

A general learning scheme for classical and quantum Ising machines

Ludwig Schmid^{1*}, Enrico Zardini² and Davide Pastorello^{3,4}

¹ Chair for Design Automation — Technical University of Munich, 80333 Munich, Germany

² Department of Information Engineering and Computer Science — University of Trento, 38123 Trento, Italy

³ Department of Mathematics — University of Bologna, 40126 Bologna, Italy

⁴ TIFPA-INFN, 38123 Povo-Trento, Italy

* ludwig.s.schmid@tum.de

Abstract

An Ising machine is any hardware specifically designed for finding the ground state of the Ising model. Relevant examples are coherent Ising machines and quantum annealers. In this paper, we propose a new machine learning model that is based on the Ising structure and can be efficiently trained using gradient descent. We provide a mathematical characterization of the training process, which is based upon optimizing a loss function whose partial derivatives are not explicitly calculated but estimated by the Ising machine itself. Moreover, we present some experimental results on the training and execution of the proposed learning model. These results point out new possibilities offered by Ising machines for different learning tasks. In particular, in the quantum realm, the quantum resources are used for both the execution and the training of the model, providing a promising perspective in quantum machine learning.

Copyright attribution to authors.

This work is a submission to SciPost Physics Core.

License information to appear upon publication.

Publication information to appear upon publication.

Received Date

Accepted Date

Published Date

1

2 Contents

3	1 Introduction	2
4	2 Ising machines	4
5	3 The proposed model	6
6	3.1 Definition	6
7	3.2 Training process	7
8	3.3 Hidden spins	10
9	3.4 Computational cost	11
10	4 Empirical evaluation	12
11	4.1 Experimental setup	12
12	4.2 Random data	13
13	4.3 Function approximation	14

14	4.4 Bars and stripes	16
15	4.5 Choice of hyperparameters	18
16	5 Conclusion	18
17	References	19
18	<hr/>	
19		

20 1 Introduction

21 Machine learning models are algorithms that provide predictions about observed phenomena
 22 by extracting information from a set of collected data (the training set). In particular, *para-*
 23 *metric models* capture all relevant information within a finite set of parameters, with the set
 24 being independent of the number of training instances [1]. A celebrated example is repre-
 25 sented by *artificial neural networks* [2–4]. In the context of quantum computers, a common
 26 approach to machine learning is to employ variational quantum circuits, which can be trained
 27 by backpropagation as done with classical feedforward neural networks [5–8]. In addition
 28 to gate-based quantum computing, *quantum annealing* has also been considered to develop
 29 machine learning algorithms [9–11]. In any case, a crucial point in quantum machine learn-
 30 ing is the implementation of quantum procedures for model training as alternatives to classical
 31 methods. An example in this sense is the quantum support vector machine, trained by running
 32 the HHL quantum algorithm [12], which, however, presents the shortcoming of an impractical
 33 implementation on the currently available quantum devices. Therefore, a general challenge in
 34 quantum machine learning is to define learning schemes that can be efficiently implemented
 35 on quantum machines of the Noisy Intermediate-Scale Quantum (NISQ) era [13]. This is the
 36 motivation behind the present proposal of a learning model for quantum annealers in which
 37 the quantum resources are used both in the model execution and in the training process. The
 38 obtained theoretical and experimental results apply also to classical implementations of the
 39 model. Indeed, the key aspect of the training and execution of the proposed learning mech-
 40 anism is the computation of the ground state of the *Ising model*, which can, in principle, be
 41 solved using classical or quantum procedures.

42 An *Ising machine* can be considered a specific-purpose computer designed to return the
 43 absolute or approximate ground state of the Ising model. The latter is described by the energy
 44 function of a spin glass system under the action of an external field, namely,

$$E(\mathbf{z}) = \sum_{i=1}^N \theta_i z_i + \sum_{(i,j)} \Gamma_{ij} z_i z_j, \quad \text{with } \mathbf{z} \in \{-1, 1\}^N, \theta_i \in \mathbb{R}, \text{ and } \Gamma_{ij} \in \mathbb{R}, \quad (1)$$

45 where the sum $\sum_{(i,j)}$ is taken over the pairs of connected spins, counting each pair only once.
 46 The ground state is the spin configuration $\mathbf{z}^* \in \{-1, 1\}^N$ that minimizes the function (1).
 47 Therefore, in practice, an Ising machine solves a combinatorial optimization problem that can
 48 be represented as a quadratic unconstrained binary optimization (QUBO) problem, which is
 49 an NP-hard problem, by means of the change of variables $x_i = \frac{z_i+1}{2} \in \{0, 1\}$. In particular, an
 50 Ising machine can be an analog computer that evolves toward the Ising ground state due to a
 51 physical process like thermal or quantum annealing. Alternatively, it can also be implemented
 52 on a digital computer in terms of simulated annealing.

53 Ising machines are conceptually related to *Boltzmann machines* in the sense that they are
 54 both defined in terms of the Ising model, with couplings among spins and the action of an

55 external field. In the case of a Boltzmann machine, the coefficients θ and Γ of the energy
 56 function (1) are tuned so that, by sampling the spin configuration over the state of the system
 57 at thermal equilibrium (at a finite temperature T), a probability distribution resembling an
 58 input distribution defined on the training set [14] is generated [14].
 59 In detail, the output distribution of a Boltzmann machine is given by

$$p_T(\mathbf{z}) = Z^{-1} \exp\left[-\frac{E(\mathbf{z})}{k_B T}\right], \quad (2)$$

60 where $Z := \sum_{\mathbf{z}} \exp\left[-\frac{E(\mathbf{z})}{k_B T}\right]$ is the partition function and k_B is the Boltzmann constant. Usually,
 61 only a subset of spins is sampled, the so-called *visible nodes*, and the output distribution is given
 62 by the marginal distribution of (2). Instead, in the ideal case, the output of an Ising machine
 63 is deterministic and corresponds to the absolute minimum of (1). However, in a realistic
 64 scenario in which the Ising machine operates by thermal annealing, the output is probabilistic
 65 and distributed according to (2) with a value of T as low as possible.

66 The difference between Boltzmann and Ising machines lies in the fact that Boltzmann
 67 machines are parametric generative models. In contrast, Ising machines are considered as
 68 solvers of combinatorial optimization problems [15–17]. However, in this paper, we propose
 69 a supervised learning model for Ising machines whose training is inspired by the training of
 70 Boltzmann machines. A peculiar aspect of a Boltzmann machine is that it can be trained by
 71 gradient descent of a loss function \mathcal{L} depending on the weights θ and Γ , like the average neg-
 72 ative log-likelihood between the input distribution and the generated distribution, iteratively
 73 changing the parameters by a step in the opposite direction of the gradient. However, the par-
 74 tial derivatives of \mathcal{L} are not explicitly calculated but are estimated by sampling the network
 75 units. For instance, let us consider the update rule $\Gamma_{ij} \rightarrow \Gamma_{ij} + \delta\Gamma_{ij}$, which updates the coupling
 76 terms toward the minimum of the average negative log-likelihood. The update step ($\delta\Gamma_{ij}$) is
 77 given by [14]:

$$\delta\Gamma_{ij} = -\eta \left(\langle z_i z_j \rangle - \sum_{\mathbf{v}} p_{data}(\mathbf{v}) \langle z_i z_j \rangle_{\mathbf{v}} \right) \quad i, j = 1, \dots, N, \quad (3)$$

78 where $\eta > 0$ is the learning rate (user-specified), the sum is taken over the visible nodes \mathbf{v} ,
 79 p_{data} is the input distribution, $\langle \rangle$ is the Boltzmann average, and $\langle \rangle_{\mathbf{v}}$ is the Boltzmann average
 80 with clamped visible nodes. In other words, both the training and the execution of a Boltz-
 81 mann machine are performed by sampling the units of the network at thermal equilibrium. A
 82 quantum version of the Boltzmann machine has also been proposed [18], and the simulations
 83 have shown that the presence of a *transverse field Hamiltonian* improves the training process
 84 with respect to the classical model, generating distributions that are closer to the input one in
 85 terms of the Kullback-Liebler divergence.

86 This paper adopts a similar viewpoint for training an Ising machine. After defining a para-
 87 metric predictive model based on the ground state of the Ising model, we prove that it can be
 88 trained by gradient descent of a mean squared error loss function, executing the model itself
 89 to obtain the gradient estimates. In particular, the structure of the model does not require
 90 that the Ising machine returns the true ground state with infinite precision, and a suboptimal
 91 output works for training and executing the predictive model. In addition, our results apply to
 92 both classical and quantum machines. However, in the second case, the impact may be more
 93 significant since the quantum annealing resources are also exploited for the training process.
 94 In this sense, the purpose is similar to that of the *parameter-shift rule*, which is used in gate-
 95 based quantum computing to train a parametric quantum circuit without explicitly calculating
 96 the partial derivatives [19].

97 The paper is structured as follows: in Section 2, we introduce generalities and elemen-
 98 tary notions about the Ising model and Ising machines, with a particular focus on quantum

99 annealing; Section 3 deals with the proposed parametric learning model, to be executed by
 100 an Ising machine, and the main theoretical result of the paper, i.e., the proof that the model
 101 can be trained by running the Ising machine itself; in Section 4, an empirical evaluation of the
 102 proposed machine learning method is provided; in Section 5, we discuss the perspectives of
 103 the proposal, and we draw our conclusions on the proposed parametric model.

104 2 Ising machines

105 This section introduces the formal definition of the Ising model and the concept of using specific
 106 Ising machines to solve the corresponding groundstate problem. Afterward, we briefly describe
 107 the two Ising machines employed in this work, namely simulated and quantum annealing.

108 The *Ising model* is a mathematical description extensively utilized in the study of ferromag-
 109 netism. Renowned for its versatility and simplicity, it stands as a fundamental paradigm in the
 110 domain of statistical mechanics [20]. In its general formulation, the Ising model is defined on
 111 a graph (V, E) , wherein each vertex represents a discrete variable $z_i \in \{-1, 1\}$. These vari-
 112 ables correspond to *spins*, with associated *biases* $\theta_i \in \mathbb{R}$ denoting the inclination of each spin
 113 toward one of the two available values. Furthermore, the weighted edges $\Gamma_{ij} \in \mathbb{R}$ connecting
 114 two spins i and j define the coupling dynamics between the spins, indicating their preference
 115 to align or oppose each other in value. This graph structure is illustrated in Figure 1. The total
 116 energy of a spin configuration $\mathbf{z} \in \{-1, 1\}^{|V|}$ is expressed as

$$E(\boldsymbol{\theta}, \boldsymbol{\Gamma}, \mathbf{z}) = \sum_{i=1}^{|V|} \theta_i z_i + \sum_{(i,j) \in E} \Gamma_{ij} z_i z_j = \boldsymbol{\theta} \mathbf{z} + \mathbf{z}^T \boldsymbol{\Gamma} \mathbf{z}, \quad (4)$$

117 where the *biases* $\theta_1, \dots, \theta_{|V|} \in \mathbb{R}$ and the *couplings* $\Gamma_{ij} \in \mathbb{R} \forall (i, j) \in E$ are conveniently con-
 118 solidated into the vector $\boldsymbol{\theta}$ and the matrix $\boldsymbol{\Gamma}$ (with $\Gamma_{ij} = 0$ when $(i, j) \notin E$), respectively.
 119 Realistically, the values of the parameters are bounded. Hence, it is possible to assume that
 120 biases and couplings take values into compact intervals of \mathbb{R} . Within the realm of statisti-
 121 cal physics, these quantities are typically referred to as the external magnetic field strength
 122 and spin interactions due to their fundamental roles in the physical manifestation of the Ising
 123 model.

124 An *Ising machine* can be defined as a non-von Neumann computer for solving combinato-
 125 rial optimization problems [21]. More precisely, its input is represented by the energy function
 126 of the Ising model (4), with biases and coupling terms properly initialized. The machine effec-
 127 tively operates by minimizing the energy function and providing the optimal spin configuration
 128 \mathbf{z}^* as the output. Actually, the quest to determine the ground state of an Ising model is of sig-
 129 nificant importance, as any problem within the NP complexity class can be formulated as an
 130 Ising problem with only a polynomial increase in complexity [22]. An elementary and abstract
 131 definition of an Ising machine, motivated by the general approach adopted in this paper, is the
 132 following:

133 **Definition 1.** Given the energy function defined in (4), an **(abstract) Ising machine** is any map
 134 $(\boldsymbol{\theta}, \boldsymbol{\Gamma}) \mapsto \mathbf{z}^* := \operatorname{argmin}_{\mathbf{z}} E(\boldsymbol{\theta}, \boldsymbol{\Gamma}, \mathbf{z})$.

135 Additionally, we can also consider the minimum value of the energy $E_0(\boldsymbol{\theta}, \boldsymbol{\Gamma}) := E(\boldsymbol{\theta}, \boldsymbol{\Gamma}, \mathbf{z}^*)$
 136 as the output of an Ising machine. This ground state energy of the Ising model is obtained by
 137 substituting the spin configuration $\mathbf{z}^* = \operatorname{argmin}_{\mathbf{z}} E(\boldsymbol{\theta}, \boldsymbol{\Gamma}, \mathbf{z})$ into (4). In this context, the Ising
 138 machine consistently yields a numerical result with a negative sign. An illustration of an Ising
 139 machine that finds the ground state of a small Ising model is shown in Figure 1.

140 Relevant examples of Ising machines as specific-purpose hardware devices are quantum
 141 annealers [23] or coherent Ising machines with optical processors [24–27]. However, an Ising

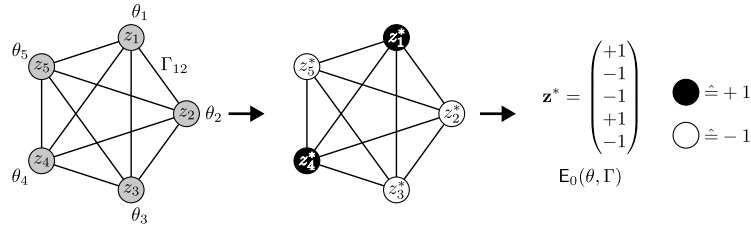


Figure 1: Ising model and Ising machine: On the left, an illustration of the graph structure of an Ising model characterized by a fully connected graph, with $|V| = 5$ spins \mathbf{z} , corresponding biases θ , and couplings Γ . An Ising machine maps the Ising model to the right-hand side of the figure, returning a $\{-1, +1\}$ assignment (illustrated as white/black nodes) to each binary variable z_i . The output is the spin configuration \mathbf{z}^* and the corresponding minimal energy $E_0(\theta, \Gamma)$.

142 machine can also be simulated on a classical digital computer. In this respect, simulated an-
 143 nealing is a standard approach and addresses the Ising model as a combinatorial optimization
 144 problem. In more detail, simulated annealing is a probabilistic metaheuristic inspired by the
 145 analogical notion of controlling the cooling process observed in physical materials [28]. The
 146 algorithm employs stochastic acceptance criteria, resembling a Boltzmann probability, to navi-
 147 gate the solution space and escape local optima. Over time, usually indicated by a temperature
 148 parameter T that mimics the cooling process, less favorable moves are increasingly rejected.
 149 In practice, simulated annealing employs random search and local exploration to converge
 150 toward near-optimal or optimal solutions. However, although the algorithm is easy to imple-
 151 ment and robust from a theoretical point of view, it may present a slow convergence rate [29].
 152 A promising alternative path is the development of analog platforms like coherent Ising ma-
 153 chines. They represent optical parametric oscillator (OPO) networks in which the collective
 154 mode of oscillation beyond a certain threshold corresponds to an optimal solution for a given
 155 large-scale Ising model [24–27]. The learning scheme proposed here is agnostic and can be
 156 implemented on this kind of Ising machines. Nevertheless, in the experimental part we have
 157 considered only simulated and quantum annealing.

158 Quantum annealing is a type of heuristic search used to solve optimization problems [23,
 159 30–32]. The procedure is implemented by the time evolution of a quantum system toward
 160 the ground state of a *problem Hamiltonian*. More precisely, let us consider the time-dependent
 161 Hamiltonian

$$H(t) = \gamma(t)H_D + H_P \quad t \geq 0, \quad (5)$$

162 where H_P is the problem Hamiltonian, H_D is the *transverse field Hamiltonian*, and $\gamma : \mathbb{R}^+ \rightarrow \mathbb{R}$
 163 is a decreasing function. Roughly speaking, H_D gives the kinetic term inducing the exploration
 164 of the solution landscape by means of quantum fluctuations, and γ attenuates the kinetic term
 165 driving the system toward the ground state of H_P . Quantum annealing can be physically re-
 166 alized by considering a network of qubits arranged on the vertices of a graph (V, E) , with
 167 $|V| = n$ and whose edges E represent the couplings among the qubits. In detail, the prob-
 168 lem Hamiltonian is defined as the following self-adjoint operator on the n -qubit Hilbert space
 169 $\mathbb{H} = (\mathbb{C}^2)^{\otimes n}$:

$$H_P = \sum_{i \in V} \theta_i \sigma_z^{(i)} + \sum_{(i,j) \in E} \Gamma_{ij} \sigma_z^{(i)} \sigma_z^{(j)}, \quad (6)$$

170 with real coefficients θ_i, Γ_{ij} , which are identified again as biases and couplings due to their
 171 similar role in the Ising model. In the computational basis, the $2^n \times 2^n$ matrix $\sigma_z^{(i)}$ acts locally
 172 as the Pauli matrix

$$\sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \quad (7)$$

173 on the i -th tensor factor and as the 2×2 identity matrix on the other tensor factors. In fact,
 174 the eigenvectors of H_P form the computational basis of \mathbb{H} , and the corresponding eigenvalues
 175 are the values of the classical energy function (4). On the other hand, for the transverse field
 176 Hamiltonian a typical form is

$$H_D = \sum_{i \in V} \theta_i \sigma_x^{(i)}, \quad (8)$$

where the local operator $\sigma_x^{(i)}$ is defined in a similar way to $\sigma_z^{(i)}$ in terms of the Pauli matrix

$$\sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

177 H_D does not commute with H_P and provides the unbiased superposition of all the conceivable
 178 solutions as the system initial state. Eventually, it is worth highlighting that quantum annealing
 179 is related to *adiabatic quantum computing* (AQC) as the solution of a given problem can be
 180 encoded into the ground state of a problem Hamiltonian. However, the two notions do not
 181 coincide. Indeed, in quantum annealing, the quantum system is not assumed to be isolated;
 182 therefore, it can be characterized by a non-unitary evolution. Another difference is that, in
 183 quantum annealing, the entire computation is not required to take place in the instantaneous
 184 ground state of the time-varying Hamiltonian like in AQC [32].

185 3 The proposed model

186 This section formally introduces the proposed parametric model, followed by an in-depth dis-
 187 cussion on the training using gradient descent and the estimation of the relevant partial deriva-
 188 tives of a quadratic loss function. The final part presents [some practical considerations required](#)
 189 [to operate and train the model in real-world scenarios](#) and [a discussion of its computational](#)
 190 [cost](#).

191 3.1 Definition

192 In the context of supervised learning, the goal of an algorithm is to approximate a function
 193 $f : X \rightarrow Y$ given a training set $\{(x_1, f(x_1)), \dots, (x_N, f(x_N))\}$, which is a collection of ele-
 194 ments in the set X with the corresponding values of f . An approximation of f can be obtained
 195 through a parametric function after an optimal choice of its parameters, generalizing the in-
 196 formation encoded into the training set. In fact, the notion of a parametric model is closely
 197 related to the existence of a parametric function that can be used to approximate the target
 198 function.

Definition 2. Let X and Y be non-empty sets respectively called **input domain** and **output domain**. A (deterministic) **parametric model** is a function

$$x \mapsto y = F(x|\Gamma) \quad x \in X, y \in Y,$$

199 with Γ being a set of real parameters.

In practice, given a training set of input-output pairs, the task consists in finding the parameters Γ such that the model assigns the correct or approximately correct output, with high probability, to any previously unseen input. The parameters are typically determined by optimizing a *loss function* such as

$$\mathcal{L}(\Gamma) = \frac{1}{N} \sum_{i=1}^N d(y_i, F(x_i|\Gamma)),$$

200 where d is a metric defined over Y , and the procedure is commonly referred to as *training*.

201 A preliminary depiction of the general problem considered in this paper is the following:
 202 given a real-valued function $f : X \rightarrow \mathbb{R}$, with $X \subset \mathbb{R}^n$ and $n \in \mathbb{N}$, the objective consists
 203 in training a predictive model F that approximates the original function f within the super-
 204 vised learning framework. This function approximation task encompasses a wide range of
 205 conventional machine learning endeavors such as regression and classification. In particular,
 206 the proposed parametric model is defined over the concept of Ising machines as introduced
 207 in Section 2. The input information is encoded into the biases θ of an Ising model, while the
 208 adjustable parameters are represented by the couplings Γ of (4). The Ising machine is then
 209 used to find the ground state of the Ising model, and the corresponding ground state energy
 210 is used as the model output. Note that the ground state energy invariably assumes a negative
 211 value, and the magnitude of the input biases significantly influences its absolute magnitude.
 212 To account for this, we introduce an ancillary scaling factor denoted as λ and an energy offset
 213 indicated as ϵ . This yields the subsequent formulation of the model.

214 Given an Ising machine, an input vector $\theta = (\theta_1, \dots, \theta_n) \in X \subset \mathbb{R}^n$, and the parameters
 215 $\{\Gamma_{ij}\}$ with $i, j = 1 \dots n$ (the nonzero Γ_{ij} are specified by the topology graph of the machine),
 216 one can define a parametric model F based on the ground state energy of an Ising model as

$$\begin{aligned} F(\theta|\Gamma, \lambda, \epsilon) &:= \lambda \min_{z \in \{-1, 1\}^n} E(\theta, \Gamma, z) + \epsilon \\ &= \lambda E_0(\theta, \Gamma) + \epsilon, \end{aligned} \quad (9)$$

217 where $\lambda \in \mathbb{R}$ and $\epsilon \in \mathbb{R}$ are additional tunable parameters that do not influence the Ising
 218 model energy. The model definition reveals a general neural approach in the sense that data
 219 are represented by the biases of the spins, which can be associated with neurons, and the
 220 parameters are the weights attached to the connections between spins (neurons). It is worth
 221 noting that, for the model execution, there is no requirement that the Ising machine returns the
 222 true ground state. More precisely, the fact that an approximated ground state does not match
 223 the exact solution of the combinatorial problem underlying the minimization is not a severe
 224 drawback for the learning process. Indeed, assuming that the deviation of the energy output
 225 from E_0 is systematic (e.g. due to the finite precision of the Ising machine), this deviation
 226 becomes a characteristic of the model itself, and the training procedure accordingly provides
 227 optimized parameters. Despite its simplicity, the model presents interesting training properties
 228 that we mathematically characterize in the next section.

229 3.2 Training process

230 Training the proposed parametric model for the approximation of a real-valued function entails
 231 minimizing the empirical risk across a provided dataset, denoted as \mathcal{D} , encompassing input-
 232 output pairs derived from the original function. To this aim, we employ the conventional
 233 approach of optimizing the model parameters to minimize the mean squared error (MSE)
 234 between the predicted output and the actual data values.

235 Given the training set $\mathcal{D} = \{(\theta^{(a)}, y^{(a)})\}_{a=1, \dots, N}$, with $y^{(a)} = f(\theta^{(a)})$, where $f : X \rightarrow \mathbb{R}$,
 236 with $X \subset \mathbb{R}^n$, is an unknown function to approximate, the model (9) can be trained by mini-
 237 mizing the MSE loss function

$$\mathcal{L}(\Gamma, \lambda, \epsilon) = \frac{1}{N} \sum_{a=1}^N [F(\theta^{(a)}|\Gamma, \lambda, \epsilon) - y^{(a)}]^2. \quad (10)$$

238 Our objective is to address this minimization task employing a gradient descent approach,
 239 iteratively updating the parameters Γ , λ , and ϵ by taking steps in the direction opposite to the

240 gradient of the loss function \mathcal{L} :

$$\delta\Gamma = -\eta\nabla_{\Gamma}\mathcal{L}, \quad \delta\lambda = -\eta\frac{\partial\mathcal{L}}{\partial\lambda}, \quad \delta\epsilon = -\eta\frac{\partial\mathcal{L}}{\partial\epsilon}, \quad (11)$$

241 where $\eta > \mathbf{0}$ is the learning rate, which controls the optimization step size. Let us remark
242 that each parameter is assumed to take values into a compact interval in \mathbb{R} ; consequently, the
243 parameter space is a hyperrectangle. On one hand, the partial derivatives of \mathcal{L} with respect to
244 λ and ϵ are well-defined and trivial to calculate. On the other hand, the following theorem,
245 which provides the update rules for the optimization of \mathcal{L} by gradient descent, implies that
246 the gradient $\nabla_{\Gamma}\mathcal{L}$ is defined almost everywhere in the parameter hyperrectangle.

247 **Theorem 3.** Let F be the parametric model defined in (9), $\mathcal{D} = \{(\theta^{(a)}, y^{(a)})\}_{a=1, \dots, N}$ be a train-
248 ing set for F , \mathcal{L} be the MSE loss function defined in (10), and $\eta > \mathbf{0}$ be the learning rate. Then,
249 the partial derivatives of F with respect to the couplings Γ are defined almost everywhere in the
250 parameter space, and the update rules for Γ , λ , ϵ for the gradient descent of \mathcal{L} are:

$$\Gamma_{ij}^{(k+1)} = \Gamma_{ij}^{(k)} - \eta \frac{2\lambda^{(k)}}{N} \sum_{a=1}^N [\lambda^{(k)} E_0(\theta^{(a)}, \Gamma^{(k)}) + \epsilon^{(k)} - y^{(a)}] z_i^* z_j^*, \quad (12)$$

251

$$\lambda^{(k+1)} = \lambda^{(k)} - \eta \frac{2}{N} \sum_{a=1}^N [\lambda^{(k)} E_0(\theta^{(a)}, \Gamma^{(k)}) + \epsilon^{(k)} - y^{(a)}] \left[\sum_{i=1}^n \theta_i^{(a)} z_i^* + \sum_{(i,j) \in E} \Gamma_{ij}^{(k)} z_i^* z_j^* \right], \quad (13)$$

252

$$\epsilon^{(k+1)} = \epsilon^{(k)} - \eta \frac{2}{N} \sum_{a=1}^N [\lambda^{(k)} E_0(\theta^{(a)}, \Gamma^{(k)}) + \epsilon^{(k)} - y^{(a)}], \quad (14)$$

253 where $\Gamma^{(k)}$, $\lambda^{(k)}$, $\epsilon^{(k)}$ are the values of the parameters within the k -th iteration of the gradient
254 descent, and $\mathbf{z}^* = \operatorname{argmin}_{\mathbf{z}} E(\theta^{(a)}, \Gamma^{(k)}, \mathbf{z})$.

255 *Proof.* By direct calculation, the partial derivative of F with respect to Γ_{ij} is

$$\frac{\partial F(\theta|\Gamma, \lambda, \epsilon)}{\partial \Gamma_{ij}} = \lambda \frac{\partial}{\partial \Gamma_{ij}} \left(\sum_{i=1}^n \theta_i z_i^* + \sum_{(i,j)} \Gamma_{ij} z_i^* z_j^* \right) = \lambda z_i^* z_j^*, \quad (15)$$

256 where z_i^* and z_j^* are the i -th and j -th components of $\mathbf{z}^* = \operatorname{argmin}_{\mathbf{z}} E(\theta, \Gamma, \mathbf{z})$, respectively.
257 Since the optimal spin configuration \mathbf{z}^* also depends on Γ (and θ), we should consider the
258 derivatives $\partial z_l^* / \partial \Gamma_{ij}$ for $l = 1, \dots, n$ in the final step outlined in (15). However, it must be
259 noted that the function $z_l^* = z_l^*(\theta, \Gamma)$ is piecewise constant. Hence, its derivative is zero almost
260 everywhere in its domain, and the remaining points, corresponding to spin flips of z_l^* , turn out
261 to be points of non-differentiability of $z_l^*(\theta, \Gamma)$. Substituting (15) into (11), we obtain the
262 following update step ($\delta\Gamma_{ij}$) for the MSE loss function (10):

$$\delta\Gamma_{ij} = -\eta \frac{\partial\mathcal{L}}{\partial\Gamma_{ij}} = -\eta \frac{2}{N} \sum_{a=1}^N [F(\theta^{(a)}|\Gamma, \lambda, \epsilon) - y^{(a)}] \frac{\partial F}{\partial\Gamma_{ij}} \quad (16)$$

$$\begin{aligned} &= -\eta \frac{2\lambda}{N} \sum_{a=1}^N [F(\theta^{(a)}|\Gamma, \lambda, \epsilon) - y^{(a)}] z_i^* z_j^* \\ &= -\eta \frac{2\lambda}{N} \sum_{a=1}^N [\lambda E_0(\theta^{(a)}, \Gamma) + \epsilon - y^{(a)}] z_i^* z_j^*. \end{aligned} \quad (17)$$

Algorithm 1: Model training

Input: dataset $\mathcal{D} = \{(\boldsymbol{\theta}^{(a)}, y^{(a)})\}_{a=1, \dots, N}$, learning rate η , optimization steps N_{epochs}
Output: trained model $F_{\text{model}}(\boldsymbol{\theta})$

- 1 Initialize the parameters Γ ;
- 2 **for** step k in N_{epochs} **do**
- 3 **for** $(\boldsymbol{\theta}^{(a)}, y^{(a)})$ in \mathcal{D} **do**
- 4 run the Ising machine to obtain $E_0(\boldsymbol{\theta}^{(a)}, \Gamma^{(k)})$ and \mathbf{z}^* ;
- 5 **end**
- 6 update $\Gamma^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\epsilon}^{(k)}$ according to (12) - (13) - (14);
- 7 **end**
- 8 **return** $F_{\text{model}}(\boldsymbol{\theta}) = F(\boldsymbol{\theta} | \Gamma^{N_{\text{epochs}}}, \boldsymbol{\lambda}^{N_{\text{epochs}}}, \boldsymbol{\epsilon}^{N_{\text{epochs}}})$;

263 Therefore, the parameter update rule for the $(k + 1)$ -th iteration turns out to be

$$\Gamma_{ij}^{(k+1)} = \Gamma_{ij}^{(k)} - \eta \frac{2\lambda^{(k)}}{N} \sum_{a=1}^N [\lambda^{(k)} E_0(\boldsymbol{\theta}^{(a)}, \Gamma^{(k)}) + \epsilon^{(k)} - y^{(a)}] \mathbf{z}_i^* \mathbf{z}_j^*, \quad (18)$$

264 wherein we have omitted the explicit dependence of \mathbf{z}_i^* and \mathbf{z}_j^* on \mathbf{a} and \mathbf{k} for the sake of
 265 brevity of notation. The update rules for $\boldsymbol{\lambda}$ and $\boldsymbol{\epsilon}$ can be derived analogously. Specifically, the
 266 partial derivatives of F with respect to $\boldsymbol{\lambda}$ and $\boldsymbol{\epsilon}$ are

$$\frac{\partial F(\boldsymbol{\theta} | \Gamma, \boldsymbol{\lambda}, \boldsymbol{\epsilon})}{\partial \boldsymbol{\lambda}} = \sum_{i=1}^n \theta_i \mathbf{z}_i^* + \sum_{(i,j)} \Gamma_{ij} \mathbf{z}_i^* \mathbf{z}_j^*, \quad \frac{\partial F(\boldsymbol{\theta} | \Gamma, \boldsymbol{\lambda}, \boldsymbol{\epsilon})}{\partial \boldsymbol{\epsilon}} = 1. \quad (19)$$

267 Then, the claims (13) and (14) follow. □

268 In this way, the model parameters can be optimized for a certain number of steps N_{epochs} . The
 269 complete training process is described as pseudocode in Algorithm 1 and illustrated as a flow
 270 diagram in Figure 2. In particular, for each training step \mathbf{k} , the model is evaluated on each
 271 $(\boldsymbol{\theta}^{(a)}, y^{(a)})$ pair in the training set \mathcal{D} and the parameters are updated according to Theorem 3.
 272 The trained model is defined by the final iteration as

$$F_{\text{model}}(\boldsymbol{\theta}) = F(\boldsymbol{\theta} | \Gamma^{N_{\text{epochs}}}, \boldsymbol{\lambda}^{N_{\text{epochs}}}, \boldsymbol{\epsilon}^{N_{\text{epochs}}}). \quad (20)$$

273 Therefore, the training process bears similarities to that of a neural network but with a note-
 274 worthy distinction. Indeed, in our model, the conventional backpropagation step for calcu-
 275 lating the partial derivatives is replaced by the Ising machine computation of E_0 and \mathbf{z}^* . In
 276 particular, we propose the usage of quantum annealing as a well-suited Ising machine, which
 277 serves a dual purpose: executing the model according to (9) and facilitating the model training
 278 through the iterative assessment of the loss function gradient. In detail, the spin configuration
 279 \mathbf{z}^* , retrieved from the annealer and representing the ground state of the qubit network, can
 280 be used to compute the parameter adjustments according to (12), (13) and (14). Instead, the
 281 corresponding energy value is used to compute the model prediction.

282 This ability to utilize the output of the Ising machine to train and evaluate the model
 283 constitutes the major distinction to other Ising machine-based models [33, 34] that require an
 284 explicit calculation of the corresponding derivatives to update the model parameters.

285 A model trained in this manner possesses the capability to predict inputs beyond those
 286 present in \mathcal{D} . Analogously to other machine learning models, this rests upon the expectation
 287 that, if the model is trained on an extensive dataset, it can assimilate and generalize from those
 288 examples, ultimately serving as an approximation of the original function within a certain
 289 value range. Moreover, although the Ising energy (4) depends only linearly on the input

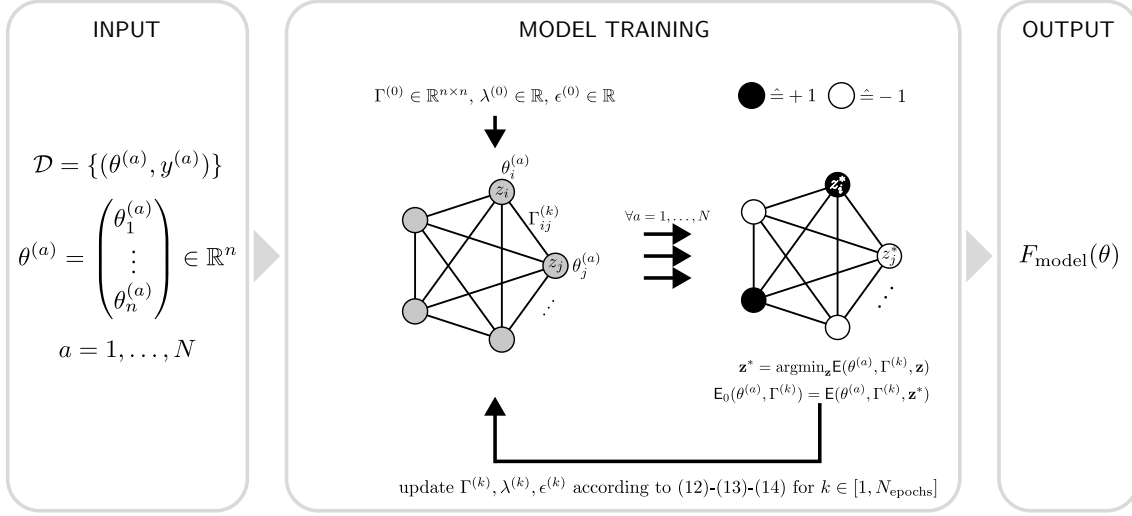


Figure 2: Model training: Illustration of the training process for the proposed model. In particular, given a dataset $\mathcal{D} = \{(\theta^{(a)}, y^{(a)})\}_{a=1, \dots, N}$, an Ising model is instantiated for each sample by setting the biases to $\theta^{(a)}$ and using the couplings Γ as free parameters. Then, for each model, an Ising machine is run in order to obtain the spin configuration \mathbf{z}^* and the corresponding model minimal energy E_0 . Finally, the collected values are used to update the couplings Γ and the two additional parameters λ and ϵ according to the rules presented in Theorem 3. This procedure is repeated N_{epochs} times until the trained model $F_{\text{model}}(\theta) = F(\theta | \Gamma^{N_{\text{epochs}}}, \lambda^{N_{\text{epochs}}}, \epsilon^{N_{\text{epochs}}})$ is returned.

290 vector θ , determining the minimum energy entails a complex interplay between the input
 291 and the model parameters Γ . Consequently, an open theoretical question regarding the class
 292 of functions that can be approximated through the proposed methodology arises. In other
 293 words, given an Ising model, what is its expressibility in terms of ground state energies by
 294 varying only the qubit couplings? From a practical perspective, the limitations of the quantum
 295 annealer architecture (number of qubits, topology connectivity, value bounds for θ and Γ)
 296 impose additional obvious constraints.

297 3.3 Hidden spins

298 In the proposed model, assuming a complete topology graph, the number of tunable paramete-
 299 rs Γ_{ij} scales quadratically with respect to the input dimension n . In practice, the number of
 300 model parameters is intrinsically fixed by the input dimensionality, akin to a neural network
 301 featuring only input and output layers. In the neural network scenario, to enhance the model
 302 expressiveness, the number of parameters is typically augmented by introducing additional
 303 hidden layers. In a similar way, we consider additional *hidden spins*, represented by addi-
 304 tional nodes in the topology graph. These additional spins increase the number of couplings
 305 and, therefore, the number of parameters of the model. This is accomplished by adding a
 306 preprocessing step,

$$h_{\text{pre}} : \mathbb{R}^n \rightarrow \mathbb{R}^{n_{\text{total}}}, \quad (21)$$

307 mapping the original input vector θ from the feature space \mathbb{R}^n to a higher-dimensional space
 308 characterized by $n_{\text{total}} = n + n_{\text{hidden}}$ dimensions, with n_{hidden} representing the number of
 309 additional hidden spins. An illustration of this preprocessing step and the increase in the
 310 number of coupling parameters is given in Figure 3.

311 The preprocessing step does not affect the training process. Indeed, the model can still
 312 be trained as described in Section 3.2. Instead, the choice of the preprocessing function ex-

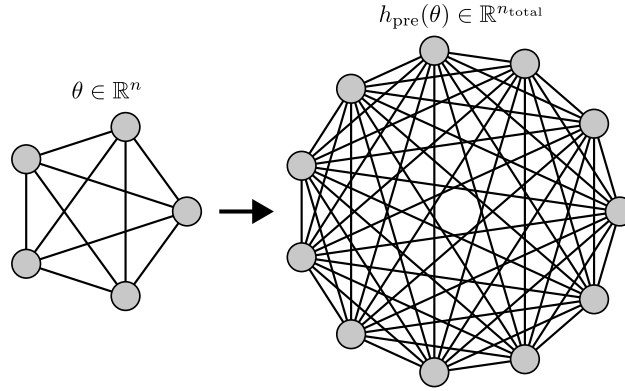


Figure 3: Hidden spins: Two exemplary Ising models with full connectivity. This comparison shows the increase in trainable coupling parameters (graph edges) when the original input θ is mapped to a higher dimensional space using a preprocessing step h_{pre} .

313 erts a significant influence on the model’s performance. For instance, let us consider a trivial
 314 preprocessing procedure that appends zero values to the input vector in order to reach the
 315 desired dimension. Although this approach would increase the number of model parameters,
 316 the hidden spins would be indistinguishable from each other, resulting in a very similar learn-
 317 ing behavior and making them redundant. In contrast, initializing the additional dimensions
 318 with random values would mitigate this issue, but these values may overshadow the original
 319 input, especially if $n_{\text{hidden}} \gg n$. In this work, we propose and evaluate a first simple scheme
 320 to initialize additional spins based on a constant real-valued offset. This *offset initialization*
 321 approach is defined as

$$\theta \in \mathbb{R}^n \rightarrow h_{\text{offset}}(\theta) = \begin{pmatrix} \theta \\ \theta + 1 \cdot d \\ \vdots \\ \theta + (l-1) \cdot d \end{pmatrix} \in \mathbb{R}^{n_{\text{total}}}, \quad (22)$$

322 where $d \in \mathbb{R}^n$, $l \in \mathbb{Z}^+$, and $n_{\text{total}} = ln$ (i.e., n_{total} is a multiple of n). This corresponds to a
 323 repeated concatenation of the original input θ with an increasing real-valued offset d .

324 3.4 Computational cost

325 To define the computational cost of the proposed machine-learning model, three aspects
 326 must be taken into consideration: the initialization of the model, including the embedding
 327 into the available Ising machine; the training of the model, with repeated calls to the Ising
 328 machine and the update of the parameters; the evaluation on the provided input. In the
 329 following, the space and time complexity of the proposed model are discussed.

330 The encoding of the problem requires n_{total} spin variables, considering the possibility of
 331 additional hidden spins (see Section 3.3). However, to achieve an all-to-all connectivity on
 332 sparse topology graphs like those of the D-Wave machines, an additional quadratic overhead
 333 has to be paid [35], resulting in a total space complexity of $\mathcal{O}(n_{\text{total}}^2)$.

334 The time complexity of the initialization (including the embedding) is $\mathcal{O}(n_{\text{total}}^2)$ [35].
 335 Instead, the training of the model requires N_{epochs} optimization steps. Specifically, in each
 336 optimization step, the Ising machine is evaluated on N samples and the $\mathcal{O}(n_{\text{total}}^2)$ model
 337 parameters are updated. To provide its output, the Ising machine has to solve the spin
 338 glass system of the model (Equation 1), which is an NP-hard problem in general [36, 37].
 339 Even for quantum annealers, the time required to find an exact solution is expected to be

340 inversely proportional to the minimum energy gap of the ground state [23], resulting in
 341 an exponential worst-case complexity [36]. Nevertheless, an approximate solution can be
 342 found by leveraging the non-adiabatic system evolution [23]. Therefore, we can assume
 343 the Ising machine runtime to be upper bounded by some τ throughout the whole training
 344 process. Given the Ising machine outputs, the parameters update can be done in time
 345 $\mathcal{O}(n_{\text{total}}^2 N)$ using Theorem 3. Overall, the time complexity of the model turns out to be
 346 $\mathcal{O}(n_{\text{total}}^2 + N_{\text{epochs}} N(\tau + n_{\text{total}}^2) + \tau)$, where the last term corresponds to the evaluation of
 347 the model using a single Ising machine execution.

348 4 Empirical evaluation

349 This section provides an initial proof of concept of the model’s capabilities. Indeed, this is
 350 neither a benchmarking exercise nor an in-depth analysis of the model’s expressiveness but a
 351 demonstration of possible use cases and applications of the model. A detailed performance
 352 evaluation of the model, entailing the necessary statistical repetitions and the comparison to
 353 alternative models, is left for future work. To simplify the usage of the model, a Python package
 354 that automates the repeated calls to the Ising machines during the training of the model and
 355 also facilitates the cross-usage with other common Python machine learning packages (such
 356 as PyTorch) was published on Github [38]. As a first experiment, the model has been trained
 357 on randomly sampled datasets to demonstrate the trainability of the model itself according
 358 to the update rules of Theorem 3. Then, as real-world demonstrations, the model has been
 359 trained for the function approximation task and also as a binary classifier for the bars and
 360 stripes dataset.

361 4.1 Experimental setup

362 As discussed in Section 3, the model supports different Ising machines. In this work, we have
 363 considered simulated annealing and quantum annealing, both provided by the D-Wave Ocean
 364 Software SDK [39]. While the former represents a software implementation of simulated
 365 annealing, the latter directly accesses the superconducting annealing hardware supplied by D-
 366 Wave. In particular, the *Advantage_system5.4* has been used here. More in detail, the quantum
 367 annealing hardware in question is characterized by 5760 qubits and is based on the Pegasus
 368 topology, with an inter-qubit connectivity of 15. To control the hardware, D-Wave provides
 369 the Ocean SDK, which includes multiple software packages facilitating the handling of the an-
 370 nealing hardware. Among them, it is worth mentioning the *minorminer* package, which has
 371 been used to embed the problems into the annealer topology. In practice, to achieve the de-
 372 sired connectivity (all-to-all in this case), multiple physical qubits are chained together to form
 373 logical qubits; the drawback lies in the reduced number of available qubits. In particular, in
 374 each run, the embedding has been computed once for a fully connected graph of the required
 375 size and reused in the subsequent calls to the annealer; for this aim, the *FixedEmbeddingCom-
 376 posite* class of the Ocean SDK has been employed. Regarding the actual annealing process, the
 377 default setup has been used, namely, automatic rescaling of bias and coupling terms to fit the
 378 available hardware ranges, chain strength settings according to *uniform_torque_compensation*,
 379 an annealing time of $20\mu\text{s}$, and a twelve-point annealing schedule. To account for the high
 380 number of calls to the annealing hardware throughout training and save hardware access time,
 381 a number of reads (sampling shots) equal to 1 has been used for each annealing process. For
 382 more information, refer to Zenodo [40], where the set of notebooks used [have has](#) been made
 383 available.

384 Concerning the model parameters, in all experiments, the couplings $\Gamma_{ij}^{(0)}$ have been initial-

385 ized to zero and updated according to (12). Instead, λ and ϵ have been kept fixed through-
 386 out the training process and considered as hyperparameters to facilitate the learning process.
 387 Specifically, the selection of the λ value has been done manually to ensure that the model
 388 output was reasonably well-aligned with the range of values of the training data. By contrast,
 389 the ϵ value has been set according to the outcomes of a the first round of sampling. In detail,
 390 the following rule has been used:

$$\epsilon = \frac{1}{N} \sum_{a=1}^N [y^{(a)} - F(\theta^{(a)} | \Gamma^{(0)}, \lambda, \mathbf{0})] = \frac{1}{N} \sum_{a=1}^N \left[y^{(a)} + \lambda \sum_{i=1}^n |\theta_i^{(a)}| \right], \quad (23)$$

391 with the last equivalence being valid only if $\Gamma_{ij}^{(0)} = \mathbf{0}$ for $i, j \in \{1, \dots, n\}$.

392 4.2 Random data

393 To demonstrate the trainability of the model, 30 distinct datasets, each comprising $N = 20$
 394 data points with input dimension $n = 10$, have been considered. In particular, the input
 395 and target output values have been randomly sampled from a uniform distribution over the
 396 interval $[-1, 1]$. In addition, in this experiment, the simulated annealing algorithm bundled
 397 in the Ocean SDK has been employed as the Ising machine for estimating the ground state and
 398 the corresponding energy value. Hence, no quantum annealing hardware has been used in this
 399 case. The parameters used for simulated annealing can be found directly in the source code
 400 at [40]. Instead, regarding the parameters of the proposed model, λ has been set to 1, and
 401 ϵ has been set according to (23) (taking a different value for each dataset). For the training
 402 process, $N_{\text{epochs}} = 50$ epochs have been executed, with $\eta = 0.2$. The MSE loss progression
 403 through the training is shown in Figure 4, where the error bars represent the standard deviation
 404 across the datasets.

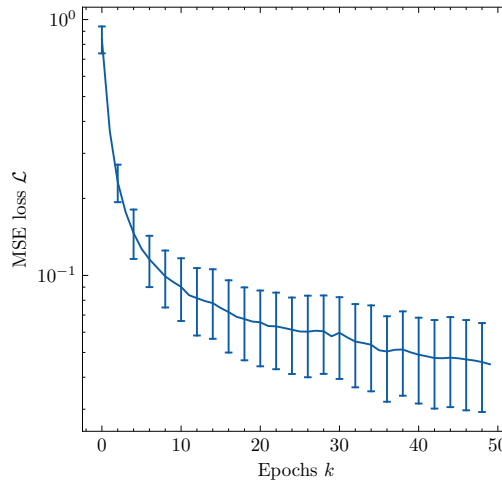


Figure 4: MSE loss on random data: Mean squared error, averaged over 30 randomly generated training sets of size $N = 20$. The MSE loss is tracked as a function of the number of epochs (with $N_{\text{epochs}} = 50$). The Ising machine in this experiment is the simulated annealing algorithm bundled in the Ocean SDK. The decreasing trend of the loss demonstrates the trainability of the model.

405 Although this particular example lacks practical significance, it serves as a simple demon-
 406 stration that the proposed Ising-machine-based parametric model can be effectively trained
 407 by utilizing its own output according to the update rules presented in Theorem 3. Further-
 408 more, it highlights the fact that the discontinuity observed in the derivative of the optimal

409 spin configuration \mathbf{z}^* , as discussed in the proof of Theorem 3, does not hinder the model’s
 410 ability to minimize the loss function. In essence, the assumption made in (15) regarding the
 411 computation of the partial derivatives proves to be sufficiently accurate.

412 4.3 Function approximation

413 In this second experiment, datasets comprising $N = 20$ data points sampled from polynomial
 414 functions have been considered. Due to the limited quantum annealing time available on the
 415 D-Wave hardware, the analysis has been limited to two straightforward cases, and no statistical
 416 repetition has been performed. Although this shortage prohibits any general conclusion on the
 417 model’s performance, it serves as a first demonstration of the possibility of using the model
 418 to approximate simple functions. Specifically, the following two polynomial functions of first
 419 and second degree, respectively, have been considered:

$$f_{\text{lin}}(x) = 2x - 6, \quad (24)$$

$$f_{\text{quad}}(x) = 1.2(x - 0.5)^2 - 2. \quad (25)$$

420 In both cases, the coefficients have been chosen manually and arbitrarily, and the input domain
 421 has been restricted to the interval $[0, 1]$. As the input dimensionality is $n = 1$, additional
 422 n_{hidden} hidden spins (see Section 3.3) have been considered. In particular, two different total
 423 sizes $n_{\text{total}} = \{50, 150\}$ have been analyzed in order to study the effect of the number of hidden
 424 spins on the model learning. Additionally, the spins have been initialized using the offset
 425 technique described in Section 3.3. Regarding the model parameters, fixed values have been
 426 manually chosen for the scaling factor λ , whereas the offset ϵ has again been set according to
 427 (23). All model parameters used for the two total sizes considered are summarized in Table 1.
 428 In this case, simulated and quantum annealing have been employed as Ising machines and
 429 compared. The simulated annealing parameters are the same as those used in Section 4.2.

Table 1: Parameters used to train the model for the function approximation task.

	n_{total}	d	λ	ϵ	N_{epochs}	η
f_{lin}	50	0.8/50	-0.3	-9.30	200	0.02
	150	0.8/150	-0.1	17.63	200	0.02
f_{quad}	50	1/50	-0.05	-2.70	200	0.25
	150	1/150	-0.0167	-4.23	200	0.25

430 The MSE loss throughout the training epochs for the two functions is shown in Figure 5.
 431 In the case of the linear function (Figure 5a), the model demonstrates a significant reduction
 432 in the mean squared error (MSE), over nearly three orders of magnitude, after approximately
 433 200 optimization steps. Instead, in the case of the quadratic function (Figure 5b), the initial
 434 loss was already low, indicating that the offset method chosen for the hidden layers was ap-
 435 propriate for this dataset. Nevertheless, the model has managed to decrease the loss by nearly
 436 additional three orders of magnitude. It is also worth noting that, in both cases, for equal
 437 model sizes, the results achieved using the quantum annealing hardware align closely with
 438 those obtained by employing the simulated annealing algorithm. Specifically, the fluctuations
 439 in the quantum annealing loss are caused by the very low number of reads (1), resulting in
 440 non-optimal solutions occasionally returned by the annealer. Finally, the higher number of
 441 hidden spins (150) has shown significant advantages only for the linear function.

442 Instead, Figure 6 displays the output of the trained models compared to the original func-
 443 tions. It is clear that the model has successfully learned to approximate the target functions.

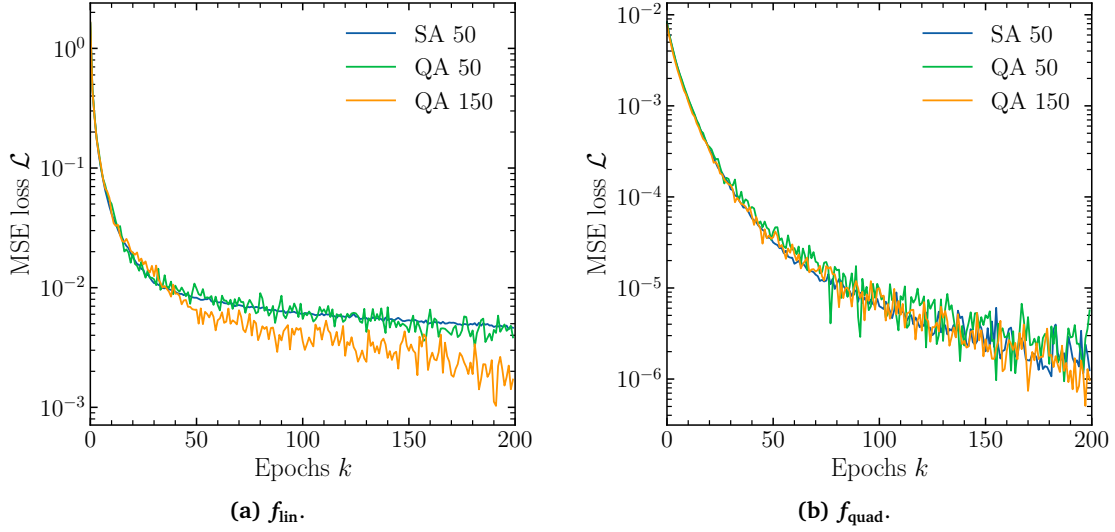


Figure 5: MSE loss in function approximation: Evolution of mean squared error loss during training for linear (a) and quadratic (b) functions. The results achieved by both simulated annealing (SA) and quantum annealing (QA) are shown, with the numeric value following the method name representing the total number of hidden spins n_{total} . SA and QA perform similarly with equal sizes, with the fluctuations of QA being caused by the very low number of reads (1). For f_{lin} , a larger number of hidden spins corresponds to better performance of QA.

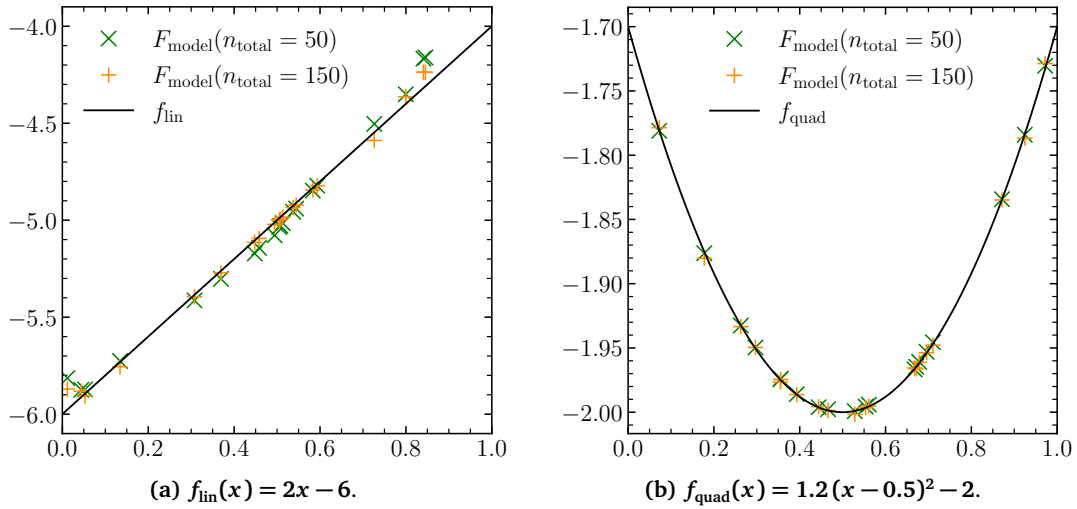


Figure 6: Trained model output: Output of the trained model F_{model} compared to the original function (black line). In both cases (linear and quadratic), for both n_{total} values, the model demonstrates the ability to approximate the function with a good accuracy, performing slightly worse for f_{lin} , especially toward the edges of the considered interval.

444 Specifically, as expected from the low final loss value, the model closely aligns with the original
 445 function in the case of the quadratic function. Instead, in the linear case, the model perfor-
 446 mance deteriorates significantly toward the interval edges, and the output values exhibit a
 447 tendency toward a shape resembling an even-degree polynomial, especially for the case with
 448 less hidden spins ($n_{\text{total}} = 50$). This behavior stems from the initialization method chosen for
 449 the hidden spins and the symmetry properties of the Ising model. At extreme bias values, lo-
 450 cated near the interval boundaries, the biases exert a dominant influence on the energy term

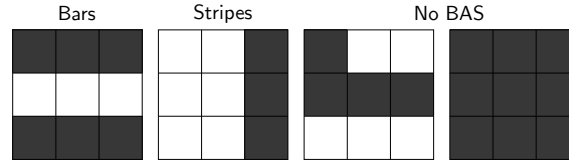


Figure 7: Bars and stripes (BAS) dataset: Illustration of exemplary 3×3 BAS and non-BAS data samples. The last two samples cannot be uniquely classified as bars or stripes and, therefore, are not part of the BAS dataset.

451 in Equation (4), causing $F(\theta) \rightarrow \infty$ as $|\theta| \rightarrow \pm\infty$. Consequently, the behavior resembles
 452 that of even polynomials, thus explaining the outliers in Figure 6a. Using more hidden spins
 453 ($n_{\text{total}} = 150$) reduces this effect by providing more trainable parameters to the model. It is
 454 also worth mentioning that different initialization methods for the hidden spins (e.g., taking
 455 the inverse values) influence this behavior.

456 4.4 Bars and stripes

457 In this last experiment, the proposed model has been applied to a different [machine learning](#)
 458 [machine-learning](#) task: binary classification. For this purpose, the well-known bars and stripes
 459 (BAS) dataset has been used. In detail, the dataset consists of square matrices with binary en-
 460 tries such that the values in the rows/columns are identical within each row/column; the
 461 resulting patterns can be identified as bars/stripes, giving the dataset its name. Actually, the
 462 cases in which all entries of the matrix are the same have been left out as the label is not
 463 unique. Some examples are shown in Figure 7. Regarding the classification task, it consists in
 464 assigning a label $l \in \{\text{bars, stripes}\}$ to each matrix, corresponding to the pattern it represents.
 465 In particular, the dataset was created by randomly deciding the label of each data point and
 466 randomly assigning one of the two binary values to each row/column. This procedure has
 467 been repeated N times, without accounting for duplicates.

468 In order to apply the proposed model to the BAS dataset, the input matrices have been
 469 flattened row-wise, and the binary values have been directly provided as input to the model.
 470 The binary labels $l \in \{\text{bars, stripes}\}$ have been encoded into y and decoded from the model
 471 output F_{model} according to

$$y = \begin{cases} 0 & , l = \text{bars} \\ 10 & , l = \text{stripes} \end{cases} \quad l_{\text{model}} = \begin{cases} \text{bars} & , F_{\text{model}} \leq 5 \\ \text{stripes} & , F_{\text{model}} > 5 \end{cases}, \quad (26)$$

472 with the factor **10** being arbitrarily chosen (different values can be used, but the λ and ϵ
 473 parameters must be adjusted accordingly). For the training, a randomly generated dataset
 474 comprising $N = 80$ data points, with each data point representing a BAS matrix of size 12×12 ,
 475 has been used. In particular, the model has been trained for $N_{\text{epochs}} = 8$ epochs, with $\eta = 0.02$,
 476 and has been evaluated on a separate test set consisting of other **80** data points. Since no
 477 additional hidden spins have been employed, $n = n_{\text{total}} = 144$ in this case. Concerning λ and
 478 ϵ , the former has been manually set to $\lambda = -0.3$, while the latter has been set to $\epsilon = -15.43$
 479 according to (23). Due to the large number of spins $n_{\text{total}} = 144$, only the quantum annealing
 480 hardware was used to train the model.

481 The results obtained are shown in Figure 8. Specifically, Figure 8a displays the model out-
 482 put during training for the training set and test set, respectively. The values shown are the
 483 average output values across all the data points with the same label, with the corresponding
 484 standard deviations indicated by the transparent envelopes. The dotted horizontal line repre-
 485 sents the classification threshold from (26). In practice, the average output value for the two
 486 labels diverges, approaching [toward](#) 0 and 10, respectively, as the number of epochs increases.

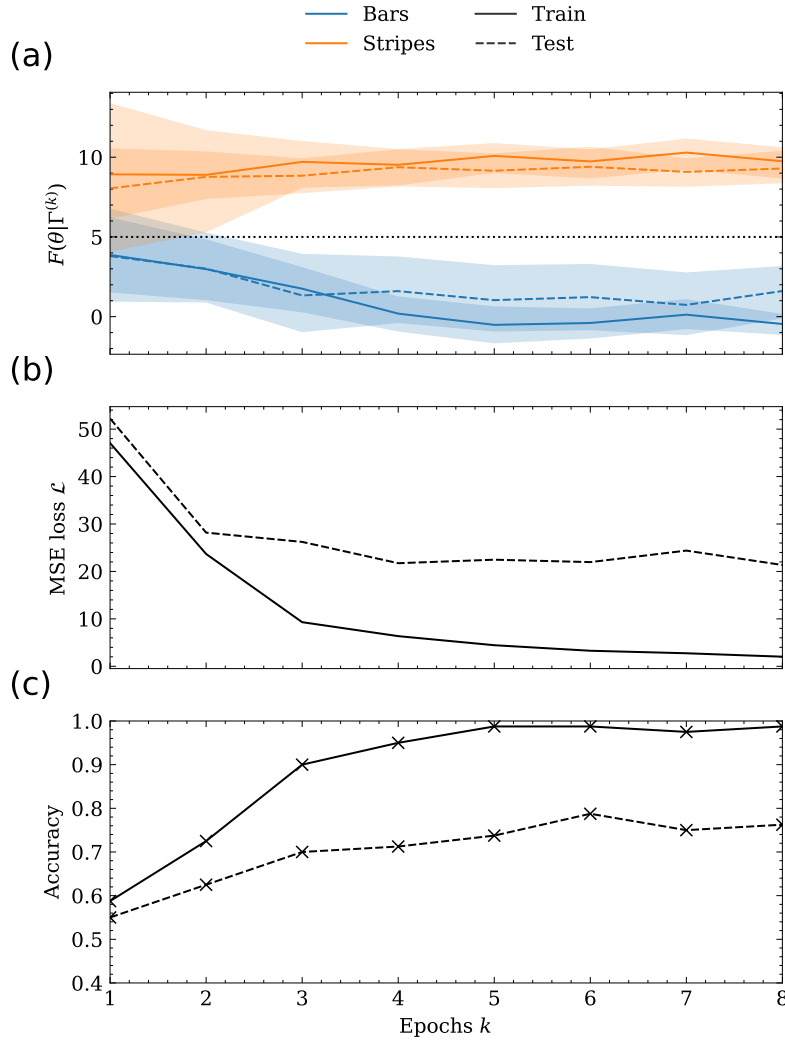


Figure 8: Results on BAS dataset: (a) Average model output value $F(\theta|\Gamma^{(k)}, \lambda, \epsilon)$ across all the data points with the same label $l \in \{\text{bars, stripes}\}$. The training (solid lines) and test (dashed lines) sets are considered independently; the envelopes represent the standard deviations, and the dotted horizontal line corresponds to the classification threshold according to (26). During training, the model learns to separate the two classes by increasing the energy for the stripes and decreasing it for the bars. (b) MSE loss for the training and test sets throughout epochs. The decreasing losses denote successful training, but the test loss stagnating after some epochs implies overfitting. (c) Accuracy for the training and test sets throughout the training. The accuracy on the training set reaches almost 1, with only one misclassified sample, while the accuracy on the test set also increases but saturates at about 75%.

487 This means that the model has learnt learned to increase the output value for stripe data points
488 and lower it for samples labeled as bars. This generalizes also to the unseen examples of the
489 test set, but the separation between the two classes is more marked for the training set. This
490 effect is also visible in Figure 8b, where the MSE loss for the training set and test set is shown.
491 In detail, the training loss decreases in a monotone way, while the test loss stagnates after a few
492 epochs. This is a typical indicator of model overfitting, which could be addressed in different
493 ways, among which increasing the number of training samples N in order to help the model
494 generalize. A similar conclusion can be drawn considering the accuracy of the model shown
495 in Figure 8c. The trained model is able to correctly classify 79 out of 80 training samples, but
496 the accuracy on the test set saturates at only about 75%.

497 In conclusion, this experiment has demonstrated the possibility of using the proposed
498 model to address also binary classification tasks by choosing an appropriate encoding-decoding
499 procedure for the model input and output. Indeed, the model has proven to be able to gener-
500 alize to unseen examples while exhibiting overfitting effects, at least for the chosen dataset.

501 4.5 Choice of hyperparameters

502 Selecting appropriate values for the model’s hyperparameters is a common issue in machine
503 learning. Multiple hyperparameters have been manually set in the experiments presented
504 in this work. These include the learning rate η , the number of epochs N_{epochs} , the problem
505 encoding (see 26), the Ising machine parameters like the number of samples per step for simu-
506 lated annealing or the embedding procedure, the annealing time, and the number of reads for
507 quantum annealing. Choosing appropriate values may reduce, for example, the fluctuations
508 observed in Figure 6a. The values used here have been selected based on observations result-
509 ing from trial and error runs; the analysis of different configurations and a more systematic
510 approach to choosing appropriate values are left for future work.

511 Among the model-related hyperparameters, the choice of the initialization strategy for the
512 additional hidden spins has a significant impact. Specifically, when the input dimension is low,
513 a large number of hidden spins $n_{\text{hidden}} \gg n$ may be necessary in order to have enough trainable
514 model parameters. However, particular care must be put in choosing the corresponding new
515 bias terms. Indeed, in preliminary experiments, it has been observed that initializing the biases
516 in the wrong way may negatively affect the performance to the point that the model is unable
517 to approximate the target function. Finding suitable ansätze for different tasks is still an open
518 question.

519 5 Conclusion

520 In this paper, we have proposed a novel parametric learning model that leverages the inher-
521 ent structure of the Ising model for training purposes. We have presented a straightforward
522 optimization procedure based on gradient descent and we have provided the rules for com-
523 puting all relevant derivatives of the mean squared error loss. Notably, if the Ising machine
524 is realized by a quantum platform, our approach allows for the utilization of quantum re-
525 sources for both the execution and the training of the model. Experimental results using a
526 D-Wave quantum annealer have demonstrated the successful training of our model on simple
527 proof-of-concept datasets, specifically for linear and quadratic function approximations and
528 binary classification. This novel approach unveils the potential of employing Ising machines,
529 particularly quantum annealers, for general learning tasks. In addition, it raises intriguing the-
530 oretical and practical questions from both computer science and physics perspectives. From a
531 theoretical standpoint, questions regarding the expressibility of the Ising model arise, as well
532 as inquiries into the classes of functions that the model can represent. These questions are

533 non-trivial due to the non-linear minimization step involved. From a practical point of view,
534 given the broad definition of the model and its similarity to other classical parametric models,
535 a wide range of machine learning tools and methods can be explored to enhance its training.
536 Advanced gradient-based optimizers and general learning techniques such as mini-batching,
537 early stopping, and dropout, among others, offer promising avenues for improvement.

538 In addition to function approximation and binary classification, we aim to investigate the
539 application of the model to other machine learning tasks, especially tasks in which the feature
540 space is large, to reduce the necessity of additional hidden spins. This study might be extended
541 with a comparison to other Ising machine-based models advancing the field of parametric
542 machine learning models utilizing Ising machines.

543 **Funding information** This work was partially supported by project SERICS (PE00000014)
544 under the MUR National Recovery and Resilience Plan funded by the European Union - NextGen-
545 erationEU. E.Z. was supported by Q@TN, the joint lab between University of Trento, FBK-
546 Fondazione Bruno Kessler, INFN-National Institute for Nuclear Physics and CNR-National Re-
547 search Council. The authors gratefully acknowledge CINECA for providing computing time on
548 the D-Wave quantum annealer within the project “Testing the learning performances of quan-
549 tum machines”, and the Jülich Supercomputing Center for providing computing time on the
550 D-Wave quantum annealer through the Jülich UNified Infrastructure of Quantum computing
551 (JUNIQ).

552 References

- 553 [1] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, Prentice Hall, 3 edn.
554 (2010).
- 555 [2] B. Cheng and D. M. Titterton, *Neural Networks: A Review from a Statistical Perspective*,
556 *Statistical Science* **9**(1), 2 (1994), doi:[10.1214/ss/1177010638](https://doi.org/10.1214/ss/1177010638).
- 557 [3] J. Zou, Y. Han and S.-S. So, *Overview of Artificial Neural Networks*, pp. 14–22, Hu-
558 mana Press, Totowa, NJ, ISBN 978-1-60327-101-1, doi:[10.1007/978-1-60327-101-1_2](https://doi.org/10.1007/978-1-60327-101-1_2)
559 (2009).
- 560 [4] L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, Y. Duan, O. Al-Shamma, J. Santa-
561 maría, M. A. Fadhel, M. Al-Amidie and L. Farhan, *Review of deep learning: concepts, cnn*
562 *architectures, challenges, applications, future directions*, *J Big Data* **53** (2021).
- 563 [5] M. Benedetti, E. Lloyd, S. Sack and M. Fiorentini, *Parameterized quantum circuits*
564 *as machine learning models*, *Quantum Science and Technology* **4**(4), 043001 (2019),
565 doi:[10.1088/2058-9565/ab4eb5](https://doi.org/10.1088/2058-9565/ab4eb5).
- 566 [6] S. Y.-C. Chen, C.-H. H. Yang, J. Qi, P.-Y. Chen, X. Ma and H.-S. Goan, *Variational*
567 *quantum circuits for deep reinforcement learning*, *IEEE Access* **8**, 141007 (2020),
568 doi:[10.1109/ACCESS.2020.3010470](https://doi.org/10.1109/ACCESS.2020.3010470).
- 569 [7] M. Schuld and F. Petruccione, *Machine Learning with Quantum Computers*, Springer
570 Cham, ISBN 9783030830977 (2021).
- 571 [8] D. Pastorello, *Concise Guide to Quantum Machine Learning*, Springer Singapore, ISBN
572 9789811968969 (2023).
- 573 [9] D. Willsch, M. Willsch, H. De Raedt and K. Michielsen, *Support vector machines on the*
574 *d-wave quantum annealer*, *Computer Physics Communications* **248**, 107006 (2020).

- 575 [10] R. K. Nath, H. Thapliyal and T. S. Humble, *A review of machine learning classification*
576 *using quantum annealing for real-world applications*, SN COMPUT. SCI **365** (2021).
- 577 [11] H. Wang, W. Wang, Y. Liu and B. Alidaee, *Integrating machine learning algorithms with*
578 *quantum annealing solvers for online fraud detection*, IEEE Access **10**, 75908 (2022),
579 doi:[10.1109/ACCESS.2022.3190897](https://doi.org/10.1109/ACCESS.2022.3190897).
- 580 [12] P. Rebentrost, M. Mohseni and S. Lloyd, *Quantum support vector machine for big data clas-*
581 *sification*, Phys. Rev. Lett. **113**, 130503 (2014), doi:[10.1103/PhysRevLett.113.130503](https://doi.org/10.1103/PhysRevLett.113.130503).
- 582 [13] J. Preskill, *Quantum Computing in the NISQ era and beyond*, Quantum **2**, 79 (2018),
583 doi:[10.22331/q-2018-08-06-79](https://doi.org/10.22331/q-2018-08-06-79).
- 584 [14] D. H. Ackley, G. E. Hinton and T. J. Sejnowski, *A learning algorithm for Boltzmann Ma-*
585 *chines*, Cognitive Science **9**(1), 147 (1985).
- 586 [15] C. Bybee, D. Kleyko, D. E. Nikonov, A. Khosrowshahi, B. A. Olshausen and F. T. Sommer,
587 *Efficient optimization with higher-order ising machines*, arXiv preprint arXiv:2212.03426
588 (2022).
- 589 [16] N. Mosheni, P. McMahon and T. Byrnes, *Ising machines as hardware solvers of combinato-*
590 *rial optimization problems*, Nat Rev Phys **4**, 363 (2022).
- 591 [17] N. Mwamsojo, F. Lehmann, K. Merghem, B.-E. Benkelfat and Y. Frignac, *Optoelectronic*
592 *coherent ising machine for combinatorial optimization problems*, Optics Letters **48**(8),
593 2150 (2023), doi:[10.1364/OL.485215](https://doi.org/10.1364/OL.485215).
- 594 [18] M. H. Amin, E. Andriyash, J. Rolfe, B. Kulchytskyy and R. Melko, *Quantum boltzmann*
595 *machine*, Phys. Rev. X **8**, 021050 (2018), doi:[10.1103/PhysRevX.8.021050](https://doi.org/10.1103/PhysRevX.8.021050).
- 596 [19] K. Mitarai, M. Negoro, M. Kitagawa and K. Fujii, *Quantum circuit learning*, Phys. Rev. A
597 **98**, 032309 (2018), doi:[10.1103/PhysRevA.98.032309](https://doi.org/10.1103/PhysRevA.98.032309).
- 598 [20] D. Sherrington and S. Kirkpatrick, *Solvable model of a spin-glass*, Phys. Rev. Lett. **35**,
599 1792 (1975), doi:[10.1103/PhysRevLett.35.1792](https://doi.org/10.1103/PhysRevLett.35.1792).
- 600 [21] S. Tanaka, Y. Matsuda and N. Togawa, *Theory of ising machines and a common software*
601 *platform for ising machines*, In *2020 25th Asia and South Pacific Design Automation Con-*
602 *ference (ASP-DAC)*, pp. 659–666, doi:[10.1109/ASP-DAC47756.2020.9045126](https://doi.org/10.1109/ASP-DAC47756.2020.9045126) (2020).
- 603 [22] A. Lucas, *Ising formulations of many NP problems*, Front. Phys. **2** (2014),
604 doi:[10.3389/fphy.2014.00005](https://doi.org/10.3389/fphy.2014.00005).
- 605 [23] P. Hauke, H. G. Katzgraber, W. Lechner, H. Nishimori and W. D. Oliver, *Perspectives of*
606 *quantum annealing: Methods and implementations*, Reports on Progress in Physics **83**(5),
607 054401 (2020), doi:[10.1088/1361-6633/ab85b8](https://doi.org/10.1088/1361-6633/ab85b8).
- 608 [24] P. L. McMahon, A. Marandi, Y. Haribara, R. Hamerly, C. Langrock, S. Tamate, T. Inagaki,
609 H. Takesue, S. Utsunomiya, K. Aihara, R. L. Byer, M. M. Fejer *et al.*, *A fully programmable*
610 *100-spin coherent ising machine with all-to-all connections*, Science **354**(6312), 614
611 (2016), doi:[10.1126/science.aah5178](https://doi.org/10.1126/science.aah5178), [https://www.science.org/doi/pdf/10.1126/](https://www.science.org/doi/pdf/10.1126/science.aah5178)
612 [science.aah5178](https://www.science.org/doi/pdf/10.1126/science.aah5178).
- 613 [25] T. Inagaki, K. Inaba, R. Hamerly, K. Inoue, Y. Yamamoto and H. Takesue, *Large-scale ising*
614 *spin network based on degenerate optical parametric oscillators*, Nature Photon **10**, 415
615 (2016).

- 616 [26] T. Inagaki, Y. Haribara, K. Igarashi, T. Sonobe, S. Tamate, T. Honjo, A. Marandi,
617 P. L. McMahon, T. Umeki, K. Enbutsu, O. Tadanaga, H. Takenouchi *et al.*, *A co-*
618 *herent ising machine for 2000-node optimization problems*, *Science* **354**(6312), 603
619 (2016), doi:[10.1126/science.aah4243](https://doi.org/10.1126/science.aah4243), [https://www.science.org/doi/pdf/10.1126/](https://www.science.org/doi/pdf/10.1126/science.aah4243)
620 [science.aah4243](https://www.science.org/doi/pdf/10.1126/science.aah4243).
- 621 [27] Y. Yamamoto, T. Leleu, S. Ganguli and H. Mabuchi, *Coherent ising machines—quantum*
622 *optics and neural network perspectives*, *Appl. Phys. Lett.* **117**, 160501 (2020),
623 doi:[10.1063/5.0016140](https://doi.org/10.1063/5.0016140).
- 624 [28] S. Kirkpatrick, C. D. Gelatt and M. P. Vecchi, *Optimization by Simulated Annealing*, *Science*
625 **220**(4598), 671 (1983), doi:[10.1126/science.220.4598.671](https://doi.org/10.1126/science.220.4598.671).
- 626 [29] T. Guilmeau, E. Chouzenoux and V. Elvira, *Simulated annealing: a review and a*
627 *new scheme*, In *2021 IEEE Statistical Signal Processing Workshop (SSP)*, pp. 101–105,
628 doi:[10.1109/SSP49050.2021.9513782](https://doi.org/10.1109/SSP49050.2021.9513782) (2021).
- 629 [30] A. B. Finnila, M. A. Gomez, C. Sebenik, C. Stenson and J. D. Doll, *Quantum annealing: A*
630 *new method for minimizing multidimensional functions*, *Chemical Physics Letters* **219**(5),
631 343 (1994), doi:[10.1016/0009-2614\(94\)00117-0](https://doi.org/10.1016/0009-2614(94)00117-0).
- 632 [31] T. Kadowaki and H. Nishimori, *Quantum annealing in the transverse Ising model*, *Physical*
633 *Review E* **58**(5), 5355 (1998), doi:[10.1103/PhysRevE.58.5355](https://doi.org/10.1103/PhysRevE.58.5355).
- 634 [32] C. C. McGeoch, *Adiabatic Quantum Computation and Quantum Annealing*, Springer
635 Cham, ISBN 978-3-031-01390-4 (2014).
- 636 [33] K. Kitai, J. Guo, S. Ju, S. Tanaka, K. Tsuda, J. Shiomi and R. Tamura, *Designing meta-*
637 *materials with quantum annealing and factorization machines*, *Physical Review Research*
638 **2**(1), 013319 (2020), doi:[10.1103/PhysRevResearch.2.013319](https://doi.org/10.1103/PhysRevResearch.2.013319).
- 639 [34] Y. Seki, R. Tamura and S. Tanaka, *Black-box optimization for integer-variable problems*
640 *using Ising machines and factorization machines*, doi:[10.48550/arXiv.2209.01016](https://doi.org/10.48550/arXiv.2209.01016) (2022),
641 [2209.01016](https://doi.org/10.48550/arXiv.2209.01016).
- 642 [35] P. Date, R. Patton, C. Schuman and T. Potok, *Efficiently embedding QUBO problems*
643 *on adiabatic quantum computers*, *Quantum Information Processing* **18**(4), 117 (2019),
644 doi:[10.1007/s11128-019-2236-3](https://doi.org/10.1007/s11128-019-2236-3).
- 645 [36] S. Mukherjee and B. Chakrabarti, *Multivariable optimization: Quantum annealing*
646 *and computation*, *The European Physical Journal Special Topics* **224**(1), 17 (2015),
647 doi:[10.1140/epjst/e2015-02339-y](https://doi.org/10.1140/epjst/e2015-02339-y).
- 648 [37] P. Date, D. Arthur and L. Pusey-Nazzaro, *QUBO formulations for training machine*
649 *learning models*, *Scientific Reports* **11**(1), 10029 (2021), doi:[10.1038/s41598-021-](https://doi.org/10.1038/s41598-021-89461-4)
650 [89461-4](https://doi.org/10.1038/s41598-021-89461-4).
- 651 [38] Github, *Ising Learning Model* (2023), [https://github.com/lsschmid/](https://github.com/lsschmid/ising-learning-model)
652 [ising-learning-model](https://github.com/lsschmid/ising-learning-model).
- 653 [39] D-Wave Systems Inc., *D-wave ocean software*, [https://docs.ocean.dwavesys.com/en/](https://docs.ocean.dwavesys.com/en/stable/)
654 [stable/](https://docs.ocean.dwavesys.com/en/stable/).
- 655 [40] L. Schmid, E. Zardini and D. Pastorello, *Evaluation data for "A general learning scheme*
656 *for classical and quantum Ising machines"*, doi:[10.5281/zenodo.10031307](https://doi.org/10.5281/zenodo.10031307) (2023).