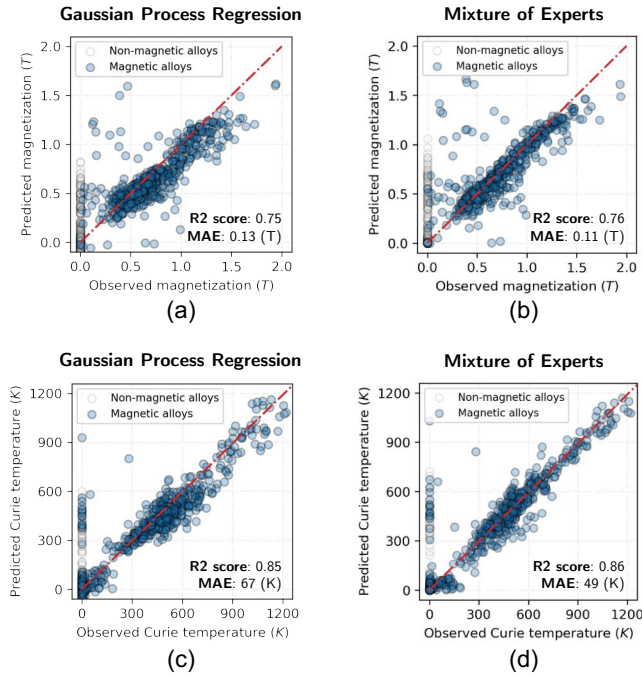


FIG. 6. Prediction accuracies for magnetization (a, b) and Curie temperature (c, d) of the alloys with 10-fold cross-validations. Prediction validation results with single gaussian process regression models for magnetization and Curie temperature are shown in sub-figures (a) and (c), respectively. Prediction validation results with mixtures of expert models for magnetization and Curie temperature are shown in sub-figures (b) and (d), respectively. Blue and white circles indicate magnetic alloys (finite magnetization) and non-magnetic alloys (zero magnetization), respectively.

1 anisms to assess the efficiency of this similarity measure.
 2 Detailed results of these experiments are summarized in
 3 the Supplementary Section II. Briefly, the eRSM demon-
 4 strates that it can effectively assess the similarity between
 5 the data instances and use the similarity for detecting
 6 outliers and a mixture of mechanisms.

7 C. Assessment of the similarity between quaternary 8 high-entropy alloys based on mechanisms of magnetization

9 The effectiveness of the eRSM in detecting outliers and
 10 identifying mixture mechanisms in the material dataset
 11 has been shown in the previous experiment. In the next
 12 two experiments, we show the potential of applying the
 13 measured similarity to design descriptors for materials.

14 Considering this experiment, we subsequently apply
 15 the eRSM to assess the similarities between 990 quater-
 16 nary high-entropy alloys comprising 14 transition met-
 17 als in the dataset $\mathcal{D}_{quaternary}^{Mag}$ based on their magnetiza-
 18 tion. To predict the magnetization of these alloys, we
 19 attempted to construct an optimal Gaussian process re-
 20 gression model using the designed descriptors. The Gaus-

21 sian process can poorly regress the magnetization with an
 22 R^2 score of 0.75 and an MAE of 0.13 (T) in the ten-fold
 23 cross-validation. The obtained results suggest that the
 24 magnetization of these alloys may not be described by a
 25 single model in the designed descriptor space. This in-
 26 dicates that the existence of outliers or mixture models
 27 of the magnetization properties of these alloys in the de-
 28 scriptor space should be considered in the analysis of this
 29 dataset.

30 Applying the eRSM, we obtain a similarity matrix
 31 $M_{quaternary}^{Mag}$ with two core groups of alloys denoted by
 32 G_1^{Mag} and G_2^{Mag} , showing high intra-group similarities
 33 and exceeding 0.5 (Fig. 5 a). Some of the alloys in G_1^{Mag}
 34 are similar to those in G_2^{Mag} ; nonetheless, the rest show
 35 apparent dissimilarities. Furthermore, one small group
 36 of alloys (Fig. 5 a, yellow region) showed dissimilarities
 37 with the others and could be considered as outliers.
 38 The remaining alloys in $\mathcal{D}_{quaternary}^{Mag}$ do not exhibit ap-
 39 parent similarities with alloys in groups G_1^{Mag} and G_2^{Mag} .
 40 Therefore, they are not assigned to any group.

41 Similar to the previous session, to validate the obtained
 42 results quantitatively, we trained three regression mod-
 43 els using data from each group, G_1^{Mag} , G_2^{Mag} , and out-
 44 liers. We monitored the prediction accuracy of the three
 45 learned regression models for data in all the groups. The
 46 confusion matrix summarizing the correlations between
 47 the observed and predicted values of the target variable
 48 using the learned regression models is shown in Fig. 5 (c).
 49 The diagonal plots illustrate the ten-fold cross-validation
 50 results of the models learned from these three groups of
 51 alloys. The off-diagonal plot shows the correlation be-
 52 tween the observed magnetization and the predictions
 53 made by the model learned from the alloys of the other
 54 groups.

55 The obtained results confirm the intragroup similarity
 56 of the alloys in groups G_1^{Mag} and G_2^{Mag} , respectively, the
 57 dissimilarity between the two groups, and the intra-group
 58 dissimilarity of the alloys considered as outliers. Specif-
 59 ically, we observe that group G_2^{Mag} consists of ferrimag-
 60 netic alloys or alloys whose magnetization is relatively
 61 smaller (magnetization < 0.1 (T)) than the others in the
 62 group G_1^{Mag} . In contrast, using the data in G_1^{Mag} , we
 63 can construct a Gaussian process regression model with
 64 a high prediction accuracy with an R^2 score of 0.992 and
 65 an MAE of 0.016 (T) in the ten-fold cross-validation.

66 Therefore, we can use the information of the con-
 67 stituent elements of each alloy to predict which group it
 68 belongs to in advance²⁰ and apply an appropriate regres-
 69 sion model to improve prediction accuracy for the alloys.
 70 We combine the similarity measured by using the eRSM
 71 with the Jaccard similarity coefficient³⁴ and apply the T-
 72 distributed Stochastic Neighbor Embedding³⁵ (t-SNE) to
 73 construct a two-dimensional embedding map (Fig. 5 c).
 74 Details of the combination method are shown in Supple-
 75 mentary Section IV. As a result, we can easily distinguish
 76 the alloys in groups G_1^{Mag} (red) and G_2^{Mag} (blue) when
 77 they form two separate regions with high density in the

1 embedding space. We apply a Gaussian mixture model³³
 2 (GMM) on the embedding space to identify groups and
 3 calculate the probability of an alloy belonging to a par-
 4 ticular identified group. Alloys in different groups are
 5 treated differently by using a mixture of experts³⁶ (MoE)
 6 approach. Figure 6 (a-b) show a reduction of the pro-
 7 posed mixture of experts in MAE of 18% compared with
 8 result of the single model, from 0.13 (T) to 0.11 (T). Fur-
 9 ther analysis shows that applying the obtained similar-
 10 ities in MOE improves the prediction accuracy for mag-
 11 netic alloys (Supplementary figure 7 a).

12 D. Assessment of the similarity between the quaternary 13 high-entropy alloys based on mechanisms of Curie 14 temperature

15 Considering this experiment, the target data are the
 16 same as in the previous section ($\mathcal{D}_{quaternary}$); however,
 17 the physical property of interest is T_C . A regression
 18 model can be constructed using a Gaussian process. This
 19 shows a rather high prediction accuracy in *ten*-fold cross-
 20 validation with an R^2 score of 0.85 and an MAE of 67
 21 (K). We also observe two distinguishable groups of qua-
 22 ternary alloys in the dataset $\mathcal{D}_{quaternary}^{T_C}$ when applying
 23 the eRSM. Figure 5 (d) illustrates the similarity matrix
 24 $M_{quaternary}^{T_C}$ with two groups of alloys denoted as $G_1^{T_C}$
 25 and $G_2^{T_C}$, showing high intra-group similarities and ex-
 26 ceeding 0.5. Some of the alloys in $G_1^{T_C}$ are similar to
 27 those in $G_2^{T_C}$. Nonetheless, the others exhibit apparent
 28 dissimilarities, which is consistent with the observation
 29 of two high-density regions (red) in the embedding map
 30 of $M_{quaternary}^{T_C}$ (Fig. 5 e). Furthermore, a small group
 31 of alloys (Fig. 5 d, yellow region) showed dissimilarities
 32 with all the others and could be considered as outliers.
 33 The remaining alloys do not show apparent similarities
 34 with alloys in groups $G_1^{T_C}$ and $G_2^{T_C}$; thus, they are not
 35 assigned to any group.

36 Following the same analysis procedure as in the previ-
 37 ous section, we trained regression models for Curie tem-
 38 perature using data from each of the three groups $G_1^{T_C}$,
 39 $G_2^{T_C}$, and outliers and monitored their prediction accu-
 40 racy on these groups. Figure 5 (f) shows the confusion
 41 matrix that summarizes the obtained results. The diag-
 42 onal plots illustrate the *ten*-fold cross-validation results
 43 of the models learned from these three groups of alloys.
 44 The off-diagonal plot shows the correlation between the
 45 observed Curie temperature and the predictions made by
 46 the regression model learned from the alloys of the other
 47 groups. We can also confirm the intra-group similarity of
 48 the alloys in groups $G_1^{T_C}$ and $G_2^{T_C}$, respectively, dissim-
 49 ilarity between the two groups, and intra-group dissim-
 50 ilarity of the alloys considered as outliers.

51 Specifically, we observe that the Curie temperatures
 52 of approximately all the alloys in group $G_2^{T_C}$ have a low
 53 T_C , which is 0 (K) or relatively smaller than that of the
 54 other alloys. Furthermore, using the data in $G_1^{T_C}$, we

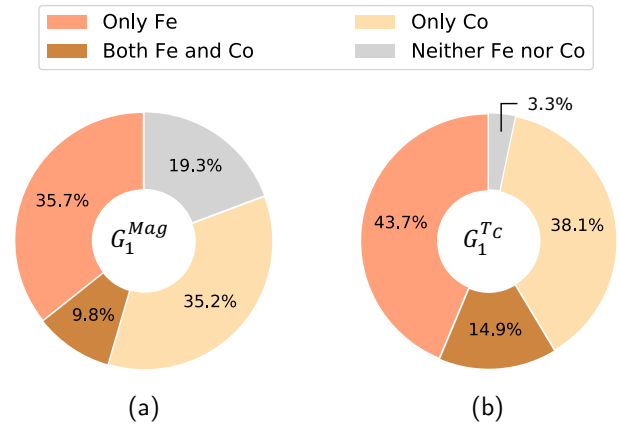


FIG. 7. Proportions of quaternary alloys containing Fe or Co in group G_1^{Mag} (a) and $G_1^{T_C}$ (b).

55 can construct a Gaussian process regression model with
 56 a high prediction accuracy with an R^2 score of 0.985 and
 57 an MAE of 19 (K) in the *ten*-fold cross-validation.

58 Therefore, we utilize the similarity information to de-
 59 sign descriptors for quaternary alloys due to the effec-
 60 tiveness of the data for detecting the mixture of multi-
 61 ple mechanisms in the dataset. We apply similar meth-
 62 ods as in the previous experiment to construct a two-
 63 dimensional embedding map (Fig. 5 f) and then learn a
 64 mixture of experts to predict Curie temperature of qua-
 65 ternary alloys in the dataset $\mathcal{D}_{quaternary}^{T_C}$. The proposed
 66 mixture of models exhibits higher prediction accuracy
 67 than the single model in *10*-folds cross-validations (Fig.
 68 6 c-d). The MAE of the proposed mixture of expert re-
 69 duces approximately 36%, from 67 (K) to 49 (K).

70 E. Discussion of the obtained similarities between materials 71 and the associated physical mechanisms

72 Regarding the experiments with the datasets
 73 $\mathcal{D}_{quaternary}^{Mag}$ and $\mathcal{D}_{quaternary}^{T_C}$ focusing on magnetiza-
 74 tion or T_C , the datasets seem to be a self-evident
 75 example where magnetization and T_C are cases sensitive
 76 to finite or zero. As we can see from the results described
 77 above (Sections III C, III D, and Supplementary Section
 78 VI), the prediction accuracy is low when considering a
 79 single regression model for the entire dataset. In this
 80 section, we pay attention to the analysis of the extracted
 81 alloys groups G_1^{Mag} , G_2^{Mag} , $G_1^{T_C}$, and $G_2^{T_C}$ to identify
 82 underlying patterns.

83 Figure 7 shows that Fe and Co, which have a large spin
 84 moment, ferromagnetic interactions with many elements
 85 and result in high magnetization or T_C , are dominant
 86 elements comprising alloys in two groups G_1^{Mag} (a) and
 87 $G_1^{T_C}$ (b). Furthermore, in the analysis that considers the
 88 proportion of the quaternary alloys fixing two of their
 89 four constituent elements concerning the extracted four

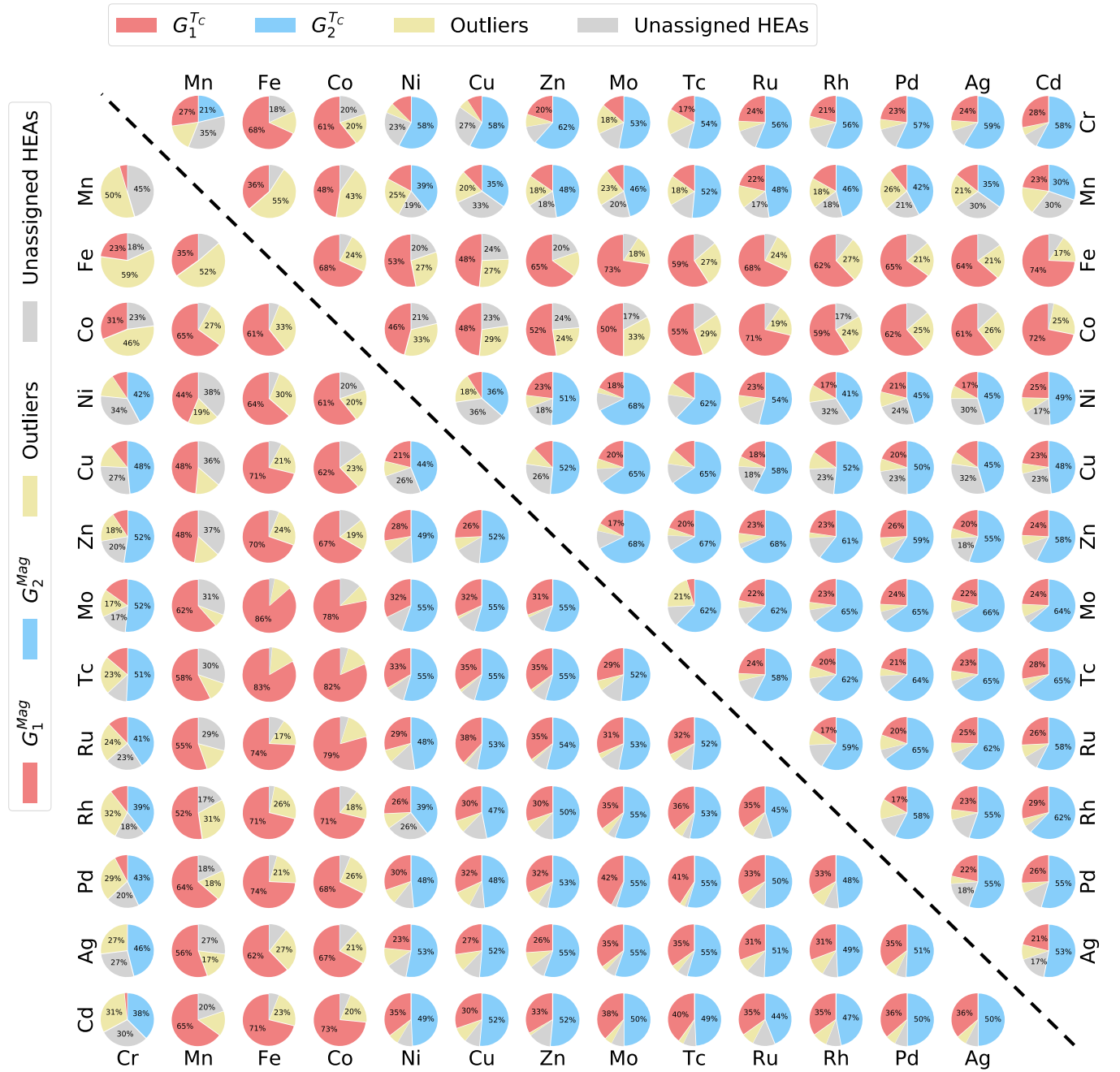


FIG. 8. Effect of coexistence of the 14 transition metals on magnetization and Curie temperature mechanisms. Each pie chart results from quaternary alloys containing the respective element pair. They show the percentages of alloys that follow the magnetization mechanisms (lower-left triangle) and Curie temperature mechanisms (upper-right triangle), as extracted by the eRSM. Red and blue areas indicate the percentages of alloys whose magnetization and T_C are finite (G_1^{Mag} and G_1^{Tc}) and zero (G_2^{Mag} and G_2^{Tc}), respectively. Yellow areas indicate the percentages of alloys that are detected as outliers. By contrast, gray regions indicate the fractions of alloys not assigned to the extracted groups.

1 groups G_1^{Mag} , G_2^{Mag} , G_1^{Tc} , and G_2^{Tc} , we observe that the
 2 proportion of Fe-containing and Co-containing alloys in
 3 two groups G_1^{Mag} (a) and G_1^{Tc} are significantly larger
 4 than other groups (Fig. 8). Thus, the prediction mod-
 5 els constructed from the data of the alloys in G_1^{Mag} or
 6 G_1^{Tc} are more suitable to predict magnetization or T_C ,

7 respectively, of alloys containing these elements. The re-
 8 maining Fe-X and Co-X (X denotes the other transition
 9 metals comprised in the alloys) alloys are considered out-
 10 liers of the extracted mechanisms or unassigned HEAs,
 11 which are not assigned to any of these mechanisms. Con-
 12 versely, Mn-X alloys exhibit similar behavior as Fe-X and

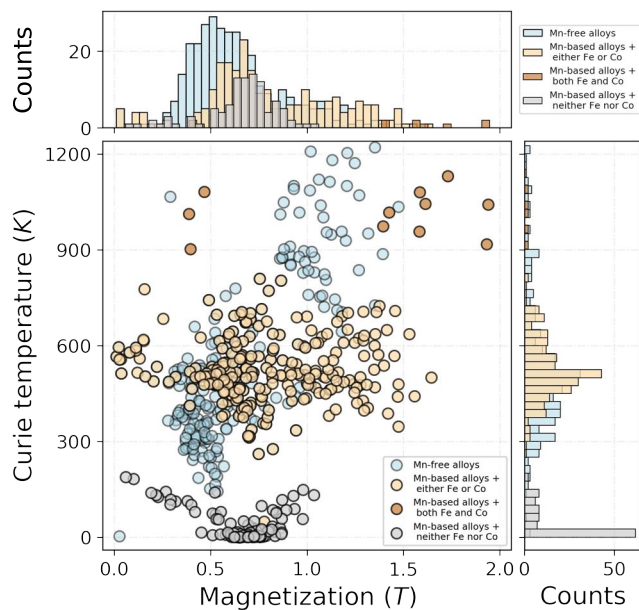


FIG. 9. Correlation between magnetization (T) and Curie temperature (K) of quaternary alloys with non-zero magnetization and non-zero Curie temperature in datasets $\mathcal{D}_{quaternary}^{Mag}$ and $\mathcal{D}_{quaternary}^{TC}$. Marginal plots show histogram of the properties of the alloys.

1 Co-X when focusing on the magnetization mechanisms.
 2 However, for the Curie temperature, the Mn-X alloys are
 3 categorized in the group G_2^{TC} of low T_C besides the other
 4 groups. Especially among the Fe-X and Co-X alloys, the
 5 percentage of Fe-Mn and Co-Mn alloys are considered as
 6 outliers of the mechanisms extracted from G_1^{TC} are rela-
 7 tively higher, 55% and 43%, respectively (Fig. 8).

8 For further investigation, we organized the raw data
 9 of the quaternary alloys by focusing on the presence or
 10 absence of Mn. Figure 9 shows the correlation between
 11 magnetization and Curie temperature of 556 (56%) al-
 12 loys with non-zero properties. The total number of data
 13 instances is 990, and the number of data instances where
 14 both T_C and magnetization are zero is 413 (42%), while
 15 there are twenty-one (2%) alloys with zero T_C but have
 16 finite magnetization. We found that the alloys contain-
 17 ing all three elements, Mn, Fe, and Co, show high Curie
 18 temperatures ($T_C > 900$ (K)). Conversely, the alloys
 19 containing either pairs of Mn-Fe or Mn-Co show moder-
 20 ate Curie temperatures. By contrast, the Mn-containing
 21 alloys without Fe or Co have low Curie temperatures
 22 ($T_C < 250$ (K)). Furthermore, the trends of these three
 23 alloy groups do not offer any significant correlation be-
 24 tween magnetization and Curie temperature. However,
 25 an apparent positive correlation between magnetization
 26 and Curie temperature can be observed for the group of
 27 Mn-free alloys.

28 To interpret the results obtained, we considered a hy-
 29 pothesis of the origin of the observed data. The esti-
 30 mated magnetization is the sum of all the local mag-

31 netic moments divided by the unit volume. The local
 32 magnetic moments are determined by the spin configura-
 33 tions of atomic sites that stabilize the structure of alloys.
 34 Conversely, given a particular structure and spin con-
 35 figuration, the T_C can be estimated from the spin-spin
 36 exchange energy. First-principles calculations show that
 37 early transition metals and late transition metals often
 38 have antiferromagnetic interactions³⁷. This interaction
 39 has also been confirmed in high-entropy alloys by using
 40 automatic exhaustive calculations³¹. Mn lies between
 41 early and late transition metals; thus, the estimation
 42 of the spin configuration (ferromagnetic or antiferromag-
 43 netic) in Mn-containing alloys should be cautiously con-
 44 sidered in different situations, especially in high-entropy
 45 alloys whose elements can stochastically exist at the same
 46 atomic site. From this consideration, we can admit a hy-
 47 pothesis that the alloys containing Mn follow a different
 48 rule for magnetization than those grouped into G_2^{Mag} .
 49 Conversely, the alloys containing Mn may follow the same
 50 rules for T_C as the alloys grouped into G_2^{TC} , albeit with
 51 a spin configuration that provides magnetization. The
 52 details are beyond the scope of this paper and will not
 53 be discussed here, but further analysis is promising.

54 IV. CONCLUSIONS

55 In this study, we developed a method that can be
 56 used to rationally transform material data from multi-
 57 ple sources into evidence of similarities between materi-
 58 als and combine the evidence to conclude the similarities
 59 between materials. The extracted similarity-dissimilarity
 60 information has significant potential for application in
 61 subgroups discovery of materials. The effectiveness of
 62 the eRSM in detecting homogenous subgroups of materi-
 63 als has been demonstrated by using two experiments on
 64 two datasets of magnetic materials. In addition, further
 65 analysis of the detected subgroups improves the existing
 66 knowledge of problems related to the applied datasets of
 67 magnetic materials. For example, we reveal the differ-
 68 ences in the mechanisms of the Curie temperature of Co-
 69 based binary alloys when using our method to a dataset
 70 of 100 transition-rare earth metal binary alloys compris-
 71 ing Ni, Mn, Co, and Fe. Moreover, we explored the
 72 mechanisms of ferrimagnetic and low Curie temperature
 73 alloys from the magnetic dataset of calculated quater-
 74 nary alloys. By measuring the similarity between ma-
 75 terials with uncertainty, the method described herein is
 76 expected to extract valuable information for describing
 77 and interpreting the underlying physical mechanisms in
 78 material datasets.

79 SUPPLEMENTARY MATERIAL

80 See supplementary materials for the following addi-
 81 tional information: 1) Explanation of the formulation
 82 modeling uncertainty, 2) Evaluation of the eRSM us-

ing the toy datasets, and 3) Features selection and pre-analysis in the dataset of quaternary high-entropy alloys.

ACKNOWLEDGMENTS

This work is supported by the Ministry of Education, Culture, Sports, Science, and Technology of Japan (MEXT) with the Program for Promoting Research on the Supercomputer Fugaku (DPMSD), JSPS KAKENHI Grants 20K05301, JP19H05815 (Grants-in-Aid for Scientific Research on Innovative Areas Interface Ionics), 21K14396 (Grant-in-Aid for Early-Career Scientists), and 20K05068, Japan.

DATA AVAILABILITY STATEMENT

Datasets related to this article are deposited to Zenodo repository²⁷.

REFERENCES

- ¹B. Kailkhura, B. Gallagher, S. Kim, A. Hiszpanski, and T. Y.-J. Han, “Reliable and explainable machine-learning methods for accelerated material discovery,” *npj Computational Materials* **5**, 108 (2019).
- ²J. B. Tenenbaum, “Learning the structure of similarity,” *Advances in Neural Information Processing Systems* **8**, 3–9 (1996).
- ³J. Tenenbaum, V. Silva, and J. Langford, “A global geometric framework for nonlinear dimensionality reduction,” *Science* **290**, 2319–2323 (2000).
- ⁴Y. Yang, F. Liang, S. Yan, Z. Wang, and T. S. Huang, “On a theory of nonparametric pairwise similarity for clustering: Connecting clustering to classification,” *Advances in Neural Information Processing Systems* **27**, 145–153 (2014).
- ⁵C. Chen, O. Li, D. Tao, A. Barnett, C. Rudin, and J. K. Su, “This looks like that: Deep learning for interpretable image recognition,” in *Advances in Neural Information Processing Systems*, Vol. 32, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc., 2019).
- ⁶B. Letham, C. Rudin, T. H. McCormick, and D. Madigan, “Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model,” *Ann. Appl. Stat.* **9**, 1350–1371 (2015).
- ⁷C. Rudin, “Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead,” *Nature Machine Intelligence* **1**, 206–215 (2019).
- ⁸B. R. Goldsmith, M. Boley, J. Vreeken, M. Scheffler, and L. M. Ghiringhelli, “Uncovering structure-property relationships of materials by subgroup discovery,” *New Journal of Physics* **19**, 013031 (2017).
- ⁹R. Ramprasad, R. Batra, G. Piliya, A. Mannodi-Kanakkithodi, and C. Kim, “Machine learning in materials informatics: recent applications and prospects,” *npj Computational Materials* **3**, 54 (2017).
- ¹⁰D.-N. Nguyen, T.-L. Pham, V.-C. Nguyen, T.-D. Ho, T. Tran, K. Takahashi, and H.-C. Dam, “Committee machine that votes for similarity between materials,” *IUCrJ* **5**, 830–840 (2018).
- ¹¹D.-N. Nguyen, T.-L. Pham, V.-C. Nguyen, H. Kino, T. Miyake, and H.-C. DAM, “Ensemble learning reveals dissimilarity between rare-earth transition binary alloys with respect to the curie temperature,” *Journal of Physics: Materials* (2019).
- ¹²J. Bardeen, L. N. Cooper, and J. R. Schrieffer, “Theory of superconductivity,” *Phys. Rev.* **108**, 1175–1204 (1957).
- ¹³A. Seko, A. Togo, and I. Tanaka, “Descriptors for machine learning of materials data,” in *Nanoinformatics*, edited by I. Tanaka (Springer Singapore, Singapore, 2018) pp. 3–23.
- ¹⁴A. Tversky, “Features of similarity,” *Psychol. Rev.* **84**, 327–352 (1977).
- ¹⁵R. L. Goldstone, D. L. Medin, and J. Halberstadt, “Similarity in context,” *Memory & Cognition* **25**, 237–255 (1997).
- ¹⁶G. Shafer, *A Mathematical Theory of Evidence* (Princeton University Press, 1976).
- ¹⁷T. Denceux, D. Dubois, and H. Prade, “Representations of uncertainty in artificial intelligence: Beyond probability and possibility,” in *A Guided Tour of Artificial Intelligence Research*, Vol. 1, edited by P. Marquis, O. Papini, and H. Prade (Springer Verlag, 2020) Chap. 4, pp. 119–150.
- ¹⁸A. P. Dempster, “Upper and lower probabilities induced by a multivalued mapping,” *The Annals of Mathematical Statistics* **38**, 325–339 (1967).
- ¹⁹N. Nu Thanh Ton, M.-Q. Ha, T. Ikenaga, A. Thakur, H.-C. Dam, and T. Taniike, “Solvent screening for efficient chemical exfoliation of graphite,” *2D Materials* **8**, 015019 (2020).
- ²⁰M.-Q. Ha, D.-N. Nguyen, V.-C. Nguyen, T. Nagata, T. Chikyow, H. Kino, T. Miyake, T. Denceux, V.-N. Huynh, and H.-C. Dam, “Evidence-based recommender system for high-entropy alloys,” *Nature Computational Science* **1**, 470–478 (2021).
- ²¹C. Williams and C. Rasmussen, “Gaussian processes for regression,” in *Advances in neural information processing systems* **8**, Max-Planck-Gesellschaft (MIT Press, Cambridge, MA, USA, 1996) pp. 514–520.
- ²²M. Lázaro-Gredilla, S. Van Vaerenbergh, and N. D. Lawrence, “Overlapping mixtures of gaussian processes for the data association problem,” *Pattern Recognition* **45**, 1386–1395 (2012).
- ²³V. Vovk, “Kernel ridge regression,” in *Empirical inference* (Springer, 2013) pp. 105–116.
- ²⁴L. Breiman, “Random forests,” *Machine Learning* **45**, 5–32 (2001).
- ²⁵A. Jain, J. Mao, and K. Mohiuddin, “Artificial neural networks: a tutorial,” *Computer* **29**, 31–44 (1996).
- ²⁶P. Smets, “Belief functions: The disjunctive rule of combination and the generalized bayesian theorem,” *International Journal of Approximate Reasoning* **9**, 1–35 (1993).
- ²⁷D. Hieu-Chi, “Datasets of binary and quaternary alloys with Curie temperature and magnetization for the eRSM,” (2023).
- ²⁸P. Villars, M. Berndt, K. Brandenburg, K. Cenzual, J. Daams, F. Hulliger, T. Massalski, H. Okamoto, K. Osaki, A. Prince, H. Putz, and S. Iwata, “The pauling file, binaries edition,” *Journal of Alloys and Compounds* **367** (2004), 10.1016/j.jallcom.2003.08.058.
- ²⁹Y. Xu, M. Yamazaki, and P. Villars, “Inorganic Materials Database for Exploring the Nature of Material,” *Jpn. J. Appl. Phys.* **50** (2011), 10.1143/JJAP.50.11RH02.
- ³⁰H. C. Dam, V. C. Nguyen, T. L. Pham, A. T. Nguyen, K. Terakura, T. Miyake, and H. Kino, “Important descriptors and descriptor groups of curie temperatures of rare-earth transition-metal binary alloys,” *Journal of the Physical Society of Japan* **87**, 113801 (2018).
- ³¹T. Fukushima, H. Akai, T. Chikyow, and H. Kino, “Automatic exhaustive calculations of large material space by korringa-kohn-rostoker coherent potential approximation method applied to equiatomic quaternary high entropy alloys,” *Phys. Rev. Materials* **6**, 023802 (2022).
- ³²Y.-Y. Ahn, J. P. Bagrow, and S. Lehmann, “Link communities reveal multiscale complexity in networks,” *Nature* **466**, 761–764 (2010).
- ³³B. G. Lindsay, “Mixture models: Theory, geometry and applications,” NSF-CBMS Regional Conference Series in Probability and Statistics **5**, i–163 (1995).
- ³⁴A. H. Murphy, “The finley affair: A signal event in the history of forecast verification,” *Weather and Forecasting* **11**, 3–20 (1996).

- ¹ ³⁵L. van der Maaten and G. Hinton, “Visualizing data using t-sne,”
² Journal of Machine Learning Research **9**, 2579–2605 (2008).
³ ³⁶T. L. Pham, H. Kino, K. Terakura, T. Miyake, and H. C. Dam,
⁴ “Novel mixture model for the representation of potential energy
⁵ surfaces,” The Journal of Chemical Physics **145**, 154103 (2016).
⁶ ³⁷H. Akai, M. Akai, S. Blügel, B. Drittler, H. Ebert, K. Terakura,
⁷ R. Zeller, and P. H. Dederichs, “Theory of Hyperfine Interactions
⁸ in Metals,” Progress of Theoretical Physics Supplement **101**, 11–
⁹ 77 (1990).