

# Implicit Bias and Invariance: How Hopfield Networks Efficiently Learn Graph Orbits

Anonymous authors

Paper under double-blind review

## Abstract

Many learning problems involve symmetries, and while invariance can be built into neural architectures, it can also emerge implicitly when training on group-structured data. We study this phenomenon in classical Hopfield networks and illustrate how they can infer the isomorphism class of a graph from a small, random sample. Our results reveal that: (i) graph isomorphism classes can be represented within a three-dimensional invariant subspace, (ii) using gradient descent to minimize energy flow (MEF) has an implicit bias toward norm-efficient solutions, which underpins a polynomial sample complexity bound for learning isomorphism classes, and (iii) across multiple learning rules, parameters converge toward the invariant subspace as sample sizes grow. Together, these findings highlight a unifying mechanism for generalization in Hopfield networks: a bias toward norm efficiency in learning drives the emergence of approximate invariance under group-structured data <sup>1</sup>.

## 1 Introduction

Here, we analyze the emergence of invariance arising implicitly during training in Hopfield networks (HNs) (Hopfield, 1982), which represent arguably the simplest example of an Associative Memory. Building on classical ideas Rosenblatt (1958); Willshaw et al. (1969); Amari (1972); Little (1974); Pastur & Figotin (1977), HNs are recurrent neural networks consisting of  $n$  linear-threshold McCulloch–Pitts neurons McCulloch & Pitts (1943) that can store binary patterns as distributed memories in the form of fixed-point attractors of recurrent dynamics. In the literature, HNs are usually associated with a particular Hebbian learning scheme called the Outer-Product Rule, but for the purposes of this work we also consider other standard training methods. This setting is intentionally minimal so that we can focus on developing novel mathematical tools for understanding generalization in a classical architecture. As data symmetry is not made explicit in this model, any invariance must arise from the interplay of the group structure in the data and the implicit bias of the learning rule in question. More specifically, and inspired by Hillar & Tran (2018); Hillar et al. (2021), this paper studies whether or not standard learning rules and objectives, notably minimization of the energy flow (MEF) Hillar et al. (2012), a tractable convex loss, *can learn the isomorphism class of a graph from a small, random subset*. Our key findings are as follows.

1. **HNs can memorize any graph isomorphism class.** We characterize the subspace of parameters invariant to edge-adjacency-preserving permutations (of which graph isomorphisms are a subset) (Lemma 4.6) and observe that this subspace aligns well with the parameters of successfully trained models (Fig. 5). Moreover, for any graph we give an explicit construction within this space that memorizes it (Lemma 4.7) as well as its isomorphism class.
2. **Implicit bias towards norm efficient solutions.** We reparameterize the MEF objective and show that gradient descent on it is directionally biased towards the solution to a hard-margin support vector machine (HSVM) problem on an induced linear representation (Corollary 3.1).

---

<sup>1</sup>**Statement on the use of LLMs.** Large language models (LLMs) were used to assist with literature search, checking and refining the clarity of writing, high level ideation and planning as well as organizing related work.

3. **Polynomial sample complexity suffices for orbit generalization.** Suppose  $D \subset \{0, 1\}^n$  is strictly memorizable with min-norm parameter  $\theta^*$ ,  $\|\mathbf{x}\|_0 \leq m$  for all  $\mathbf{x} \in D$ , and let  $\mathcal{D}$  be a distribution supported on  $D$ . We prove  $N = \tilde{\Omega}(n\|\theta^*\|^2 m \epsilon^{-2})$  random samples suffices for both HSVM and MEF to memorize new samples with probability at least  $1 - \epsilon$  (Theorem 3.2). These theoretical results corroborate the empirical “few-shot-to-orbit” phenomenon we observe (Figs. 1, 2, 6, 8, 9, 10). Moreover, specializing to isomorphism classes this result implies a polynomial sample complexity in the number of vertices  $v$ , supporting a conjecture in Hillar et al. (2021) for the case of cliques. In this case, the critical number of samples for generalization appears to coincide with when these networks solve the Hidden Clique Problem (HCP) Dekel et al. (2014); see Fig. 6 in Appendix C.3. These critical sample counts can also be used to differentiate between combinatorial classes of graphs (see Fig. 2).
4. **Emergence of invariance.** We observe that as the sample size  $N$  grows the learned parameters concentrate near the invariant subspace (Figs. 3, 5, 6, 9). For a simplified average HSVM surrogate, we prove the sample solution converges to the invariant set at rate  $\tilde{O}(v^2/\sqrt{N})$  (Lemma 4.10), where  $v$  is the number of vertices of the graph.

## 1.1 Related work

**Capacity of Hopfield networks.** The capacity of a Hopfield network depends on the learning rule and the structure of the data. For dense, uncorrelated random patterns under Hebbian learning, the statistical–mechanics analysis of (Amit et al., 1985) gives the classic linear law: reliable retrieval up to approximately 0.138 patterns per neuron with subsequent refinements via replica methods (Gardner, 1988; Krauth & Mézard, 1989). Coding-theoretic analyses further show that, for Hebbian constructions and exact recovery of randomly chosen patterns, one typically cannot exceed  $n/(4 \ln n)$  (McEliece et al., 1987) memories. More generally, Cover’s classical bound Cover (1965) restricts the capacity for exact storage of dense, random data to only  $2n$ . Nonetheless, superlinear capacity is achievable for certain structured datasets. For example, sparse data having few active neurons can yield an increase of capacity to nearly quadratic in  $n$  (Tsodyks & Feigel’man, 1988; Amari, 1989). Additionally, robust exponential memory has been observed for particular examples of group structured data (Hillar & Tran, 2018; Hillar et al., 2021); in particular, for storing all  $k$ -clique graphs and their hypergraph analogues. Our work builds on the observations of (Hillar & Tran, 2018; Hillar et al., 2021) by proving that all graph isomorphism classes are memorizable.

**Modern Hopfield Networks.** A line of recent investigation has sought to increase the capacity and retrieval properties of Hopfield networks by changing the energy function. Dense Associative Memories (DAMs) replace the classical quadratic energy with higher-order polynomial interactions, resulting in a capacity that scales polynomially with neuron count (Krotov & Hopfield, 2016; Horn & Usher, 1988). Building on this, Modern Hopfield Networks (MHNs) introduced a log-sum-exp energy function that allows the capacity to grow exponentially (Ramsauer et al., 2021; Demircigil et al., 2017). Our work provides a complementary perspective to these advancements by showing classical HNs can achieve exponential capacity by capturing the symmetry of the data in its parameters. Preliminary experiments suggest that DAMs do not generalize well in this setting (see Fig. 7, Appendix C.4).

**Generalization beyond the training set for HNs.** Theoretical study of classical HNs primarily focuses on storage limits, basins of attraction, and noise robustness around memorized patterns, rather than sample–complexity guarantees for generalization to patterns outside the training set. Earlier analyses of concept generalization in classical HNs investigate when networks capture latent data regularities (Fontanari, 1990). More recently, Random-Features Hopfield Models (RFHMs) exhibit learning phase transitions and even retrieval of previously unseen examples (Negri et al., 2023; Kalaj et al., 2025). These results complement but are distinct from our own; in particular, they do not provide sample–complexity bounds or analyze the emergence of invariance induced by a symmetry present in the data. In addition, while these results primarily use techniques from statistical physics, here we leverage tools from statistical learning theory.

**Emergent invariance through data augmentation and feature averaging.** One perspective on data augmentation is as orbit averaging over a symmetry group; in particular, empirical risk minimization on augmented data is equivalent to averaging features or predictions along group orbits. This has been

shown to induce approximate invariance and reduces estimator variance (Chen et al., 2020). From a kernel perspective, augmentation decomposes into first-order invariant feature averaging plus a second-order variance regularizer (Dao et al., 2019). Enforcing invariance through this averaging method yields provable generalization benefits in the context of invariant kernel regression (Elesedy & Zaidi, 2021). Beyond static estimators, recent results show that with full group augmentation deep ensembles become equivariant in expectation at all training times in the infinite-width limit (Gerken & Kessel, 2024) and that the expected predictions of group-convolutional networks match those of data-augmented conventional networks throughout training (Misof et al., 2025). Note a formal comparison between group equivariant architectures and data augmentation, without resorting to the infinite-width limit, is provided in Nordenfors et al. (2024). Finally, in regard to signatures of feature learning in the context of group structured data, then under certain conditions if a neural network is invariant to a finite group its trained weights recover the Fourier transform on that group (Marchetti et al., 2024). While these results assume explicit invariance, either through architectural design or by averaging over the full group orbit, here we ask whether simple learning rules can implicitly recover approximate invariance from a small, random sample of elements.

**Implicit bias.** A large body of work shows that, even without explicit regularization, certain learning dynamics have a preference for particular solutions. In particular, for classification using the logistic loss, gradient descent drives the parameter norm to infinity while the parameter direction converges to the max-margin classifier (Soudry et al., 2018; Ji & Telgarsky, 2018; Nacson et al., 2019). Our work leverages these results in order to show that standard learning rules for HNs are implicitly biased towards learning invariant representations when trained on group data.

## 2 Preliminaries

**Notation:** we use capitalized boldface characters to denote matrices, bold lowercase characters to denote vectors and non-bold lowercase characters to denote scalar values. If  $\mathbf{x} \in \mathbb{R}^n$  is a vector then  $x_i$  denotes the  $i$ th entry of  $\mathbf{x}$ . If  $\mathbf{X} \in \mathbb{R}^{N \times n}$  then  $\mathbf{x}_i \in \mathbb{R}^n$  denotes the  $i$ th row of  $\mathbf{X}$  and to access individual entries of  $\mathbf{X}$  we use the notation  $x_{ij}$  or  $[\mathbf{X}]_{ij}$ , whichever is clearer in context. Whether a matrix, vector or scalar is deterministic or random is also inferred from context. We use  $\Pi_a$  to denote the set of permutations on  $[a]$  and  $\mathcal{P}_a$  to denote the set of  $a \times a$  permutation matrices. Overloading our notation, we also refer to  $\mathcal{P}_a$  as the group of permutation matrices. Finally, if  $\mathcal{H}$  is a group that acts on a set  $\mathcal{A}$ , then the orbit of  $a \in \mathcal{A}$  under  $\mathcal{H}$  is denoted  $\text{Orb}(a, \mathcal{H}) = \{ha \in \mathcal{A} : h \in \mathcal{H}\}$ .

**Associative Memory:** we consider a Hopfield network Hopfield (1982) with asynchronous dynamics but do not restrict ourselves to Hebbian learning. To this end, let  $\text{Sym}_0^n \subset \mathbb{R}^{n \times n}$  denote the set of symmetric, real,  $n \times n$  matrices whose diagonal entries are zero, and let  $\Theta = \text{Sym}_0^n \times \mathbb{R}^n$ . Clearly  $\Theta$  is a convex vector space. We introduce the energy function  $E : \{0, 1\}^n \times \Theta \rightarrow \mathbb{R}$  defined as

$$E(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{2} \mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{b}^T \mathbf{x}, \quad (1)$$

where  $\boldsymbol{\theta} = (\mathbf{W}, \mathbf{b}) \in \Theta$ . Given an input binary vector  $\mathbf{x} \in \{0, 1\}^n$ , the Hopfield network generates a sequence of binary vectors  $(\mathbf{x}(t))_{t \geq 0}$  through the following recurrence dynamics: with  $\mathbf{x}(0) = \mathbf{x}$  then

$$x_j(t) = \begin{cases} \mathbb{1}(-\mathbf{w}_j^T \mathbf{x}(t-1) > b_j) & t \equiv j \pmod{n}, \\ x_j(t-1) & \text{otherwise} \end{cases} \quad (2)$$

for all  $t \geq 1$  and  $j \in [n]$ . For any input  $\mathbf{x} \in \{0, 1\}^n$ , this sequence converges in finite time to a fixed point Bruck (1990). We define the input-output map of the Hopfield network, denoted  $H : \{0, 1\}^n \times \Theta \rightarrow \{0, 1\}^n$  as follows: given parameters  $\boldsymbol{\theta}$  and an input  $\mathbf{x}$ , the output  $H_\theta(\mathbf{x}; \boldsymbol{\theta})$  is the attractor or fixed point of 2 reached when initialized with  $\mathbf{x}(0) = \mathbf{x}$ . If  $H(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{x}$ , then  $\mathbf{x}$  is a fixed point of the recurrence dynamics, furthermore in this setting we say that the function  $H_\theta(\cdot) := H(\cdot; \boldsymbol{\theta})$  has *memorized*  $\mathbf{x}$ . A sufficient condition for  $H_\theta$  to memorize  $\mathbf{x}$  is  $E(\mathbf{x}; \boldsymbol{\theta}) < E(\mathbf{x}'; \boldsymbol{\theta})$  for all  $\mathbf{x}' \in \mathcal{N}(\mathbf{x})$ , where here  $\mathcal{N}(\mathbf{x})$  denotes the set of all binary vectors a Hamming distance of exactly one from  $\mathbf{x}$ . Under this condition we say that  $H_\theta$  *strictly memorizes*  $\mathbf{x}$ . We also denote the action of a permutation matrix  $\mathbf{P} \in \mathcal{P}_n$  on the parameters of a Hopfield network as  $\mathbf{P}\boldsymbol{\theta} := (\mathbf{P}^T \mathbf{W} \mathbf{P}, \mathbf{P}^T \mathbf{b})$ .

**Training and memorization:** let  $\mathcal{S} \subset \{0, 1\}^n$ , we say that  $H_\theta$  memorizes  $\mathcal{S}$  if it memorizes all  $\mathbf{x} \in \mathcal{S}$ . There are many methods Hertz et al. (1991); Tolmachev & Manton (2020) that have been proposed to train networks to memorize a set  $\mathcal{S}$  (Appendix A.2). We focus on minimization of the energy flow (MEF) Hillar et al. (2012); Hillar & Tran (2018); Hillar et al. (2021) and study its implicit bias. If  $\mathbf{x} \in \{0, 1\}^n$  and  $j \in [n]$ , let  $\mathbf{x}^{(j)} \in \{0, 1\}^n$  satisfy  $x_l \neq x_l^{(j)}$  iff  $l = j$ . We define the *energy flow* loss as:

$$L(\theta; \mathcal{S}) = \sum_{\mathbf{x} \in \mathcal{S}} \sum_{j=1}^n \exp \left( E(\mathbf{x}; \theta) - E(\mathbf{x}^{(j)}; \theta) \right). \quad (3)$$

For any given set of points  $\mathcal{S}$ , note that  $L$  is nonnegative, infinitely differentiable, and is convex in  $\Theta$ . As a result, minimizing  $L$  is a convex problem to which a wide variety of numerical techniques can be applied, including, but not limited to, variants of gradient descent (GD) as well as (approximate) second order methods, such as L-BFGS Liu & Nocedal (1989). As long as  $\mathcal{S}$  can be memorized, then sufficiently minimizing 3 will result in a network which memorizes  $\mathcal{S}$ . For further details on MEF learning for discrete recurrent neural networks, we refer the reader to Hillar et al. (2021) and its inspiration from the theory of density estimation Sohl-Dickstein et al. (2011).

### 3 Implicit bias, minimum norm memorizers and generalization

In this section, and prior to specializing to study datasets drawn from isomorphism classes of graphs, we connect memorization to solving a linear program and identify the implicit bias of MEF. This enables us to provide generalization guarantees for memorization in Hopfield networks as per Theorem 3.2. Given  $\theta = (\mathbf{W}, \mathbf{b}) \in \Theta$ , for any  $j \in [n]$  we define the vector  $\theta_j = [\mathbf{w}_j, b_j] \in \mathbb{R}^{n+1}$ . Overloading our notation, we also use  $\theta = [\theta_1, \theta_2 \dots \theta_n] \in \mathbb{R}^{n, n+1}$  to refer to the flattened vector of all the network parameters  $(\mathbf{W}, \mathbf{b})$ . For any  $\mathbf{x} \in \{0, 1\}^n$  let  $\mathbf{z}(\mathbf{x}) = [\mathbf{x}, 1] \in \{0, 1\}^{n+1}$  and  $y_j(\mathbf{x}) = 1 - 2x_j \in \{\pm 1\}$ . Using this notation it is well known that the energy difference between a point and one of its neighbors is

$$E(\mathbf{x}^{(j)}; \theta) - E(\mathbf{x}; \theta) = y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle, \quad (4)$$

see Appendix A.1 for further details. As a consequence, parameters which strictly memorize a set  $\mathcal{S} \subseteq \{0, 1\}^n$  must satisfy a system of linear inequalities: in particular, there must exist some  $\epsilon > 0$ , referred to as the functional margin, such that  $y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle \geq \epsilon$  for all  $\mathbf{x} \in \mathcal{S}$  and  $j \in [n]$ . Clearly the energy function (1) is quadratic in the inputs  $\mathbf{x}$  but linear in the parameters,  $E(\mathbf{x}; a\theta) = aE(\mathbf{x}, \theta)$ . Moreover, this implies the energy is positively homogeneous of degree 1 in the parameters and as a result the set of attractors of a Hopfield network is invariant under positive rescaling of the parameters. Without loss of generality, we therefore select a functional margin of one and define the feasible set of parameters up to positive rescaling as

$$\mathcal{F}_\theta(\mathcal{S}) = \{ \theta \in \Theta : y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle \geq 1 \quad \forall \mathbf{x} \in \mathcal{S}, \forall j \in [n] \} \quad (5)$$

In addition, the inequality constraints that define  $\mathcal{F}(\mathcal{S})$  can be written with respect to a single vector of unique parameters, which we denote  $\omega$ . To this end let  $p = \frac{n(n-1)}{2}$  and  $q = \frac{n(n+1)}{2}$ . There exists a  $\mathbf{V} \in \{0, 1, 1/\sqrt{2}\}^{n(n+1) \times q}$  such that for any  $\theta \in \Theta$  there exists an  $\mathbf{a} \in \mathbb{R}^p$  with  $\omega = [\sqrt{2}\mathbf{a}, \mathbf{b}] \in \mathbb{R}^q$ , such that  $\theta = \mathbf{V}\omega$ . In short,  $\mathbf{V}$  copies the unique elements, i.e., the upper triangular elements of  $\mathbf{W}$  and the biases  $\mathbf{b}$ , into their appropriate locations in the flattened vector  $\theta$ . For any  $j \in [n]$  let  $\mathbf{V}_j \in \mathbb{R}^{(n+1) \times q}$  denote the submatrix of rows of  $\mathbf{V}$  such that  $\theta_j = \mathbf{V}_j\omega$ . For any  $\mathbf{x} \in \{0, 1\}^n$  and  $j \in [n]$  let  $\mathbf{u}_j(\mathbf{x}) = y_j(\mathbf{x})\mathbf{V}_j^T \mathbf{z}(\mathbf{x})$ . Then each constraint can be re-written as  $y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle = \langle \mathbf{u}_j(\mathbf{x}), \omega \rangle$  and thus we can equivalently define the feasible set as

$$\mathcal{F}_\omega(\mathcal{S}) = \{ \omega \in \mathbb{R}^q : \langle \mathbf{u}_j(\mathbf{x}), \omega \rangle \geq 1 \quad \forall \mathbf{x} \in \mathcal{S}, \forall j \in [n] \} \quad (6)$$

Inspecting (6), clearly strict memorization of a dataset is equivalent to solving a linear program (LP) and therefore any algorithm which successfully memorizes  $\mathcal{S}$  is implicitly solving an LP. Moreover, these algorithms may have an *implicit bias* towards feasible points or solutions which satisfy other conditions or criteria. A popular and well studied example is the feasible point with the smallest norm: identifying this requires solving a quadratic program (QP) or, more specifically, a hard margin support vector machine (HSVM) problem. In

particular, note that if  $\boldsymbol{\theta} = \mathbf{V}\boldsymbol{\omega}$  where  $\boldsymbol{\omega} = [\sqrt{2}\mathbf{a}, \mathbf{b}]$  then  $\|\boldsymbol{\theta}\|^2 = \|\mathbf{W}\|_F^2 + \|\mathbf{b}\|^2 = \|\sqrt{2}\mathbf{a}\|^2 + \|\mathbf{b}\|^2 = \|\boldsymbol{\omega}\|^2$ . As a result, finding the minimum norm feasible point for a set  $\mathcal{S} \subset \{0, 1\}^n$  is equivalent to solving

$$\text{HSVM}(\mathcal{S}) = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}^q} \|\boldsymbol{\omega}\|^2 \quad \text{s.t.} \quad \boldsymbol{\omega} \in \mathcal{F}_{\boldsymbol{\omega}}(\mathcal{S}). \quad (7)$$

The key takeaway of this section is that minimizing (3) with gradient descent (GD) is implicitly biased in direction towards the solution of (7), i.e., norm-minimization. To this end, first observe that (3) can be re-parameterized as

$$L(\boldsymbol{\theta}; \mathcal{S}) = \sum_{\mathbf{x} \in \mathcal{S}} \sum_{j=1}^n \exp\left(E(\mathbf{x}; \boldsymbol{\theta}) - E(\mathbf{x}^{(j)}; \boldsymbol{\theta})\right) = \sum_{\mathbf{x} \in \mathcal{S}} \sum_{j=1}^n \exp(-\langle \mathbf{u}_j(\mathbf{x}), \boldsymbol{\omega} \rangle) =: L(\boldsymbol{\omega}; \mathcal{S}).$$

Consider now updates to the parameters of the Hopfield network using GD: in particular, given initial parameters  $\boldsymbol{\omega}(0)$  and step-size  $\eta > 0$ , for all  $t \geq 0$  let

$$\boldsymbol{\omega}(t+1) = \boldsymbol{\omega}(t) + \eta \sum_{\mathbf{x} \in \mathcal{S}} \sum_{j=1}^n \exp(-\mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega}(t)) \mathbf{u}_j(\mathbf{x}). \quad (8)$$

Using (Soudry et al., 2018, Thm.3), it can be proved that this sequence of GD iterates converges in direction to the solution of (7). Indeed, (Soudry et al., 2018, Thm.3) directly implies the following in our setting.

**Corollary 3.1.** *((Soudry et al., 2018, Thm.3)) Assume  $\mathcal{S}$  can be strictly memorized, let  $\boldsymbol{\omega}^* = \text{HSVM}(\mathcal{S})$ ,  $\boldsymbol{\omega}(0) \in \mathbb{R}^q$  be arbitrary and  $\boldsymbol{\omega}(t)$  be generated for all  $t \in \mathbb{N}_{\geq 1}$  as per (8). There exists a choice of step size  $\eta$  such that  $\boldsymbol{\omega}(t) = \boldsymbol{\omega}^* \log(t) + \rho(t)$  for all  $t \in \mathbb{N}_{\geq 1}$ , where  $\rho(t)$  grows as  $\|\rho(t)\| = O(\log(\log(t)))$ . Moreover,  $\lim_{t \rightarrow \infty} \frac{\boldsymbol{\omega}(t)}{\|\boldsymbol{\omega}(t)\|} = \frac{\boldsymbol{\omega}^*}{\|\boldsymbol{\omega}^*\|}$ .*

Informally, Corollary 3.1 states that the solution returned by minimizing the energy flow with gradient descent (MEF-GD) after exponentially many iterations is a close approximation directionally to the solution returned by solving the HSVM problem (7). We now derive generalization bounds both for the HSVM solution and MEF with GD.

**Theorem 3.2.** *Let  $D \subset \{0, 1\}^n$  be a set which can be strictly memorized and assume  $\|\mathbf{x}\|_0 \leq m \in \mathbb{N}_{\geq 1}$  for all  $\mathbf{x} \in D$ . Let  $\mathcal{D}$  be a probability distribution on  $D$ , and consider a random sample  $\mathcal{S}_N = (\mathbf{x}_i)_{i \in [N]}$ , where  $\mathbf{x}_i \sim \mathcal{D}$  are mutually i.i.d. Let  $\hat{\boldsymbol{\omega}} = \text{HSVM}(\mathcal{S}_N)$ ,  $\boldsymbol{\omega}^* = \text{HSVM}(D)$ ,  $\boldsymbol{\omega}(0) \in \mathbb{R}^q$  be arbitrary and  $\boldsymbol{\omega}(t)$  be generated for all  $t \in \mathbb{N}_{\geq 1}$  as per (8),  $\delta, \epsilon \in (0, 1)$  and assume  $\mathbf{x} \sim \mathcal{D}$  is sampled independent of  $\mathcal{S}_N$ . If  $N \gtrsim \epsilon^{-2} n \|\boldsymbol{\omega}^*\|^2 m \log(1/\delta)$  then both*

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\hat{\boldsymbol{\omega}}) \neq \mathbf{x}) \leq \epsilon \quad \text{and} \quad \mathbb{P}(H(\mathbf{x}; \mathbf{V}\boldsymbol{\omega}(t)) \neq \mathbf{x}) = \tilde{O}\left(\frac{\sqrt{m}\|\boldsymbol{\omega}^*\|}{\log(t)}\right) + \epsilon$$

hold with probability at least  $1 - \delta$  over the sample  $\mathcal{S}_N$ .

Note, the  $\tilde{O}(\cdot)$  notation hides a factor of  $\log(\log(t))$  arising in (Soudry et al., 2018, Thm.5). To prove Theorem 3.2 we combine a vector contraction inequality (Maurer, 2016, Corollary 1) with Rademacher bounds, see e.g., (Mohri et al., 2018, Theorem 3.3), we refer the reader to Appendix B.1 for a full proof. It is worth emphasizing that Theorem 3.2 implies that any dataset  $D$  which can be strictly memorized, can be at least *nearly* strictly memorized using only a polynomial number of samples. In Section 4.3 we take preliminary steps towards relaxing this statement, i.e., from memorizing samples drawn from  $\mathcal{D}$  with high probability, to memorizing the set  $D$  itself with high probability. Finally, we comment that the MEF bound implies gradient descent may require exponentially many iterations to converge directionally to the max-margin solution. In practice we emphasize that this is a tail phenomenon: once the weights are approximately aligned, all points are classified with a significant margin and their loss contributions become exponentially small.

## 4 Storing isomorphism classes of graphs

### 4.1 Graph encoding

Let  $\mathcal{G}_v$  denote the set of all simple, undirected graphs on  $v \in \mathbb{N}$  vertices. Recall two graphs  $G = (V, E)$ ,  $G' = (V', E')$  are isomorphic, which we denote  $G \cong G'$ , if there exists a bijection  $\phi : V \rightarrow V'$  such that

$(\phi(\nu_1), \phi(\nu_2)) \in E'$  if and only if  $(\nu_1, \nu_2) \in E$ . We refer to such a  $\phi$  as an isomorphism between  $G$  and  $G'$ . Note when  $V = V'$ , as is the case here since  $V = V' = [v]$ , then  $\phi$  is a permutation. The *isomorphism class* of a graph  $G \in \mathcal{G}_v$  is defined as  $\mathcal{I}(G) := \{G' \in \mathcal{G}_v : G' \cong G\}$ . Let  $\mathcal{V}_2$  denote the set of unordered pairs of  $[v]$ ,  $n := |\mathcal{V}_2| = \binom{v}{2}$  and  $\text{Ind} : \mathcal{V}_2 \rightarrow [n]$  be a bijection which indexes the elements of  $\mathcal{V}_2$ .

**Definition 4.1** (Edge representation of a graph). Let  $\mathcal{E}_{rep} : \mathcal{G}_v \rightarrow \{0, 1\}^n$  be defined as follows: if  $G = ([v], E) \in \mathcal{G}_v$  and  $\mathbf{x} = \mathcal{E}_{rep}(G)$  then, for all  $j \in [n]$ ,  $x_j := \mathbb{1}(\text{Ind}^{-1}(j) \in E)$ . We refer to  $\mathbf{x}$  as the edge representation of  $G$ .

To be clear,  $\mathcal{E}_{rep}$  is a bijection which assigns each graph to a binary vector of dimension  $n$  whose support defines the vertex pairs present in the edge set of the graph in question. Any vertex permutation induces an edge permutation.

**Definition 4.2** (Edge permutation induced by a vertex permutation). Let  $\phi : [v] \rightarrow [v]$  be a permutation. The edge permutation  $\pi_\phi : [n] \rightarrow [n]$  induced by  $\phi$  is defined as follows: if  $\text{Ind}^{-1}(j) = (\nu_1, \nu_2)$  then  $\pi_\phi(j) = \text{Ind}((\phi(\nu_1), \phi(\nu_2)))$ . We denote the subset of these edge permutations as  $\Pi_n^\Phi$  and the corresponding subset of permutation matrices as  $\Phi_n$ .

We now claim the following: first  $\Phi_n$  is a subgroup of  $\mathcal{P}_n$ , second if two graphs  $G, G' \in \mathcal{G}_v$  are isomorphic then there is a vertex induced edge permutation which maps between their edge representations, and third, for any  $G \in \mathcal{G}_v$  we have  $\mathcal{E}_{rep}(\mathcal{I}(G)) = \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$ . For proofs of these claims we refer the reader to Appendix A.3. Two edges are said to be *adjacent* if they share a vertex in common: more specifically, if  $j, l \in [n]$  then  $j$  and  $l$  are *adjacent*, which we denote  $j \sim l$ , if and only if  $|\text{Ind}^{-1}(j) \cap \text{Ind}^{-1}(l)| = 1$ , otherwise we say  $j$  and  $l$  are not adjacent, which we denote  $j \not\sim l$ . We now identify a subset of edge permutations which are characterized by preserving edge adjacency.

**Definition 4.3** (Edge adjacency preserving permutation). A permutation  $\pi : [n] \rightarrow [n]$  *preserves edge adjacency* if  $\pi(j) \sim \pi(l)$  if and only if  $j \sim l$ . We denote such permutations as  $\Pi_n^\mathcal{Q}$  and the corresponding set of permutation matrices as  $\mathcal{Q}_n$ .

Similar to  $\Phi_n$ , this subset forms a subgroup of  $\mathcal{P}_n$ . Moreover  $\Phi_n$  is a subgroup of  $\mathcal{Q}_n$  and as a result  $\mathcal{E}_{rep}(\mathcal{I}(G)) \subset \text{Orb}(\mathcal{E}_{rep}(G), \mathcal{Q}_n)$ . Again we refer the reader to Appendix A.3 for further details. As a result, the edge representations of the isomorphism class of a graph are a subset of the orbit of the edge representation of the graph in question under edge adjacency preserving permutations.

## 4.2 Experiments on graph data

To experimentally assess storage across isomorphism classes we study several classes of graphs; namely clique, bipartite, Paley and Johnson graphs. These families are standard examples in graph theory and algebraic graph theory, see for example Godsil & Royle (2001). Clique graphs, or more specifically  $k$ -cliques, contain a fully connected subset of  $k$  vertices while the remaining  $v - k$  vertices are isolated. Bipartite graphs split the  $v$  vertices into two equally sized groups with all possible inter-group edges present and no intra-group edges. Paley graphs connect vertices  $l$  and  $j$  when  $(l - j)$  is a quadratic residue mod  $v$ , as per NetworkX Developers; cf. Bollobás (2001). For integers  $n, r$ , the Johnson graph  $J(n, r)$  has as vertices the  $r$ -element subsets of  $[n]$ , with two vertices adjacent when the corresponding subsets differ in exactly one element<sup>2</sup>. These graphs are natural test cases because they are regular, highly symmetric and distance-regular while still exhibiting a nontrivial combinatorial structure that is distinct from cliques, bipartite graphs and Paley graphs. We remark that extensive experiments for cliques are already provided in Hillar & Tran (2018); Hillar et al. (2021); we include them here again for comparison and completeness. We also remark that we selected these classes due to the ease with which we are able to sample from and enumerate them, however, we emphasize that these families are representative rather than special. Indeed, we observe similar behavior for many other graph isomorphism classes, including random graphs. We refer the reader to our reproducibility statement which includes a link to our code. In addition, further experiments and preliminary results on topics including the Hidden Clique Problem (Fig. 6), comparison versus Dense Associative Memories (DAMs) (Fig. 7), orbit generalization for other graph families (Figs 8, 9), and double descent phenomena (Fig. 10) can be found in Appendix C.

<sup>2</sup>Equivalently, when their intersection has size  $(r - 1)$ .

Fig. 1 shows test accuracy versus training sample size, with mean and min-max over 10 trials, for Hopfield networks trained by MEF, Perceptron, and Delta (the latter two used only as baselines; see Appendix A.2). For small graphs ( $v = 8$ ) we enumerate the full isomorphism class and report the true accuracy, i.e., the fraction of the class memorized. For larger graphs ( $v = 20$ ), accuracy is estimated on an independent random sample of 1000 graphs. We highlight two observations: (i) MEF appears to reach higher test accuracy with fewer samples relative to the other methods, despite all methods perfectly memorizing the training set. This suggests differing implicit biases or implicit bias strengths. (ii) For MEF and Delta, the sample size needed to memorize an isomorphism class is tiny relative to the class size, aligning with the findings in Hillar & Tran (2018); Hillar et al. (2021). Furthermore, the number of iterations was capped at 1000, suggesting that the exponential dependency in Theorem 3.2 is highly pessimistic. We note that analogous experiments for other graph types are provided in Appendix C.

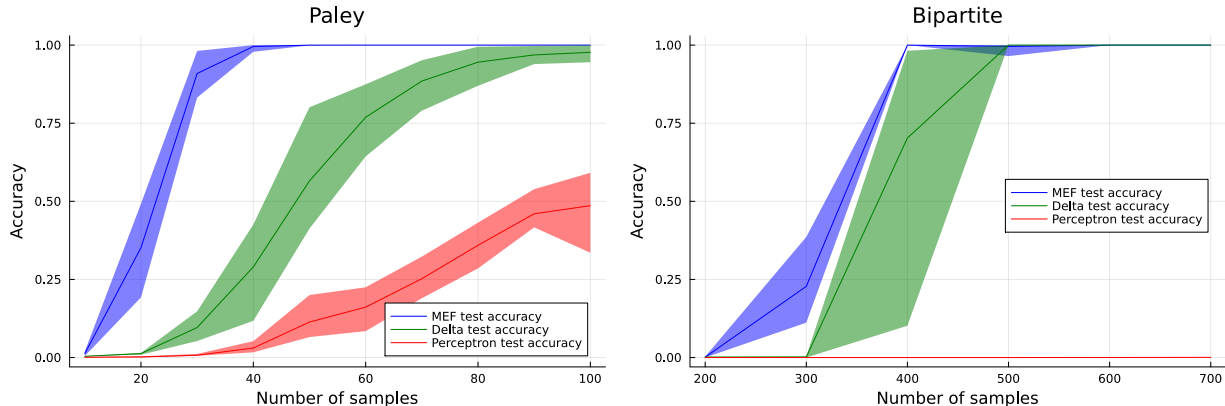


Figure 1: **Test accuracy vs. training sample size for isomorphism classes at two scales.** Top row:  $v = 8$ , isomorphism class size for Paley is 2520. Bottom row:  $v = 20$ , isomorphism class size for bipartite is 92,378. Curves show mean and min-max over 10 trials. Networks are trained with Perceptron, Delta (MSE), and MEF learning rules.

Fig. 2 estimates and compares the specific polynomial sample complexity of learning  $k$ -cliques versus Paley graphs. We do this in order to highlight that different isomorphism classes may be harder or easier to learn depending on their connectivity structure. For each graph size  $v$  we record  $s_{50}$ , which we define as the smallest training sample size for which MEF attains  $\geq 50\%$  average test accuracy on test samples of size 1000 averaged over 10 trials. The left subplot shows  $s_{50}$  vs.  $v$ . The right subplot shows a log-log fit. Assuming  $s_{50} = Cv^p$  for some constant  $C \in \mathbb{R}_{>0}$ , this allows us to estimate  $p$  via linear regression. While Paley graphs need  $N = \tilde{\Omega}(v^{2+\epsilon})$ , cliques require  $N = \tilde{\Omega}(v^{1+\epsilon})$  (note here we use  $\epsilon \in (0, 1)$  to denote a small error term). Thus, although Hopfield networks can memorize all classes (Lemma 4.7), the specific sample complexity appears to vary with graph connectivity.

Fig. 3 shows histograms of the elements of weight matrices of Hopfield networks trained on isomorphism classes of graphs, namely Bipartite and Johnson graphs, for three different sample sizes. Indeed, it is apparent that as the sample size grows the parameters returned by the optimizer converge onto a distinct subspace: we identify this subspace as the parameters invariant to the underlying group action of the data in Lemma 4.6 below. Heatmaps further supporting this conclusion can be observed in Figure 5 in Appendix C.

### 4.3 Invariant parameters

In what follows let  $\Gamma_n$  denote an arbitrary subgroup of  $\mathcal{P}_n$ . For any  $Q \in \mathcal{P}_n$  and  $\theta = (\mathbf{W}, \mathbf{b}) \in \Theta$  recall that we define the action of  $Q$  on the parameter  $\theta$  as  $Q\theta := (Q^T \mathbf{W} Q, Q^T \mathbf{b})$ . Let  $Q \in \mathcal{P}_n$ ,  $\theta = (\mathbf{W}, \mathbf{b}) \in \Theta$  and  $\theta' = Q\theta = (\mathbf{W}', \mathbf{b}')$ . Note  $\mathbf{W}'^T = (Q^T \mathbf{W} Q)^T = Q^T \mathbf{W} Q = \mathbf{W}'$  and for all  $j \in [n]$  we have  $W'_{jj} = e_{\pi(j)}^T \mathbf{W} e_{\pi(j)} = W_{\pi(j)\pi(j)} = 0$ . As a result  $\mathbf{W}' \in \text{Sym}_0^n$ , in addition trivially  $\mathbf{b}' \in \mathbb{R}^n$  and therefore  $\theta' \in \Theta$ . As a result, the action of  $\mathcal{P}_n$ , or any subgroup  $\Gamma_n$  of  $\mathcal{P}_n$ , on  $\Theta$  is well defined.

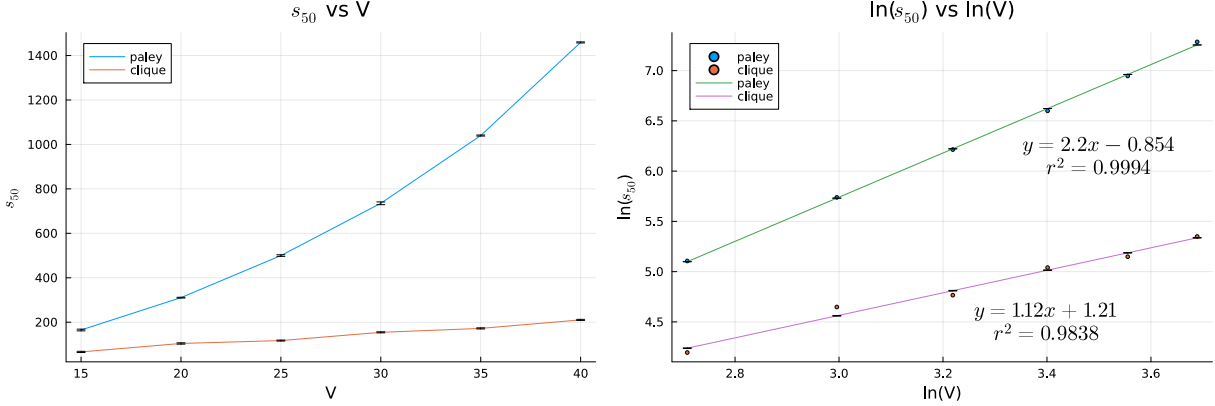


Figure 2: **Sample complexity scaling.** Plots showing the number of samples  $s_{50}$  required for a Hopfield network trained via MEF using accelerated gradient descent to memorize 50% of a random sample of 1000 graphs drawn from clique and Paley graph isomorphism classes on  $v$  vertices. On the right we plot  $\ln(s_{50})$  vs  $\ln(v)$  and compute the lines of best fit.

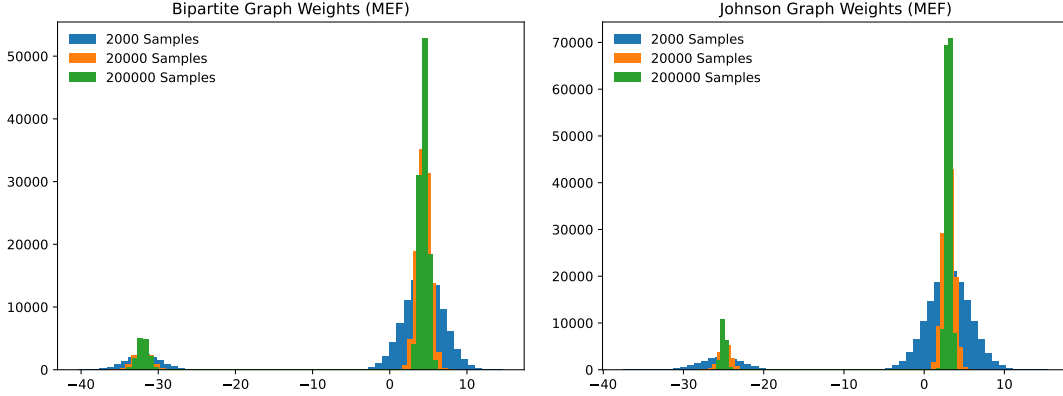


Figure 3: **Learned weights approach invariant subspace as sample size increases.** Histograms of MEF-learned weights over 3 sample counts. Top row: Bipartite graph with  $v = 32$ ,  $k = 16$  (496-bit). Bottom row: Johnson graph  $J(7, 3)$  on  $v = 35$  vertices (595-bit). Parameters are normalized so that thresholds (not shown) have mean absolute value 1.

**Definition 4.4** (Parameter invariance to the action of a subgroup). A parameter  $\theta \in \Theta$  is invariant with respect to  $\Gamma_n$  iff for all  $Q \in \Gamma_n$  then  $Q\theta = \theta$ . We denote the set of these parameters as  $\Psi(\Gamma_n)$ .

Recall  $E(Qx; \theta) = \frac{1}{2}x^T(Q^T W Q)x + (Q^T b)^T x$ . If  $\theta \in \Psi(\Gamma_n)$ , then for any  $x \in \{0, 1\}^n$  and  $Q \in \Gamma_n$  we have

$$E(Qx; \theta) = E(x; Q\theta) = E(x; \theta). \quad (9)$$

We refer to (9) as the intertwining property of the energy function. Using this property, the following lemma extends energy difference bounds between neighbors from a point to an orbit.

**Lemma 4.5.** *Let  $x_0 \in \{0, 1\}^n$  and  $\theta \in \Psi(\Gamma_n)$ . For  $\delta \in \mathbb{R}$ , if  $E(x_0^{(j)}; \theta) - E(x_0; \theta) \geq 1 - \delta$  for all  $j \in [n]$ , then for all  $x \in \text{Orb}(x_0, \Gamma_n)$  it follows that  $E(x^{(j)}; \theta) - E(x; \theta) \geq 1 - \delta$  for all  $j \in [n]$ .*

For a proof of this lemma, as well as the other results presented in this section, we refer the reader to Appendix B.2. A key implication of Lemma 4.5 is that if  $\theta \in \Psi(\Gamma_n)$  strictly memorizes  $x_0 \in \{0, 1\}^n$ , then  $\theta$  also strictly memorizes  $\text{Orb}(x_0, \Gamma_n)$ . We now show that invariance with respect to edge adjacency preserving permutations, of which graph isomorphisms are a subset, corresponds to a particular rank three subspace of the parameters.

**Lemma 4.6.** *Let  $F : \mathbb{R}^3 \rightarrow \Theta$  denote the linear map defined as follows: if  $(\mathbf{W}, \mathbf{b}) = F(\beta)$  then for all  $i, j \in [n]$  we have  $w_{ij} = 0$  if  $i = j$ ,  $w_{ij} = \beta_1$  if  $i \sim j$ ,  $w_{ij} = \beta_2$  if  $i \approx j$  and  $b_j = \beta_3$ . Then  $\Psi(\mathcal{Q}_n) = F(\mathbb{R}^3)$  where  $F(\mathbb{R}^3)$  denotes the image of  $F$ .*

By inspection, the parameter patterns observed in Figs. 3, 5, 6b, 9bd for MEF appear to approximately lie on the invariant subspace identified in Lemma 4.6. This suggests, given a sufficiently large training sample, that there is an implicit bias not just towards small norm solutions, but also those that are at least approximately invariant. Following this observation, a natural question to ask is whether or not parameters lying on this subspace can memorize any graph.

**Lemma 4.7.** *For  $m \in [0, n]$ , let  $\beta = [2, 2, 1 - 2m] \in \mathbb{R}^3$  and  $\theta = F(\beta)$ . Then  $E(\mathbf{x}^{(j)}; \theta) - E(\mathbf{x}; \theta) \geq 1$  for all  $j \in [n]$  and for all  $\mathbf{x} \in \{0, 1\}^n$  satisfying  $\|\mathbf{x}\|_0 = m$ .*

Combining Lemmas 4.7 and 4.5 we conclude that any graph isomorphism class can be strictly memorized by a Hopfield network. We also note that the only statistic used by the construction in Lemma 4.7 is the sparsity of the representation: in fact, this construction memorizes  $\mathbf{x} \in \{0, 1\}^n$  iff  $\|\mathbf{x}\|_0 = m$ . As a result, if our goal is to memorize an isomorphism class while avoiding spurious memories this is a poor parameter candidate. In addition, assuming  $m = \Theta(n)$ , the norm of this construction grows as  $\Theta(v^2)$ . For specific graph isomorphism classes solutions with a far smaller norm exist. As an example we consider  $k$ -cliques: for typographical ease we denote the set of binary representations of  $k$ -cliques on  $v$  vertices as  $\mathcal{C}_{v,k}$ .

**Lemma 4.8.** *If  $\beta = [-10/k, 38/k^2, 0] \in \mathbb{R}^3$ ,  $\theta = F(\beta)$ , and  $k \geq 5$ , then the following hold.*

1.  $E(\mathbf{x}^{(j)}; \theta) - E(\mathbf{x}; \theta) \geq 1$  for all  $\mathbf{x} \in \mathcal{C}_{v,k}$  and all  $j \in [n]$ .
2. If  $k = cv$  for some constant  $c \in (0, 1]$ , then there exists a constant  $C > 0$  such that  $\|\theta\|^2 \leq Cv$ .

Lemma 4.8 shows that a parameter exists which strictly memorizes  $\mathcal{C}_{v,k}$  with norm only  $O(\sqrt{v})$  rather than  $\Theta(v^2)$ . The construction in Lemma 4.8 also does not memorize all  $m$ -sparse binary vectors. For example, fixing some  $j \in [n]$ , suppose  $\mathbf{x}$  is such that  $\|\mathbf{x}\|_0 = m \leq 2(v - 2)$  and for all  $l \in \text{supp}(\mathbf{x})$  we have  $l \sim j$ . Then  $\mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega} = -10m/k$  and as a result  $\mathbf{x}$  is not strictly memorized. We speculate that perhaps a correlation between the size of the norm and the number of spurious memories exists, but we leave a proper investigation to future work.

Before proceeding we pause to reflect on the implications of our results with respect to (Hillar et al., 2021, Conjecture 1). Under a linear density regime  $k = cv$  where  $c \in (0, 1)$  is a constant, then together Theorem 3.2, Lemma 4.7 and Lemma 4.8 imply that  $k$ -cliques on  $v$  vertices can be strictly memorized with high probability as long as  $N = \tilde{\Omega}(v^5)$ . For simplicity assuming  $cv$  is an integer, then using Stirling’s approximation the critical ratio satisfies  $\tilde{O}(v^5)/\binom{v}{cv} = \tilde{O}(2^{-vH(c)}v^{5.5})$ , where here at the risk of confusion we use  $H$  to denote the binary entropy function. Clearly the critical ratio decays to zero at a rate which as  $v \rightarrow \infty$  is dominated by the exponential term.

Our experiments and results thus far suggest that memorization of a graph isomorphism class occurs when the training sample is sufficiently large that the optimizer is forced to return a solution lying close to the invariant set  $\Psi(\Phi_n) \subset F(\mathbb{R}^3)$ . The following lemma establishes that the HSVM solution on the full isomorphism class, which as  $N \rightarrow \infty$  is equivalent to the training sample with probability one, must be graph isomorphism invariant.

**Lemma 4.9.** *Let  $\mathbf{x}_0 \in \{0, 1\}^n$  and  $\Gamma_n$  denote a subgroup of  $\mathcal{P}_n$  and assume  $\text{Orb}(\mathbf{x}_0, \Gamma_n)$  can be strictly memorized. If  $\theta^* = \mathbf{V}\boldsymbol{\omega}^*$  where  $\boldsymbol{\omega}^* = \text{HSVM}_{\Theta}(\text{Orb}(\mathbf{x}_0, \Gamma_n))$  then  $\theta^* \in \Psi(\Gamma_n)$ .*

Following Lemma 4.9, we ask how many samples do we require in order to achieve at least approximate invariance? Deriving a sample complexity result is challenging, primarily due to the fact that the feasible set of the HSVM problem changes non-smoothly with respect to the training sample. Instead, and in order to gain intuition, we conclude this section by analyzing a related but simpler problem, which we refer to as the

average hard-margin support vector machine (AHSVM) problem. To this end, we define the following,

$$\mathcal{F}_A(\mathcal{S}) = \{\boldsymbol{\omega} \in \mathbb{R}^q : \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \langle \bar{\mathbf{u}}(\mathbf{x}), \boldsymbol{\omega} \rangle \geq 1\},$$

$$\text{AHSVM}(\mathcal{S}) = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}^q} \frac{1}{2} \|\boldsymbol{\omega}\|^2 \text{ s.t. } \boldsymbol{\omega} \in \mathcal{F}_A(\mathcal{S}).$$

The following lemma bounds the difference between the sample AHSVM solution and the population AHSVM solution in the context of a uniform distribution over an arbitrary  $\mathcal{O} \subseteq \{0, 1\}^n$ .

**Lemma 4.10.** *Let  $\mathcal{O} \subseteq \{0, 1\}^n$  satisfy  $\|\mathbf{x}\|_0 \leq m \in \mathbb{N}_{\geq 2}$  for all  $\mathbf{x} \in \mathcal{O}$  and assume  $\boldsymbol{\omega}^* = \text{HSVM}(\mathcal{O})$  is feasible. Consider a random sample  $\mathcal{S} = (\mathbf{x}_i)_{i=1}^N$  where  $\mathbf{x}_i \sim U(\mathcal{O})$  are mutually i.i.d. and define  $\boldsymbol{\omega}_{\mathcal{O}} = \text{AHSVM}(\mathcal{O})$  and  $\boldsymbol{\omega}_{\mathcal{S}} = \text{AHSVM}(\mathcal{S})$ . For  $\delta \in (0, 1]$  and  $\epsilon \in \mathbb{R}_{>0}$ , if  $N \gtrsim \epsilon^{-2} \|\boldsymbol{\omega}^*\|^4 m \log(1/\delta)$  then  $\|\boldsymbol{\omega}_{\mathcal{S}} - \boldsymbol{\omega}_{\mathcal{O}}\| \leq \epsilon$  with probability at least  $1 - \delta$ .*

Now let  $\text{Proj}_{\Psi(\Phi_n)}^\perp(\boldsymbol{\theta})$  denote the projection onto the subspace orthogonal to  $\Psi(\Phi_n)$ . Together, Lemmas 4.8, 4.10 and B.6 characterize proximity of the AHSVM solution for a  $k$ -clique sample to the invariant subspace.

**Corollary 4.11.** *Assume  $k = cv \geq 5$  for some constant  $c \in (0, 1)$  and let  $\mathcal{S}_N = (\mathbf{x}_i)_{i \in [N]}$ , where  $\mathbf{x}_i \sim U(\mathcal{C}_{v,k})$  are mutually i.i.d. Let  $\boldsymbol{\omega} = \text{AHSVM}(\mathcal{S}_N)$  and  $\boldsymbol{\theta} = \mathbf{V}\boldsymbol{\omega}$ . For  $\delta \in (0, 1)$ , if  $N \gtrsim \epsilon^{-2} v^4 \log(1/\delta)$  then  $\|\text{Proj}_{\Psi(\Phi_n)}^\perp(\boldsymbol{\omega})\| \leq \epsilon$  with probability at least  $1 - \delta$ .*

Corollary 4.11 illustrates that, at least for the AHSVM problem, we can get arbitrarily close to the invariant subspace with high probability using a sample size cubic in the number of vertices. We emphasize that even in the AHSVM setting bounding the distance from the learned parameters to the invariant subspace is challenging. We leave further refinement of these results as well as a derivation of an analogous one for the HSVM problem to future work.

## 5 Conclusions, Limitations and Future Work

In summary, this work shows that classical Hopfield networks can exploit symmetry in graph-structured data. By reformulating strict memorization as a linear feasibility problem, we identify norm efficiency as a mechanism linking MEF training, max-margin solutions, and sample-efficient orbit generalization. For graph isomorphism classes, the invariant parameter subspace provides a compact explanation for the observed few-shot-to-orbit phenomenon: as training samples increasingly constrain the feasible set, low-norm solutions are driven toward approximately invariant representations. These results suggest that even in minimal associative-memory models, implicit bias can substitute in part for architectural invariance.

This work has limits: while the full-orbit HSVM solution is invariant and, since the orbit is finite, an i.i.d. sample eventually contains the entire orbit almost surely, we do not yet provide a quantitative finite-sample convergence rate showing how quickly the HSVM or MEF solutions approach the invariant subspace. We also do not yet explain why some isomorphism classes appear to be easier to learn than others. Future work should quantify spurious fixed points and basin robustness, treat other subgroups and unions of orbits, handle noisy or non-uniform group data, extend the analysis to hypergraphs, and involve continuous and modern HNs. Towards achieving these goals, we highlight preliminary experimental findings detailed in Appendix C concerning invariant subspace convergence (Fig. 5), the Hidden Clique Problem (Fig. 6), comparison to DAMs (Fig. 7), generalization in other graph families (Figs. 8, 9), double descent phenomena (Fig. 10), and isomorphic graph checking (Algorithm 1).

**Reproducibility Statement:** Code able to reproduce our experimental results can be found in the following anonymous repository, <https://github.com/hopnetorbit/HopfieldNetworksIsomorphism>.

**Ethics Statement:** This work uses only synthetic, non-sensitive data, involves no human or animal subjects, carries minimal dual-use or environmental risks (modest compute).

**Broader Impact Statement:** This work contributes to the mathematical and conceptual foundations of machine learning and intelligent systems. Although the results are primarily theoretical, a better understanding of the underlying principles governing learning, generalization, and representation can inform the long-term development of more reliable, interpretable, and robust algorithms.

## References

- Shun-ichi Amari. Learning patterns and pattern sequences by self-organizing nets of threshold elements. *IEEE Transactions on computers*, 100(11):1197–1206, 1972.
- Shun-Ichi Amari. Characteristics of sparsely encoded associative memory. *Neural networks*, 2(6):451–457, 1989.
- Daniel J Amit, Hanoach Gutfreund, and Haim Sompolinsky. *Storing infinite numbers of patterns in a spin-glass model of neural networks*, volume 55. APS, 1985.
- Béla Bollobás. *Explicit Constructions*, pp. 348–382. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2001.
- J. Bruck. On the convergence properties of the Hopfield model. *Proceedings of the IEEE*, 78(10):1579–1585, 1990. doi: 10.1109/5.58341.
- Shuxiao Chen, Edgar Dobriban, and Jane H. Lee. A group-theoretic framework for data augmentation. *Journal of Machine Learning Research*, 21(245):1–71, 2020. URL <https://jmlr.org/papers/v21/20-163.html>.
- Thomas M. Cover. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. *IEEE Transactions on Electronic Computers*, EC-14(3):326–334, 1965. doi: 10.1109/PGEC.1965.264137.
- Tri Dao, Albert Gu, Alexander Ratner, Virginia Smith, Christopher De Sa, and Christopher Ré. A kernel theory of modern data augmentation. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, volume 97 of *PMLR*, pp. 1528–1537, 2019. URL <https://proceedings.mlr.press/v97/dao19b.html>.
- Yael Dekel, Ori Gurel-Gurevich, and Yuval Peres. Finding hidden cliques in linear time with high probability. *Combinatorics, Probability and Computing*, 23(1):29–49, 2014.
- Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Upgang, and Franck Vermet. On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168(2):288–299, 2017.
- NetworkX Developers. URL [https://networkx.org/documentation/stable/reference/generated/networkx.generators.expanders.paley\\_graph.html#networkx.generators.expanders.paley\\_graph](https://networkx.org/documentation/stable/reference/generated/networkx.generators.expanders.paley_graph.html#networkx.generators.expanders.paley_graph). Python package for the creation, manipulation, and study of the structure, dynamics, and functions of complex networks.
- Bryn Elesedy and Sheheryar Zaidi. Provably strict generalisation benefit for invariance in kernel ridge regression. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/8fe04df45a22b63156ebabbb064fcd5e-Abstract.html>.
- JF Fontanari. Generalization in a Hopfield network. *Journal de Physique*, 51(21):2421–2430, 1990.
- Elizabeth Gardner. The space of interactions in neural network models. *Journal of Physics A: Mathematical and General*, 21(1):257, 1988.
- Jan E. Gerken and Pan Kessel. Emergent equivariance in deep ensembles. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, volume 235 of *PMLR*, pp. 15438–15465, 2024. URL <https://proceedings.mlr.press/v235/gerken24a.html>.
- Chris Godsil and Gordon F. Royle. *Algebraic Graph Theory*. Number Book 207 in Graduate Texts in Mathematics. Springer, 2001. ISBN 9781461301639 1461301637. doi: 10.1007/978-1-4613-0163-9.
- Donald O. Hebb. *The Organization of Behavior: A Neuropsychological Theory*. Wiley, New York, 1949.
- JA Hertz, A Krogh, and RG Palmer. *Introduction to the theory of neural computation*. Addison-Wesley, 1991.

- Christopher Hillar, Jascha Sohl-Dickstein, and Kilian Koepsell. Efficient and optimal binary Hopfield associative memory storage using minimum probability flow. In *4th Neural Information Processing Systems (NeurIPS) Workshop on Discrete Optimization in Machine Learning (DISCML): structure and scalability*, pp. 1–6, 2012.
- Christopher Hillar, Tenzin Chan, Rachel Taubman, and David Rolnick. Hidden hypergraphs, error-correcting codes, and critical learning in Hopfield networks. *Entropy*, 23(11), 2021.
- Christopher J Hillar and Ngoc M Tran. Robust exponential memory in Hopfield networks. *The Journal of Mathematical Neuroscience*, 8:1–20, 2018.
- ME Hoff and B Widrow. Adaptive switching circuits. In *1960 IRE WESCON Convention Record, Part 4*, pp. 96–104. IRE New York, NY, USA, 1960.
- John J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, 1982. doi: 10.1073/pnas.79.8.2554.
- David Horn and Marius Usher. Capacities of multiconnected memory models. *Journal de Physique*, 49(3): 389–395, 1988.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018. URL <https://arxiv.org/abs/1803.07300>.
- Silvio Kalaj, Clarissa Lauditi, Gabriele Perugini, Carlo Lucibello, Enrico M. Malatesta, and Matteo Negri. Random features hopfield networks generalize retrieval to previously unseen examples. *Physica A: Statistical Mechanics and its Applications*, 678:130946, 2025. ISSN 0378-4371. doi: <https://doi.org/10.1016/j.physa.2025.130946>. URL <https://www.sciencedirect.com/science/article/pii/S0378437125005989>.
- Werner Krauth and Marc Mézard. Storage capacity of memory networks with binary couplings. *Journal de Physique*, 50(20):3057–3066, 1989.
- Dmitry Krotov and John J Hopfield. Dense associative memory for pattern recognition. *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 2016.
- William A Little. The existence of persistent states in the brain. *Mathematical biosciences*, 19(1-2):101–120, 1974.
- Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical programming*, 45(1):503–528, 1989.
- Giovanni Luca Marchetti, Christopher J Hillar, Danica Kragic, and Sophia Sanborn. Harmonics of learning: Universal fourier features emerge in invariant networks. In Shipra Agrawal and Aaron Roth (eds.), *Proceedings of Thirty Seventh Conference on Learning Theory*, volume 247 of *Proceedings of Machine Learning Research*, pp. 3775–3797. PMLR, 30 Jun–03 Jul 2024. URL <https://proceedings.mlr.press/v247/marchetti24a.html>.
- Andreas Maurer. A vector-contraction inequality for rademacher complexities. In Ronald Ortner, Hans Ulrich Simon, and Sandra Zilles (eds.), *Algorithmic Learning Theory*, pp. 3–17, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46379-7.
- Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- Robert J. McEliece, Edward C. Posner, Eugene R. Rodemich, and Santosh S. Venkatesh. The capacity of the Hopfield associative memory. *IEEE Trans. Inform. Theory*, 33(4):461–482, 1987. ISSN 0018-9448,1557-9654. doi: 10.1109/TIT.1987.1057328. URL <https://doi.org/10.1109/TIT.1987.1057328>.

- Philipp Misof, Pan Kessel, and Jan E Gerken. Equivariant neural tangent kernels. In Aarti Singh, Maryam Fazel, Daniel Hsu, Simon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri Wagstaff, and Jerry Zhu (eds.), *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 44470–44503. PMLR, 13–19 Jul 2025. URL <https://proceedings.mlr.press/v267/misof25a.html>.
- Mehryar Mohri, Afshin Rostamizadeh, and Amee Talwalkar. *Foundations of Machine Learning*. The MIT Press, 2nd edition, 2018. ISBN 0262039400.
- Mor Shpigel Nacson, Jason D. Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 89 of *Proceedings of Machine Learning Research*, pp. 3420–3428. PMLR, 2019. URL <https://proceedings.mlr.press/v89/nacson19b.html>.
- Kaoru Nakano. Associatron—a model of associative memory. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-2(3):380–388, 1972. doi: 10.1109/TSMC.1972.4309133.
- Matteo Negri, Clarissa Lauditi, Gabriele Perugini, Carlo Lucibello, and Enrico Malatesta. Storage and learning phase transitions in the random-features Hopfield model. *arXiv preprint arXiv:2303.16880*, 2023.
- Oskar Nordenfors, Fredrik Ohlsson, and Axel Flinth. Optimization dynamics of equivariant and augmented neural networks, 2024. URL <https://arxiv.org/abs/2303.13458>.
- Leonid A Pastur and Alexander L Figotin. Exactly soluble model of a spin glass. *Soviet Journal of Low Temperature Physics*, 3(6):378–383, 1977.
- L. Personnaz, I. Guyon, and G. Dreyfus. Information storage and retrieval in spin-glass like neural networks. *Journal de Physique Lettres*, 46(8):359–365, 1985.
- L. Personnaz, I. Guyon, and G. Dreyfus. Collective computational properties of neural networks: New learning mechanisms. *Phys. Rev. A*, 34:4217–4228, Nov 1986. doi: 10.1103/PhysRevA.34.4217. URL <https://link.aps.org/doi/10.1103/PhysRevA.34.4217>.
- Iosif Pinelis. Optimum bounds for the distributions of martingales in banach spaces. *The Annals of Probability*, 22(4):1679–1706, 1994. ISSN 00911798, 2168894X. URL <http://www.jstor.org/stable/2244912>.
- Hubert Ramsauer, Bernhard Schöfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Lukas Gruber, Markus Holzleitner, Thomas Adler, David Kreil, Michael K Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. Hopfield networks is all you need. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=tL89RnzIiCd>.
- Robert A. Rescorla and Allan R. Wagner. A theory of pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In Abraham H. Black and William F. Prokasy (eds.), *Classical Conditioning II: Current Research and Theory*, pp. 64–99. Appleton-Century-Crofts, New York, 1972.
- Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386–408, 1958.
- Jascha Sohl-Dickstein, Peter B. Battaglino, and Michael R. DeWeese. New method for parameter estimation in probabilistic models: Minimum probability flow. *Phys. Rev. Lett.*, 107:220601, Nov 2011. doi: 10.1103/PhysRevLett.107.220601. URL <https://link.aps.org/doi/10.1103/PhysRevLett.107.220601>.
- Daniel Soudry, Elad Hoffer, and Nathan Srebro. The implicit bias of gradient descent on separable data. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=r1q7n9gAb>.
- A.J. Storkey and R. Valabregue. The basins of attraction of a new hopfield learning rule. *Neural Networks*, 12(6):869–876, 1999. ISSN 0893-6080. doi: [https://doi.org/10.1016/S0893-6080\(99\)00038-6](https://doi.org/10.1016/S0893-6080(99)00038-6). URL <https://www.sciencedirect.com/science/article/pii/S0893608099000386>.

Amos Storkey. Increasing the capacity of a hopfield network without sacrificing functionality. In Wulfram Gerstner, Alain Germond, Martin Hasler, and Jean-Daniel Nicoud (eds.), *Artificial Neural Networks — ICANN'97*, pp. 451–456, Berlin, Heidelberg, 1997. Springer Berlin Heidelberg. ISBN 978-3-540-69620-9.

Pavel Tolmachev and Jonathan H. Manton. New insights on learning rules for Hopfield networks: Memory and objective function minimisation. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2020.

M. V. Tsodyks and M. V. Feigel'man. The enhanced storage capacity in neural networks with low activity level. *Europhysics Letters*, 6(2):101–105, May 1988. doi: 10.1209/0295-5075/6/2/002.

David J Willshaw, Oliver P Buneman, and HC Longuet-Higgins. Non-holographic associative memory. *Nature*, 222(5197):960–962, 1969.

## A Background

### A.1 Energy gap for binary vectors a hamming distance one apart

As discussed in Section 2, memorization of a point is equivalent to ensuring an energy gap between it and its neighbors. Recall we define  $\mathbf{x}^{(j)} \in \{0, 1\}^n$  as the vector that differs from  $\mathbf{x} \in \{0, 1\}^n$  only at the  $j$ th location, and  $\mathbf{z}(\mathbf{x}) = [\mathbf{x}, 1] \in \mathbb{R}^{n+1}$ .

**Lemma A.1.**  $E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) = y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \boldsymbol{\theta}_j \rangle$ .

*Proof.* By definition  $x_l^{(j)} \neq x_l$  iff  $l = j$ . In addition,  $x_j^{(j)} - x_j = 1 - 2x_j$  and, if  $r \neq l$ , then  $x_l x_r \neq x_l^{(j)} x_r^{(j)}$  iff either  $l = j$  and  $r \neq j$ , or  $l \neq j$  and  $r = j$ . Furthermore, recall  $\mathbf{W}$  is symmetric and  $W_{jj} = 0$  for all  $j \in [n]$ . As a result,

$$\begin{aligned}
 E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) &= \frac{1}{2} \sum_{l,r \in [n]} W_{rl} (x_r^{(j)} x_l^{(j)} - x_r x_l) + \sum_{l=1}^n b_l (x_l^{(j)} - x_l) \\
 &= \frac{1}{2} \sum_{r \in [n], r \neq j} W_{rj} (x_r^{(j)} x_j^{(j)} - x_r x_j) + \frac{1}{2} \sum_{l \in [n], l \neq j} W_{jl} (x_j^{(j)} x_l^{(j)} - x_j x_l) + b_j (1 - 2x_j) \\
 &= \sum_{l \in [n], l \neq j} W_{jl} (x_j^{(j)} x_l^{(j)} - x_j x_l) + b_j (1 - 2x_j) \\
 &= \sum_{l \in [n], l \neq j} W_{jl} (x_j^{(j)} - x_j) x_l + b_j (1 - 2x_j) \\
 &= (1 - 2x_j) \left( \sum_{l \in [n]} W_{jl} x_l + b_j \right) \\
 &= y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \boldsymbol{\theta}_j \rangle.
 \end{aligned}$$

as claimed. □

### A.2 Learning algorithms for Hopfield Networks

We briefly describe several classical learning rules Hertz et al. (1991) that can be applied to find parameters in Hopfield networks. These methods typically trade off between biological plausibility and performance. We remark that this list is far from exhaustive; see Tolmachev & Manton (2020) for a recent summary.

- **Outer-Product Rule** (Hebb, 1949; Nakano, 1972; Amari, 1972; Hopfield, 1982). In the attractor neural network case Hopfield (1982), this Hebbian rule constructs weights as the normalized sum of training pattern outer products. This rule is simple, biologically motivated, and local in nature, but it is often observed to suffer from spurious attractors, shallow basins of attraction, and overall limited capacity.
- **Perceptron Rule** (Rosenblatt, 1958). The difference between the desired response – that the network dynamics should fix a training sample – and the actual linear-threshold response of a neuron gives a learning signal to update parameters.
- **Delta** (Hoff & Widrow, 1960; Rescorla & Wagner, 1972). The delta rule, also called the Mean Squared Error (MSE) or Least Mean Square (LMS) rule, considers a relaxation and follows a gradient to minimize the squared error between the linear output activations of neurons and the training pattern to memorize.
- **Projection Rule** (Personnaz et al., 1985; 1986). The weight matrix is obtained by projecting onto the span of the training data and then zeroing the diagonal entries.

- **Storkey Rule** (Storkey, 1997; Storkey & Valabregue, 1999). A modification of the Hebbian update that reduces interference between patterns by accounting for previously stored ones. This rule achieves higher storage capacity than Hebbian learning and reduces spurious minima.

### A.3 Encoding simple, undirected graphs as binary vectors

First we recap some of our notation. Let  $\mathcal{G}_v$  denote the set of all simple, undirected graphs on  $v \in \mathbb{N}$  vertices. Recall two graphs  $G = (V, E)$ ,  $G' = (V', E')$  are isomorphic, which we denote  $G \cong G'$ , if there exists a bijection  $\phi : V \rightarrow V'$  such that  $(\phi(\nu_1), \phi(\nu_2)) \in E'$  if and only if  $(\nu_1, \nu_2) \in E$ . We refer to such a  $\phi$  as an isomorphism between  $G$  and  $G'$ . Furthermore,  $\phi$  is a permutation when  $V = V'$ : in our setting we consider  $V = V' = [v]$  and therefore we shall discuss only permutations moving forward. The *isomorphism class* of a graph  $G \in \mathcal{G}_v$  is defined as  $\mathcal{I}(G) := \{G' \in \mathcal{G}_v : G' \cong G\}$ . An automorphism of a graph  $G = (V, E)$  is a permutation  $\phi : V \rightarrow V$  such that  $(\nu_1, \nu_2) \in E$  implies  $(\phi(\nu_1), \phi(\nu_2)) \in E$ . In short, while an isomorphism preserves the vertex adjacency structure of a graph an automorphism preserves not just the vertex adjacency structure but also the vertex labels. Recall that  $\Phi_n \subset \mathcal{P}_n$  refers to the set of edge permutation matrices induced by permutations of the vertices, see Definition 4.2.

**Lemma A.2.**  $\Phi_n$  is a subgroup of  $\mathcal{P}_n$ .

*Proof.* Clearly this is equivalent to showing that  $\Pi_n^\Phi$  is a subgroup of  $\Pi_n$ . It is easy to check that  $I \in \Pi_n^\Phi$ ,  $\pi_\phi \in \Pi_n^\Phi$  implies  $\pi_\phi^{-1} \in \Pi_n^\Phi$  and  $\pi_\phi, \pi_{\phi'} \in \Pi_n^\Phi$  implies  $\pi_\phi \circ \pi_{\phi'} \in \Pi_n^\Phi$ , therefore  $\Pi_n^\Phi$  is a subgroup of  $\Pi_n$ .  $\square$

The following lemmas establish a straightforward equivalence between isomorphism classes of graphs and certain orbits of binary vectors. First, Lemma A.3 shows that if two graphs  $G, G' \in \mathcal{G}_v$  are isomorphic then there is a vertex induced edge permutation which maps between their edge representations.

**Lemma A.3.** Suppose  $G = ([v], E), G' = ([v], E') \in \mathcal{G}_v$  and  $\mathbf{x} = \mathcal{E}_{rep}(G), \mathbf{x}' = \mathcal{E}_{rep}(G')$ . Then  $G \cong G'$  iff there exists a  $\mathbf{P} \in \Phi_n$  such that  $\mathbf{P}\mathbf{x} = \mathbf{x}'$ .

*Proof.* Assume  $G \cong G'$ . Then there exists a permutation  $\phi : [v] \rightarrow [v]$  such that  $(\nu_1, \nu_2) \in E$  implies  $(\phi(\nu_1), \phi(\nu_2)) \in E'$ . Let  $\pi_\phi : [n] \rightarrow [n]$  be the edge permutation induced by  $\phi$  and  $\mathbf{P} \in \Phi_n$  the corresponding permutation matrix. By construction  $x_j = x'_{\pi_\phi(j)}$  for all  $j \in [n]$ , equivalently, if  $\mathbf{P} \in \Phi_n$  is the permutation matrix associated with  $\pi_\phi$  then  $\mathbf{P}\mathbf{x} = \mathbf{x}'$ . Now suppose there exists a  $\mathbf{P} \in \Phi_n$  such that  $\mathbf{P}\mathbf{x} = \mathbf{x}'$ . Then there exists a vertex permutation  $\phi : [v] \rightarrow [v]$  which induces an edge permutation  $\pi_\phi : [n] \rightarrow [n]$  such that  $x_j = x'_{\pi_\phi(j)}$ . By construction, if  $j = \text{Ind}((\nu_1, \nu_2))$  then this implies  $\pi_\phi(j) = \text{Ind}((\phi(\nu_1), \phi(\nu_2)))$ . Therefore, by the definition of  $\mathcal{E}_{rep}$  we have  $(\phi(\nu_1), \phi(\nu_2)) \in E'$  iff  $(\nu_1, \nu_2) \in E$ . Therefore  $\phi$  is an isomorphism between  $G$  and  $G'$  and  $G \cong G'$ .  $\square$

Building on Lemma A.3, the following lemma characterizes the isomorphism class of a graph in terms of the orbit of its edge representation under vertex induced edge permutations.

**Lemma A.4.** For any  $G \in \mathcal{G}_v$  we have  $\mathcal{E}_{rep}(\mathcal{I}(G)) = \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$ .

*Proof.* Let  $\mathbf{x} = \mathcal{E}_{rep}(G)$ , then

$$\text{Orb}(\mathcal{E}_{rep}(G), \Phi_n) = \{\mathbf{P}\mathbf{x} : \mathbf{P} \in \Phi_n\}.$$

Suppose  $\mathcal{E}_{rep}(\mathcal{I}(G)) \not\subset \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$ . Then there exists a  $G' \in \mathcal{I}(G)$  such that  $\mathbf{x}' := \mathcal{E}_{rep}(G') \notin \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$ . Therefore there does not exist a  $\mathbf{P} \in \Phi_n$  such that  $\mathbf{P}\mathbf{x} = \mathbf{x}'$ . However,  $G \cong G'$  which implies a contradiction by Lemma A.3, therefore  $\mathcal{E}_{rep}(\mathcal{I}(G)) \subset \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$ . Now suppose  $\text{Orb}(\mathcal{E}_{rep}(G), \Phi_n) \not\subset \mathcal{E}_{rep}(\mathcal{I}(G))$ , then there exists a  $\mathbf{x}' \in \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$  such that  $G' = \mathcal{E}_{rep}^{-1}(\mathbf{x}') \notin \mathcal{I}(G)$ . However, as  $\mathbf{x}' \in \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$  then there exists a  $\mathbf{P} \in \Phi_n$  such that  $\mathbf{P}\mathbf{x} = \mathbf{x}'$ , but using Lemma A.3 this implies  $G \cong G'$  which is a contradiction. Therefore  $\text{Orb}(\mathcal{E}_{rep}(G), \Phi_n) \subset \mathcal{E}_{rep}(\mathcal{I}(G))$ . Combining these two observations we conclude that  $\mathcal{E}_{rep}(\mathcal{I}(G)) = \text{Orb}(\mathcal{E}_{rep}(G), \Phi_n)$ .  $\square$

**Lemma A.5.**  $\mathcal{Q}_n$  is a subgroup of  $\mathcal{P}_n$ .

*Proof.* Trivially it suffices to show that  $\Pi_n^{\mathcal{Q}}$  is a subgroup of  $\Pi_n$ . It is easy to check that  $I \in \Pi_n^{\mathcal{Q}}$ ,  $\pi \in \Pi_n^{\mathcal{Q}}$  implies  $\pi^{-1} \in \Pi_n^{\mathcal{Q}}$  and  $\pi, \pi' \in \Pi_n^{\mathcal{Q}}$  implies  $\pi \circ \pi' \in \Pi_n^{\mathcal{Q}}$ . Therefore  $\Pi_n^{\mathcal{Q}}$  is a subgroup of  $\Pi_n$ .  $\square$

The following lemma states that the vertex induced edge permutations form a subgroup of the edge adjacency preserving subgroup of permutations. As a result, the edge representations of the isomorphism class of a graph are a subset of the orbit of the edge representation of the graph in question under edge adjacency preserving permutations.

**Lemma A.6.**  $\Phi_n$  is a subgroup of  $\mathcal{Q}_n$  and  $\mathcal{E}_{rep}(\mathcal{I}(G)) \subset \text{Orb}(\mathcal{E}_{rep}(G), \mathcal{Q}_n)$ .

*Proof.* Trivially it suffices to show that  $\Pi_n^{\phi}$  is a subgroup of  $\Pi_n^{\mathcal{Q}}$ , we proceed to show that any vertex induced permutation is an edge adjacency preserving permutation. Consider two edge indices  $i, j \in [n]$  and let  $\nu_1, \nu_2, \nu_3, \nu_4 \in [v]$  be distinct. Suppose  $i \sim j$ , then without loss of generality let  $\text{Ind}^{-1}(i) = (\nu_1, \nu_2)$  and  $\text{Ind}^{-1}(j) = (\nu_2, \nu_3)$ . Then  $\text{Ind}^{-1}(\pi_{\phi}(i)) = (\phi(\nu_1), \phi(\nu_2))$  and  $\text{Ind}^{-1}(\pi_{\phi}(j)) = (\phi(\nu_2), \phi(\nu_3))$ , therefore  $i \sim j$  implies  $\pi_{\phi}(i) \sim \pi_{\phi}(j)$ . Suppose now  $i \not\sim j$ , if  $i = j$  then trivially  $\pi_{\phi}(i) = \pi_{\phi}(j)$  and therefore  $i = j$  implies  $\pi_{\phi}(i) \approx \pi_{\phi}(i)$ . Otherwise, and again without loss of generality, let  $\text{Ind}^{-1}(i) = (\nu_1, \nu_2)$  and  $\text{Ind}^{-1}(j) = (\nu_3, \nu_4)$ . Then  $\text{Ind}^{-1}(\pi_{\phi}(i)) = (\phi(\nu_1), \phi(\nu_2))$  and  $\text{Ind}^{-1}(\pi_{\phi}(j)) = (\phi(\nu_3), \phi(\nu_4))$ , as  $\phi$  is bijection then this implies  $\pi_{\phi}(i) \not\approx \pi_{\phi}(j)$ . As a result,  $\pi_{\phi}(i) \sim \pi_{\phi}(i)$  if and only if  $i \sim j$ . Finally as  $\Phi_n$  is a group and it is a subset of  $\mathcal{Q}_n$  then it must be a subgroup of  $\mathcal{Q}_n$ . As a result  $\text{Orb}(\mathcal{E}_{rep}(G), \Phi_n) \subset \text{Orb}(\mathcal{E}_{rep}(G), \mathcal{Q}_n)$   $\square$

#### A.4 Bounded representations

In order to establish the connection between Hopfield networks and SVMs discussed in Section 3, we identified and defined a certain feature map for the inputs to the underlying linear classification problem. Recall there exists a matrix  $\mathbf{V} \in \{0, 1, 1/\sqrt{2}\}^{n(n+1) \times q}$  such that for any  $\boldsymbol{\theta} \in \Theta$  there exists a  $\mathbf{a} \in \mathbb{R}^p$ ,  $\boldsymbol{\omega} = [\sqrt{2}\mathbf{a}, \mathbf{b}]$  such that  $\boldsymbol{\theta} = \mathbf{V}\boldsymbol{\omega}$ . Recall also that we define  $\mathbf{V}_j \in \{0, 1, 1/\sqrt{2}\}^n$  as the matrix which satisfies  $\boldsymbol{\theta}_j = \mathbf{V}_j\boldsymbol{\omega}$ , where  $\boldsymbol{\theta}_j = [\mathbf{w}_j, b_j] \in \mathbb{R}^{n+1}$ . In addition, for any  $\mathbf{x} \in \{0, 1\}$  then we let  $\mathbf{z}(\mathbf{x}) = [\mathbf{x}, 1] \in \mathbb{R}^{n+1}$ ,  $\mathbf{u}_j(\mathbf{x}) = \mathbf{V}_j^T \mathbf{z}(\mathbf{x})$  for all  $j \in [n]$  and  $\bar{\mathbf{u}}(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n \mathbf{u}_j(\mathbf{x})$ . The following lemma bounds the norm of these representations.

**Lemma A.7.** For any  $\mathbf{x} \in \{0, 1\}^n$  then

$$\|\bar{\mathbf{u}}(\mathbf{x})\|^2 \leq \|\mathbf{u}_j(\mathbf{x})\|^2 = \frac{1}{2} (\|\mathbf{x}\|_0 - x_j) + 1$$

for all  $j \in [n]$ .

*Proof.* Let  $\delta_n(j) \in \{0, 1\}^n$  denote the one hot vector such that  $\text{supp}(\delta_n(j)) = j$ . In addition, let  $\phi_j : [q] \rightarrow [n]$  denote the injective mapping between the indices of  $\mathbf{a}$  and their respective positions in  $\mathbf{w}_j$ , and let  $\mathbf{B}_j \in \{0, 1\}^{n \times p}$  be the associated matrix which copies the elements of  $\mathbf{a}$  into their positions in  $\boldsymbol{\theta}_j$ . Therefore, we can write

$$\boldsymbol{\theta}_j = \begin{bmatrix} \mathbf{w}_j \\ b_j \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}}\mathbf{B}_j & \mathbf{0}_{n \times n} \\ \mathbf{0}_{1 \times n} & \delta_n(j)^T \end{bmatrix} \begin{bmatrix} \sqrt{2}\mathbf{a} \\ \mathbf{b} \end{bmatrix} = \mathbf{V}_j\boldsymbol{\omega}.$$

By definition

$$\mathbf{u}_j(\mathbf{x}) = \mathbf{V}_j^T \mathbf{z}(\mathbf{x}) = \begin{bmatrix} \frac{1}{\sqrt{2}}\mathbf{B}_j^T & \mathbf{0}_{n \times 1} \\ \mathbf{0}_{n \times n} & \delta_n(j) \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}}\mathbf{B}_j^T \mathbf{x} \\ \delta_n(j) \end{bmatrix},$$

therefore

$$\|\mathbf{u}_j(\mathbf{x})\|^2 = \frac{1}{2} \mathbf{x}^T \mathbf{B}_j \mathbf{B}_j^T \mathbf{x} + 1.$$

Recall  $W_{jj} = 0$  and each other element of  $\mathbf{w}_j$  corresponds to exactly one element of  $\mathbf{a}$ , therefore  $\mathbf{B}_j$  has one nonzero per row other than the  $j$ th row, which we let be zero. Moreover, by the injectivity of  $\phi_j$  then  $\mathbf{B}_j$  has at most one nonzero per column. As a result,

$$\mathbf{B}_j \mathbf{B}_j^T = \mathbf{I}_n - \mathbf{e}_j \mathbf{e}_j^T.$$

This implies

$$\mathbf{x}^T \mathbf{B}_j \mathbf{B}_j^T \mathbf{x} = \mathbf{x}^T (\mathbf{I}_n - \mathbf{e}_j \mathbf{e}_j^T) \mathbf{x} = \|\mathbf{x}\|_0 - x_j$$

for all  $j \in [n]$ . As

$$\|\bar{\mathbf{u}}(\mathbf{x})\| = \left\| \frac{1}{n} \sum_{j=1}^n \mathbf{u}_j(\mathbf{x}) \right\| \leq \frac{1}{n} \sum_{j=1}^n \|\mathbf{u}_j(\mathbf{x})\| \leq \|\mathbf{u}_j(\mathbf{x})\|,$$

then

$$\|\bar{\mathbf{u}}(\mathbf{x})\|^2 \leq \|\mathbf{u}_j(\mathbf{x})\|^2 = \frac{1}{2} (\|\mathbf{x}\|_0 - x_j) + 1$$

for all  $j \in [n]$  as claimed.  $\square$

### A.5 Euclidean distance bounds between normalized vectors

Here we recall some basic bounds pertaining to normalized vectors.

**Lemma A.8.** *Define  $f(\mathbf{x}) = \frac{\mathbf{x}}{\|\mathbf{x}\|^2}$  for all  $\mathbf{x} \in \mathbb{R}^q$ . Suppose without loss of generality that  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^q$  and  $\|\mathbf{y}\| \geq \|\mathbf{x}\| > 0$ , then*

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq \frac{3\|\mathbf{x} - \mathbf{y}\|}{\|\mathbf{x}\|^2}.$$

*Proof.* First observe

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{y}) &= \frac{\mathbf{x}}{\|\mathbf{x}\|^2} - \frac{\mathbf{y}}{\|\mathbf{y}\|^2} \\ &= \frac{\mathbf{x}}{\|\mathbf{x}\|^2} - \frac{\mathbf{y}}{\|\mathbf{y}\|^2} + \frac{\mathbf{y}}{\|\mathbf{x}\|^2} - \frac{\mathbf{y}}{\|\mathbf{x}\|^2} \\ &= \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x}\|^2} + \frac{\mathbf{y}(\|\mathbf{y}\|^2 - \|\mathbf{x}\|^2)}{\|\mathbf{x}\|^2 \|\mathbf{y}\|^2}. \end{aligned}$$

Taking the norm on both sides and applying the triangle inequality we have

$$\|f(\mathbf{x}) - f(\mathbf{y})\| = \frac{\|\mathbf{x} - \mathbf{y}\|}{\|\mathbf{x}\|^2} + \frac{\|\mathbf{y}\|(|\|\mathbf{y}\|^2 - \|\mathbf{x}\|^2|)}{\|\mathbf{x}\|^2 \|\mathbf{y}\|^2}.$$

By assumption  $\|\mathbf{y}\| \geq \|\mathbf{x}\|$ , therefore

$$\left| \|\mathbf{y}\|^2 - \|\mathbf{x}\|^2 \right| = |\langle \mathbf{y} - \mathbf{x}, \mathbf{y} + \mathbf{x} \rangle| \leq \|\mathbf{y} - \mathbf{x}\| \|\mathbf{y} + \mathbf{x}\| \leq \|\mathbf{y} - \mathbf{x}\| (\|\mathbf{y}\| + \|\mathbf{x}\|) \leq 2\|\mathbf{y}\| \|\mathbf{y} - \mathbf{x}\|.$$

This implies

$$\begin{aligned} \|f(\mathbf{x}) - f(\mathbf{y})\| &\leq \frac{\|\mathbf{x} - \mathbf{y}\|}{\|\mathbf{x}\|^2} + \frac{2\|\mathbf{y}\|^2 \|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|^2 \|\mathbf{y}\|^2} \\ &= \frac{3\|\mathbf{x} - \mathbf{y}\|}{\|\mathbf{x}\|^2} \end{aligned}$$

as claimed.  $\square$

### A.6 Hoeffding's inequality in Hilbert space

Lemma 4.10 rests on the application of the following concentration bound for sums of independent, mean zero, bounded random vectors. We remark that this is a specialization of more general results for martingales in 2-smooth Banach spaces.

**Lemma A.9.** [*Specialization of (Pinelis, 1994, Thm. 3.5)*] *For  $i \in [N]$  and  $R \in \mathbb{R}_{>0}$ , let  $\mathbf{x}_i \in \mathbb{R}^q$  be independent, mean zero random vectors which satisfy  $\|\mathbf{x}_i\| \leq R$  almost surely. Let  $S_N = \sum_{i=1}^N \mathbf{x}_i$ , then for  $t \in \mathbb{R}_{\geq 0}$  we have*

$$\mathbb{P}(\|S_N\| \geq t) \leq \exp\left(-\frac{t^2}{2NR^2}\right).$$

## B Proofs of results

### B.1 Proof of Theorem 3.2

**Theorem 3.2.** *Let  $D \subset \{0, 1\}^n$  be a set which can be strictly memorized and assume  $\|\mathbf{x}\|_0 \leq m \in \mathbb{N}_{\geq 1}$  for all  $\mathbf{x} \in D$ . Let  $\mathcal{D}$  be a probability distribution on  $D$ , and consider a random sample  $\mathcal{S}_N = (\mathbf{x}_i)_{i \in [N]}$ , where  $\mathbf{x}_i \sim \mathcal{D}$  are mutually i.i.d. Let  $\hat{\boldsymbol{\omega}} = \text{HSVM}(\mathcal{S}_N)$ ,  $\boldsymbol{\omega}^* = \text{HSVM}(D)$ ,  $\boldsymbol{\omega}(0) \in \mathbb{R}^q$  be arbitrary and  $\boldsymbol{\omega}(t)$  be generated for all  $t \in \mathbb{N}_{\geq 1}$  as per (8),  $\delta, \epsilon \in (0, 1)$  and assume  $\mathbf{x} \sim \mathcal{D}$  is sampled independent of  $\mathcal{S}_N$ . If  $N \gtrsim \epsilon^{-2} n \|\boldsymbol{\omega}^*\|^2 m \log(1/\delta)$  then both*

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\hat{\boldsymbol{\omega}}) \neq \mathbf{x}) \leq \epsilon \quad \text{and} \quad \mathbb{P}(H(\mathbf{x}; \mathbf{V}\boldsymbol{\omega}(t)) \neq \mathbf{x}) = \tilde{O}\left(\frac{\sqrt{m}\|\boldsymbol{\omega}^*\|}{\log(t)}\right) + \epsilon$$

hold with probability at least  $1 - \delta$  over the sample  $\mathcal{S}_N$ .

*Proof.* For any  $t \in \mathbb{R}$  define the margin loss  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  as

$$\phi(t) = \begin{cases} 0, & 1 \leq t, \\ 1 - t, & 0 \leq t \leq 1, \\ 1, & t \leq 0, \end{cases}$$

and note trivially for all  $t \in \mathbb{R}$  that  $\mathbb{1}(t \leq 0) \leq \phi(t)$ . We note on occasion we overload this notation and apply  $\phi$  to vectors by applying it elementwise. Observe for any  $\mathbf{x} \in \mathbb{R}^n$  and  $\boldsymbol{\omega} \in \mathbb{R}^q$  with  $\boldsymbol{\theta} = \mathbf{V}\boldsymbol{\omega}$ , that

$$\mathbb{1}(H(\mathbf{x}; \boldsymbol{\theta}) \neq \mathbf{x}) \leq \mathbb{1}(\exists j \in [n] : \mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega} \leq 0) = \mathbb{1}\left(\min_{j \in [n]} \mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega} \leq 0\right) \leq \phi\left(\min_{j \in [n]} \mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega}\right) = \max_{j \in [n]} \phi(\mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega}).$$

For any  $\mathbf{z} \in \mathbb{R}^n$ , let  $\ell(\mathbf{z}) = \max_{j \in [n]} \phi(z_j)$ . Using the fact that  $\phi$  is 1-Lipschitz, for any  $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^n$  we have

$$|\ell(\mathbf{z}) - \ell(\mathbf{z}')| \leq \max_{j \in [n]} |\phi(z_j) - \phi(z'_j)| = \|\phi(\mathbf{z}) - \phi(\mathbf{z}')\|_\infty \leq \|\mathbf{z} - \mathbf{z}'\|_\infty \leq \|\mathbf{z} - \mathbf{z}'\|_2.$$

Therefore  $\ell$  is 1-Lipschitz with respect to the Euclidean norm. Let  $\mathbf{U}(\mathbf{x}) \in \mathbb{R}^{n \times q}$  denote the matrix whose  $j$ th row is  $\mathbf{u}_j(\mathbf{x})^T \in \mathbb{R}^{1 \times q}$  for all  $j \in [n]$ . Furthermore, for some  $\Lambda \in \mathbb{R}_{>0}$ , define

$$\mathcal{H}_\Lambda = \{\mathbf{x} \mapsto \mathbf{U}(\mathbf{x})\boldsymbol{\omega} : \boldsymbol{\omega} \in \mathbb{R}^q, \|\boldsymbol{\omega}\| \leq \Lambda\}$$

and let

$$\mathcal{G}_\Lambda = \{\mathbf{x} \mapsto (\ell \circ h)(\mathbf{x}) : h \in \mathcal{H}_\Lambda\}.$$

Note by construction that  $g \in \mathcal{G}_\Lambda$  implies  $g : \mathbb{R}^n \rightarrow [0, 1]$ . We now compute the empirical Rademacher complexity of  $\mathcal{G}_\Lambda$  on a sample  $\mathcal{S}_N = (\mathbf{x}_i)_{i \in [N]}$ , to this end let  $\boldsymbol{\sigma} \in \{\pm 1\}^N$  and  $\boldsymbol{\epsilon} \in \{\pm 1\}^{N \times n}$  be a random vector and matrix respectively whose entries are mutually i.i.d. and distributed uniformly on  $\{\pm 1\}$ . As  $\ell$  is 1-Lipschitz in the  $\ell_2$  norm, then applying a vector contraction inequality (Maurer, 2016, Corollary 1) we have

$$\begin{aligned} \tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{G}_\Lambda) &= \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{h \in \mathcal{H}_\Lambda} \frac{1}{N} \sum_{i=1}^N \sigma_i (\ell \circ h)(\mathbf{x}_i) \right] \\ &\leq \sqrt{2} \mathbb{E}_{\boldsymbol{\epsilon}} \left[ \sup_{\|\boldsymbol{\omega}\| \leq \Lambda} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i)^T \boldsymbol{\omega} \right] \\ &\leq \sqrt{2} \mathbb{E}_{\boldsymbol{\epsilon}} \left[ \sup_{\|\boldsymbol{\omega}\| \leq \Lambda} \langle \boldsymbol{\omega}, \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i) \rangle \right] \\ &\leq \frac{\sqrt{2}\Lambda}{N} \mathbb{E}_{\boldsymbol{\epsilon}} \left[ \left\| \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i) \right\| \right]. \end{aligned}$$

Let  $Z = \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i)$ , as  $t \mapsto \sqrt{t}$  is concave then  $\mathbb{E}_Z[\sqrt{\|Z\|^2}] \leq \sqrt{\mathbb{E}_Z[\|Z\|^2]}$  by Jensen's inequality. As a result

$$\begin{aligned} \mathbb{E}_\epsilon \left[ \left\| \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i) \right\| \right] &\leq \sqrt{\mathbb{E}_\epsilon \left[ \left\langle \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i), \sum_{l=1}^N \sum_{k=1}^n \epsilon_{lk} \mathbf{u}_k(\mathbf{x}_l) \right\rangle \right]} \\ &= \sqrt{\sum_{i=1}^N \sum_{j=1}^n \sum_{l=1}^N \sum_{k=1}^n \mathbb{E}_\epsilon[\epsilon_{ij} \epsilon_{lk}] \mathbf{u}_j(\mathbf{x}_i)^T \mathbf{u}_k(\mathbf{x}_l)}. \end{aligned}$$

The Rademacher random variables are mutually i.i.d., therefore

$$\mathbb{E}_\epsilon[\epsilon_{ij} \epsilon_{lk}] = \begin{cases} 1, & (i=l) \wedge (j=k), \\ 0, & \text{otherwise.} \end{cases}$$

Recall also from Lemma A.7 that for any  $i \in [N]$  and for all  $j \in [n]$

$$\|\mathbf{u}_j(\mathbf{x}_i)\|^2 = \frac{1}{2} (\|\mathbf{x}_i\|_0 - x_{ij}) + 1.$$

Under the assumption  $\|\mathbf{x}_i\|_0 \leq m$  for all  $i \in [N]$ , then

$$\begin{aligned} \tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{G}_\Lambda) &\leq \frac{\sqrt{2}\Lambda}{N} \mathbb{E}_\epsilon \left[ \left\| \sum_{i=1}^N \sum_{j=1}^n \epsilon_{ij} \mathbf{u}_j(\mathbf{x}_i) \right\| \right] \\ &\leq \frac{\sqrt{2}\Lambda}{N} \sqrt{\sum_{i=1}^N \sum_{j=1}^n \|\mathbf{u}_j(\mathbf{x}_i)\|^2} \\ &\leq \frac{\Lambda}{N} \sqrt{\sum_{i=1}^N \sum_{j=1}^n (\|\mathbf{x}_i\|_0 - x_{ij} + 2)} \\ &\leq \Lambda \sqrt{\frac{n(m+2)}{N}}. \end{aligned}$$

Let  $\delta \in \mathbb{R}_{>0}$ . Applying a Rademacher complexity bound, e.g., (Mohri et al., 2018, Thm 3.3), then with probability at least  $1 - \delta$  over the random sample  $\mathcal{S}_N$

$$\begin{aligned} \mathbb{E}[g(\mathbf{x})] &\leq \frac{1}{N} \sum_{i=1}^N g(\mathbf{x}_i) + 2\tilde{\mathfrak{R}}_{\mathcal{S}}(\mathcal{G}_\Lambda) + 3\sqrt{\frac{\log(2/\delta)}{2N}} \\ &\leq \frac{1}{N} \sum_{i=1}^N g(\mathbf{x}_i) + 2\Lambda \sqrt{\frac{n(m+2)}{N}} + 3\sqrt{\frac{\log(2/\delta)}{2N}} \end{aligned}$$

for all  $g \in \mathcal{G}_\Lambda$ . In what follows let  $\Lambda = \|\boldsymbol{\omega}^*\|$ . Then, with probability at least  $1 - \delta$  over the random sample  $\mathcal{S}_N$ , for any  $\boldsymbol{\omega} \in \mathbb{R}^q$  such that  $\|\boldsymbol{\omega}\| \leq \|\boldsymbol{\omega}^*\|$  we have

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\boldsymbol{\omega}) \neq \mathbf{x}) \leq \frac{1}{N} \sum_{i=1}^N \phi(\min_{j \in [n]} \mathbf{u}_j(\mathbf{x}_i)^T \boldsymbol{\omega}) + 2\sqrt{\frac{\|\boldsymbol{\omega}^*\|^2 n(m+2)}{N}} + 3\sqrt{\frac{\log(2/\delta)}{2N}} \quad (10)$$

First we consider  $\hat{\boldsymbol{\omega}}$ : for any sample  $\mathcal{S}_N$  trivially  $\text{set}(\mathcal{S}_N) \subseteq D$ , therefore  $\boldsymbol{\omega}^* \in \mathcal{F}_\omega(\mathcal{S}_N)$ . As a result, with probability one  $\|\hat{\boldsymbol{\omega}}\| \leq \|\boldsymbol{\omega}^*\|$ . Conditioning on this event, observe also that  $\min_{j \in [n]} \mathbf{u}_j(\mathbf{x}_i)^T \hat{\boldsymbol{\omega}} \geq 1$  for  $i \in [N]$ . As a result, with probability at least  $1 - \delta$  over  $\mathcal{S}_N$  we have

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\hat{\boldsymbol{\omega}}) \neq \mathbf{x}) \leq \sqrt{\frac{4\|\boldsymbol{\omega}^*\|^2 n(m+2)}{N}} + \sqrt{\frac{9\log(2/\delta)}{2N}}.$$

As  $(a+b)^2 \leq 2(a^2+b^2)$  for  $a, b \in \mathbb{R}$  and assuming  $m \geq 1$ , then for any  $\epsilon \in \mathbb{R}_{>0}$ , if  $N \gtrsim \epsilon^{-2} n \|\boldsymbol{\omega}^*\|^2 m \log(1/\delta)$

$$\mathbb{P}_{\mathbf{x}}(H(\mathbf{x}; \boldsymbol{\theta}) \neq \mathbf{x}) \leq \epsilon$$

with probability at least  $1 - \delta$  over the sample  $\mathcal{S}_N$ .

We now turn our attention to  $\boldsymbol{\omega}(t)$ : recall for any  $a \in \mathbb{R}_{>0}$  as  $E(\mathbf{x}; a\boldsymbol{\theta}) = aE(\mathbf{x}; \boldsymbol{\theta})$  then

$$a(E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta})) \geq 0 \iff E(\mathbf{x}^{(j)}; a\boldsymbol{\theta}) - E(\mathbf{x}; a\boldsymbol{\theta}) \geq 0.$$

Indeed, this implies the set of memories in a Hopfield network is invariant under positive re-scalings of the parameters. As a consequence, for any distribution  $\mathcal{D}$  on  $\{0, 1\}^n$ ,  $a \in \mathbb{R}_{>0}$  and  $\boldsymbol{\theta} \in \Omega$ , if  $\mathbf{x} \sim \mathcal{D}$  we have

$$\mathbb{P}(H(\mathbf{x}; \boldsymbol{\theta}) \neq \mathbf{x}) = \mathbb{P}(H(\mathbf{x}; a\boldsymbol{\theta}) \neq \mathbf{x}).$$

Therefore, if we define

$$\bar{\boldsymbol{\omega}}(t) = \frac{\|\hat{\boldsymbol{\omega}}\|}{\|\boldsymbol{\omega}(t)\|} \boldsymbol{\omega}(t)$$

it follows that

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\boldsymbol{\omega}(t)) \neq \mathbf{x}) = \mathbb{P}(H(\mathbf{x}; \mathbf{V}\bar{\boldsymbol{\omega}}(t)) \neq \mathbf{x}).$$

From (Soudry et al., 2018, Theorem 5),

$$\left\| \frac{\bar{\boldsymbol{\omega}}(t)}{\|\bar{\boldsymbol{\omega}}\|} - \frac{\hat{\boldsymbol{\omega}}}{\|\hat{\boldsymbol{\omega}}\|} \right\| = \left\| \frac{\boldsymbol{\omega}(t)}{\|\boldsymbol{\omega}(t)\|} - \frac{\hat{\boldsymbol{\omega}}}{\|\hat{\boldsymbol{\omega}}\|} \right\| = O\left(\frac{\log(\log(t))}{\log(t)}\right).$$

Recalling  $\|\hat{\boldsymbol{\omega}}\| \leq \|\boldsymbol{\omega}^*\|$  this implies

$$\|\bar{\boldsymbol{\omega}}(t) - \hat{\boldsymbol{\omega}}\| = O\left(\frac{\|\boldsymbol{\omega}^*\| \log(\log(t))}{\log(t)}\right).$$

As a result, for all  $i \in [N]$  we have

$$\begin{aligned} \min_{j \in [n]} \mathbf{u}_j(\mathbf{x}_i)^T \bar{\boldsymbol{\omega}}(t) &= \min_{j \in [n]} (\mathbf{u}_j(\mathbf{x}_i)^T \hat{\boldsymbol{\omega}} + \mathbf{u}_j(\mathbf{x}_i)^T (\bar{\boldsymbol{\omega}}(t) - \hat{\boldsymbol{\omega}})) \\ &\geq 1 - \max_{j \in [n]} \|\mathbf{u}_j(\mathbf{x}_i)\| \|\bar{\boldsymbol{\omega}}(t) - \hat{\boldsymbol{\omega}}\| \\ &\geq 1 - O\left(\frac{\sqrt{m} \|\boldsymbol{\omega}^*\| \log(\log(t))}{\log(t)}\right), \end{aligned}$$

where the final inequality follows from Lemma A.7 and the assumption  $m \geq 1$ . By the definition of  $\phi$  it follows that

$$\phi(\min_{j \in [n]} \mathbf{u}_j(\mathbf{x}_i)^T \bar{\boldsymbol{\omega}}(t)) = O\left(\frac{\sqrt{m} \|\boldsymbol{\omega}^*\| \log(\log(t))}{\log(t)}\right)$$

for all  $i \in [N]$ . Using (10), this implies with probability at least  $1 - \delta$  over  $\mathcal{S}_N$

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\boldsymbol{\omega}(t)) \neq \mathbf{x}) \leq O\left(\frac{\sqrt{m} \|\boldsymbol{\omega}^*\| \log(\log(t))}{\log(t)}\right) + \sqrt{\frac{4 \|\boldsymbol{\omega}^*\|^2 n m}{N}} + \sqrt{\frac{9 \log(2/\delta)}{N}}.$$

As a result, if  $N \gtrsim \epsilon^{-2} n \|\boldsymbol{\omega}^*\|^2 m \log(1/\delta)$  then

$$\mathbb{P}(H(\mathbf{x}; \mathbf{V}\boldsymbol{\omega}(t)) \neq \mathbf{x}) = \tilde{O}\left(\frac{\sqrt{m} \|\boldsymbol{\omega}^*\|}{\log(t)}\right) + \epsilon.$$

with probability at least  $1 - \delta$  over the sample  $\mathcal{S}_N$ . □

## B.2 Properties of invariant parameters

### B.2.1 Energy bounds for invariant parameters across orbits

The following result states, in the context of invariant parameters, that energy bound differences for a point in an orbit extend to the entire orbit.

**Lemma 4.5.** *Let  $\mathbf{x}_0 \in \{0, 1\}^n$  and  $\boldsymbol{\theta} \in \Psi(\Gamma_n)$ . For  $\delta \in \mathbb{R}$ , if  $E(\mathbf{x}_0^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}_0; \boldsymbol{\theta}) \geq 1 - \delta$  for all  $j \in [n]$ , then for all  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  it follows that  $E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \geq 1 - \delta$  for all  $j \in [n]$ .*

*Proof.* As  $\boldsymbol{\theta} \in \Psi(\Gamma_n)$ , recall the intertwining relation (9), for  $\mathbf{x} \in \{0, 1\}^n$  and  $\mathbf{Q} \in \Gamma_n$  then

$$E(\mathbf{Q}\mathbf{x}; \boldsymbol{\theta}) = E(\mathbf{x}; \mathbf{Q}\boldsymbol{\theta}) = E(\mathbf{x}; \boldsymbol{\theta}).$$

As for any  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  there exists a  $\mathbf{Q} \in \Gamma_n$ , corresponding to a permutation  $\pi \in \Pi_n^{\mathcal{Q}}$ , such that  $\mathbf{x}_0 = \mathbf{Q}\mathbf{x}$ , this implies

$$1 - \delta \leq E(\mathbf{x}_0^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}_0; \boldsymbol{\theta}) = E((\mathbf{Q}\mathbf{x})^{(j)}; \boldsymbol{\theta}) - E(\mathbf{Q}\mathbf{x}; \boldsymbol{\theta}) = E(\mathbf{x}^{\pi(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}).$$

As  $\pi$  is a bijection then for all  $j \in [n]$  this implies

$$E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \geq 1 - \delta$$

as claimed.  $\square$

### B.2.2 Characterizing the invariant set for edge adjacency preserving permutations

The following lemma identifies the invariant parameters with respect to the set of edge adjacency preserving permutations as a particular rank three subspace.

**Lemma 4.6.** *Let  $F : \mathbb{R}^3 \rightarrow \Theta$  denote the linear map defined as follows: if  $(\mathbf{W}, \mathbf{b}) = F(\boldsymbol{\beta})$  then for all  $i, j \in [n]$  we have  $w_{ij} = 0$  if  $i = j$ ,  $w_{ij} = \beta_1$  if  $i \sim j$ ,  $w_{ij} = \beta_2$  if  $i \approx j$  and  $b_j = \beta_3$ . Then  $\Psi(\mathcal{Q}_n) = F(\mathbb{R}^3)$  where  $F(\mathbb{R}^3)$  denotes the image of  $F$ .*

*Proof.* First we show that  $F(\mathbb{R}^3) \subseteq \Psi(\mathcal{Q}_n)$ . Suppose  $\boldsymbol{\theta} \in F(\mathbb{R}^3)$ , then there exists  $\boldsymbol{\beta} \in \mathbb{R}^3$  such that  $\boldsymbol{\theta} = F(\boldsymbol{\beta})$ . Let  $\mathbf{Q} \in \mathcal{Q}_n$  and  $\pi \in \Pi_n^{\mathcal{Q}}$  be the corresponding permutation. Then  $b_i = b_{\pi(i)} = \beta_3$  for all  $i \in [n]$  and as a result  $\mathbf{Q}^T \mathbf{b} = \mathbf{b}$ . Furthermore, if  $i, j \in [n]$  then  $\pi(i) \sim \pi(j)$  if and only if  $i \sim j$ . Therefore if  $i \sim j$  then  $W_{ij} = \beta_1 = W_{\pi(i)\pi(j)}$ . Otherwise, if  $i \approx j$  then either  $i = j$ , which implies  $W_{jj} = 0 = W_{\pi(j)\pi(j)}$ , or  $i \neq j$  and then  $W_{ij} = \beta_2 = W_{\pi(i)\pi(j)}$ . As a result it follows that  $\mathbf{Q}^T \mathbf{W} \mathbf{Q} = \mathbf{W}$ , this implies  $\mathbf{Q}\boldsymbol{\theta} = \boldsymbol{\theta}$  and so  $\boldsymbol{\theta} \in \Psi(\mathcal{Q}_n)$ . We therefore conclude that  $F(\mathbb{R}^3) \subseteq \Psi(\mathcal{Q}_n)$ .

Now assume  $\boldsymbol{\theta} \in \Psi(\mathcal{Q}_n)$ , to prove there exists a  $\beta_3 \in \mathbb{R}$  such that  $b_i = \beta_3$  for all  $i \in [n]$  it suffices to show  $b_i = b_j$  for all  $i, j \in [n]$ . Similarly, to show there exist  $\beta_1, \beta_2 \in \mathbb{R}$  as per the statement of the lemma, it suffices to show  $w_{ij} = w_{ab}$  whenever either of the following hold: i)  $i \sim j$  and  $a \sim b$  or ii)  $i \approx j$  and  $a \approx b$ . It therefore suffices to prove the following two statements.

1. For any  $i, j \in [n]$  there exists a  $\pi \in \Pi_n^{\mathcal{Q}}$  such that  $\pi(i) = j$ . Note, as  $\boldsymbol{\theta} \in \Psi(\mathcal{Q}_n)$  this implies  $b_i = b_{\pi(i)} = b_j$ .
2. For any  $i, j, a, b \in [n]$  satisfying  $i \sim j$  and  $a \sim b$ , or  $i \approx j$  and  $a \approx b$ , there exists a  $\pi \in \Pi_n^{\mathcal{Q}}$  such that  $\pi(i) = a$  and  $\pi(j) = b$ . Note, and again as  $\boldsymbol{\theta} \in \Psi(\mathcal{Q}_n)$ , this implies  $w_{ij} = w_{\pi(i)\pi(j)} = w_{ab}$ .

In all that follows, for  $l \in [8]$  let  $\nu_l \in [v]$ . To prove the first statement, let  $i, j \in [n]$  and suppose  $i = \text{Ind}(\{\nu_1, \nu_2\})$  and  $j = \text{Ind}(\{\nu_3, \nu_4\})$ . Consider the permutation  $\pi \in \Pi_n^{\Phi} \subseteq \Pi_n^{\mathcal{Q}}$  which swaps the indices of the unordered vertex pairs involving  $\nu_1$  with the corresponding pair involving  $\nu_3$ , likewise for  $\nu_2$  and  $\nu_4$ , and is identity otherwise. To be clear, this is the permutation satisfying for  $t \in [2]$  the identities  $\pi(\text{Ind}(\{\nu_t, \nu\})) = \text{Ind}(\{\nu_{t+2}, \nu\})$  and  $\pi(\text{Ind}(\{\nu_{t+2}, \nu\})) = \text{Ind}(\{\nu_t, \nu\})$  for all  $\nu \in [v] \setminus \{\nu_t, \nu_{t+2}\}$ , and  $\pi(\text{Ind}^{-1}(\{\nu, \nu'\})) = \text{Ind}^{-1}(\{\nu, \nu'\})$  for all  $\nu, \nu' \in [v] \setminus \{\nu_1, \dots, \nu_4\}$ . Note if  $i \sim j$  then we can without loss of generality assume  $\nu_2 = \nu_4$  and this

permutation reduces to swapping a single pair and treating the rest with identity. By construction this permutation preserves adjacency, moreover  $\pi(i) = j$ . As a result, for all  $i, j \in [n]$  there exists a  $\pi \in \Pi_n^Q$  such that  $b_i = b_{\pi(i)} = b_j$ .

To prove the second statement, let  $i, j, a, b \in [n]$  and suppose  $i = \text{Ind}(\{\nu_1, \nu_2\})$ ,  $j = \text{Ind}(\{\nu_3, \nu_4\})$ ,  $a = \text{Ind}(\{\nu_5, \nu_6\})$  and  $b = \text{Ind}(\{\nu_7, \nu_8\})$ . Now consider the permutation  $\pi \in \Pi_n^\Phi \subseteq \Pi_n^Q$  which swaps the indices of the unordered vertex pairs involving  $\nu_1$  with those of  $\nu_5$ ,  $\nu_2$  with those of  $\nu_6$ ,  $\nu_3$  with those of  $\nu_7$ ,  $\nu_4$  with those of  $\nu_8$  and acts as identity on the indices of all other edges. To be clear, this is the permutation satisfying for  $t \in [4]$  the identities  $\pi(\text{Ind}^{-1}(\{\nu_t, \nu\})) = \text{Ind}^{-1}(\{\nu_{t+4}, \nu\})$  and  $\pi(\text{Ind}^{-1}(\{\nu_{t+4}, \nu\})) = \text{Ind}^{-1}(\{\nu_t, \nu\})$  for all  $\nu \in [v] \setminus \{\nu_t, \nu_{t+4}\}$ , and  $\pi(\text{Ind}^{-1}(\{\nu, \nu'\})) = \text{Ind}^{-1}(\{\nu, \nu'\})$  for all  $\nu, \nu' \in [v] \setminus \{\nu_1, \dots, \nu_8\}$ . By construction this permutation preserves adjacency and  $\pi(i) = a$ ,  $\pi(j) = b$ . Moreover as  $\pi \in \Pi_n^Q$  then  $i \sim j$  implies  $a \sim b$  and  $i \approx j$  implies  $a \approx b$ . As a result  $w_{ij} = w_{ab}$  for all  $i, j, a, b \in [n]$  if either  $i \sim j$  and  $a \sim b$  or  $i \approx j$  and  $a \approx b$ .

With both statements proved we conclude  $\Psi(\mathcal{Q}_n) \subseteq F(\mathbb{R}^3)$ . Finally, as  $\Psi(\mathcal{Q}_n) \subseteq F(\mathbb{R}^3)$  and  $F(\mathbb{R}^3) \subseteq \Psi(\mathcal{Q}_n)$ , then  $\Psi(\mathcal{Q}_n) = F(\mathbb{R}^3)$  as claimed.  $\square$

The next lemma states that parameters invariant to edge adjacency preserving permutations can memorize any binary vector. In combination with Lemma 4.5 this implies any graph isomorphism class is strictly memorizable.

**Lemma 4.7.** *For  $m \in [0, n]$ , let  $\beta = [2, 2, 1 - 2m] \in \mathbb{R}^3$  and  $\theta = F(\beta)$ . Then  $E(\mathbf{x}^{(j)}; \theta) - E(\mathbf{x}; \theta) \geq 1$  for all  $j \in [n]$  and for all  $\mathbf{x} \in \{0, 1\}^n$  satisfying  $\|\mathbf{x}\|_0 = m$ .*

*Proof.* Let  $\mathbf{x} \in \{0, 1\}^n$  satisfy  $|\text{supp}(\mathbf{x})| = m \in [0, n]$ . Recall from Lemma A.1 that for any  $j \in [n]$

$$E(\mathbf{x}^{(j)}; \theta) - E(\mathbf{x}; \theta) = y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle,$$

where  $y_j(\mathbf{x}) = 1 - 2x_j$ ,  $\mathbf{z}(\mathbf{x}) = [\mathbf{x}, 1]$  and  $\theta_j = [\mathbf{w}_j, 1]$ . Furthermore,

$$\begin{aligned} \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle &= \mathbf{x}^T \mathbf{w}_j + \beta_3 \\ &= \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \neq j) w_{jl} + \beta_3 \\ &= \beta_1 \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \sim j) + \beta_2 \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \approx j \wedge j \neq l) + \beta_3. \end{aligned}$$

Let  $c_j(\mathbf{x}) = \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \sim j)$  denote the number of edges of the graph adjacent to the  $j$ th edge. Then as

$$\begin{aligned} m &= \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge j \neq l) + \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge j = l) \\ &= \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \sim j \wedge j \neq l) + \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \approx j \wedge j \neq l) + \mathbb{1}(j \in \text{supp}(\mathbf{x})) \\ &= \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \sim j) + \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \approx j \wedge j \neq l) + x_j, \end{aligned}$$

it follows that

$$\sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \approx j \wedge j \neq l) = m - c_j(\mathbf{x}) - x_j.$$

As a result, the condition that must be satisfied for all  $j \in [n]$  is

$$y_j(\mathbf{x}) \langle \mathbf{z}(\mathbf{x}), \theta_j \rangle = (1 - 2x_j) (c_j(\mathbf{x})(\beta_1 - \beta_2) + \beta_2(m - x_j) + \beta_3) \geq 1.$$

If  $\beta_1 = \beta_2$  then the left-hand side simplifies to an expression which depends only on the sparsity of the representation of the graph. Under this assumption, it suffices to find a  $\beta_2, \beta_3 \in \mathbb{R}$  such that

$$(1 - 2x_j)(\beta_2(m - x_j) + \beta_3) \geq 1.$$

Let  $\beta_3 = 1 - \beta_2 m$ , if  $x_j = 0$  then

$$(1 - 2x_j)(\beta_2(m - x_j) + \beta_3) = \beta_2 m + \beta_3 = 1$$

while if  $x_j = 1$  then

$$-(\beta_2(m - 1) + \beta_3) = 1 - \beta_2.$$

Therefore, with  $\beta_1 = \beta_2 = 2$  and  $\beta_3 = 1 - 2m$  we have

$$E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \geq 1$$

for all  $j \in [n]$ . □

We now make a few remarks in regard to the the construction used in the previous lemma. First,  $F(2, 2, 1 - 2m)$  memorizes  $\mathbf{x} \in \{0, 1\}^n$  iff  $\|\mathbf{x}\|_0 = m$ . Indeed, the only if aspect can be demonstrated as follows: if  $\mathbf{x}'$  satisfies  $\|\mathbf{x}'\| = m + \delta$  for  $\delta \in \mathbb{N}_{\geq 0}$ , the required inequalities become  $2\delta + 1 \geq 1$  for  $j \in \text{supp}(\mathbf{x}')$  and  $-2\delta + 1 \geq 1$  for  $j \notin \text{supp}(\mathbf{x}')$ . These inequalities can only simultaneously hold if  $\delta = 0$ . Second

$$\|F(2, 2, 1 - 2m)\|^2 = 2v(v + 1) + (1 - 2m)^2,$$

therefore when the sparsity  $m$  is proportional to  $n$  then the norm scales like  $\Theta(n)$ .

### B.3 Constructing an invariant, small norm parameter which memorizes $k$ -cliques

Our goal in this section is to show that small norm parameters exist which can memorize specific isomorphism classes. In particular, we consider the case of  $k$ -cliques: recall that a  $k$ -clique graph has a fully connected subset of  $k$  vertices while the remaining  $v - k$  vertices are isolated. We denote the set of representations of  $k$ -cliques on  $v$  vertices as  $\mathcal{C}_{v,k}$  and trivially note  $|\mathcal{C}_{v,k}| = \binom{v}{k}$ . Towards constructing low-norm invariant parameters that strictly memorize all  $k$ -cliques, the following lemma derives specific expressions for the energy difference derived in Lemma A.1. To state this result, for  $\mathbf{x} \in \mathcal{C}_{v,k}$  let  $\text{Clique}(\mathbf{x}) \subset [v]$  denote the subset of the vertices of the graph which are in the fully connected subset.

**Lemma B.1.** *Let  $\boldsymbol{\beta} \in \mathbb{R}^3$  and suppose  $\boldsymbol{\theta} = F(\boldsymbol{\beta}) = \mathbf{V}\boldsymbol{\omega}$ , for some  $\boldsymbol{\omega} \in \mathbb{R}^q$ . For  $\mathbf{x} \in \mathcal{C}_{v,k}$  and any  $j \in [n]$ , define  $r = |\text{Clique}(\mathbf{x}) \cap \text{Ind}^{-1}(j)| \in \{0, 1, 2\}$  as the number of vertices in the  $j$ th vertex pair which are also in the clique of  $\mathbf{x}$ . Then for any  $j \in [n]$*

$$\mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega} = y_j(\mathbf{x}) \mathbf{z}(\mathbf{x})^T \boldsymbol{\theta}_j = \begin{cases} \frac{\beta_2}{2} k^2 - \frac{\beta_2}{2} k + \beta_3, & r = 0, \\ \frac{\beta_2}{2} k^2 + (\beta_1 - \frac{3}{2}\beta_2) k + (\beta_2 + \beta_3 - \beta_1), & r = 1, \\ -\left(\frac{\beta_2}{2} k^2 + (2\beta_1 - \frac{5}{2}\beta_2) k + (3\beta_2 - 4\beta_1 + \beta_3)\right), & r = 2. \end{cases}$$

*Proof.* By definition

$$\begin{aligned} \mathbf{z}(\mathbf{x})^T \boldsymbol{\theta}_j &= \mathbf{w}_j^T \mathbf{x} + b_j \\ &= \sum_{l=1}^n w_{jl} \mathbb{1}(l \in \text{supp}(\mathbf{x})) + \beta_3 \\ &= \beta_1 \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \sim j) + \beta_2 \sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \approx j \wedge l \neq j) + \beta_3 \end{aligned}$$

Observe each  $j \in [n]$  can be placed in one of three distinct categories with respect to  $\mathbf{x}$ : in particular, either both, one or neither of the vertices of  $j$  are in the  $k$ -clique of the graph represented by  $\mathbf{x}$ . Fixing an arbitrary  $j \in [n]$ , we denote these events in turn as  $\Phi_r$  for  $r \in \{0, 1, 2\}$ , where

$$\Phi_r(\mathbf{x}) = \{j \in [n] : |\text{Ind}^{-1}(j) \cap \text{Clique}(\mathbf{x})| = r\}.$$

Note  $\Phi_2(\mathbf{x}) = \text{supp}(\mathbf{x})$ . If  $j \in \Phi_0(\mathbf{x})$  then neither of the vertices of  $j$  are in the clique of  $\mathbf{x}$ , as a result the  $j$ th edge cannot be adjacent to any edge in the clique. If  $j \in \Phi_1(\mathbf{x})$  then exactly one vertex of  $j$  is in the clique, furthermore there are  $k - 1$  other vertices in the clique this vertex is connected to via an edge. Finally, if  $j \in \Phi_2(\mathbf{x})$  then both of its vertices are connected via edges to  $k - 2$  other vertices in the clique. As a result,

$$\sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \sim j) = \begin{cases} 0, & j \in \Phi_0(\mathbf{x}), \\ k - 1, & j \in \Phi_1(\mathbf{x}), \\ 2(k - 2), & j \in \Phi_2(\mathbf{x}). \end{cases}$$

Moreover, as there are  $\binom{k}{2}$  edges in total in a  $k$ -clique, and as if  $l \in \text{supp}(\mathbf{x})$  then  $j = l$  can be true only if  $j \in \text{supp}(\mathbf{x}) \in \Phi_2(\mathbf{x})$ , then

$$\sum_{l=1}^n \mathbb{1}(l \in \text{supp}(\mathbf{x}) \wedge l \approx j \wedge l \neq j) = \begin{cases} \binom{k}{2}, & j \in \Phi_0(\mathbf{x}), \\ \binom{k}{2} - k + 1, & j \in \Phi_1(\mathbf{x}), \\ \binom{k}{2} - 2k + 3, & j \in \Phi_2(\mathbf{x}). \end{cases}$$

As a result: if  $j \in \Phi_0(\mathbf{x})$  then

$$\mathbf{z}(\mathbf{x})^T \boldsymbol{\theta}_j = \left( \beta_2 \frac{k(k-1)}{2} + \beta_3 \right) = \frac{\beta_2}{2} k^2 - \frac{\beta_2}{2} k + \beta_3.$$

If  $j \in \Phi_1(\mathbf{x})$  then

$$\mathbf{z}(\mathbf{x})^T \boldsymbol{\theta}_j = \beta_1(k-1) + \beta_2 \frac{k^2 - 3k + 2}{2} + \beta_3 = \frac{\beta_2}{2} k^2 + \left( \beta_1 - \frac{3}{2} \beta_2 \right) k + (\beta_2 + \beta_3 - \beta_1).$$

Finally, if  $j \in \Phi_2(\mathbf{x})$  then

$$\mathbf{z}(\mathbf{x})^T \boldsymbol{\theta}_j = \beta_1(2k-4) + \beta_2 \frac{k^2 - 5k + 6}{2} + \beta_3 = \frac{\beta_2}{2} k^2 + \left( 2\beta_1 - \frac{5}{2} \beta_2 \right) k + (3\beta_2 - 4\beta_1 + \beta_3).$$

To conclude, observe  $y_j(\mathbf{x}) = (1 - 2x_{ij}) = -1$  iff  $j \in \Phi_2(\mathbf{x})$ . □

We now derive a simple bound on the norm of parameters which are invariant to edge adjacency preserving permutations.

**Lemma B.2.** *Let  $\beta \in \mathbb{R}^3$  and  $\boldsymbol{\theta} = F(\beta^3)$ . Then*

$$\|\boldsymbol{\theta}\|^2 \leq \beta_2^2 v^4 + 2\beta_1^2 v^3 + \beta_3^2 v^2.$$

*Proof.* For any fixed edge index  $r \in [n]$ , note as each vertex of this edge is a member of  $v - 2$  other vertex pairs then

$$\sum_{c=1}^n \mathbb{1}(c \sim r) = 2(v - 2)$$

Moreover, as there are  $\binom{v}{2}$  unordered vertex pairs in total

$$\sum_{c=1}^n \mathbb{1}(c \approx r \wedge c \neq r) = \binom{v}{2} - 2(v - 2) - 1 = \frac{v^2 - 3v + 6}{2}.$$

Therefore, and also noting that  $n = \binom{v}{2} \leq v^2$ , we have

$$\begin{aligned}
\|\boldsymbol{\theta}\|^2 &= \|\mathbf{W}\|_F^2 + \|\mathbf{b}\|^2 \\
&= \sum_{r=1}^n \left( \beta_1^2 \sum_{c=1}^n \mathbb{1}(c \sim r) + \beta_2^2 \sum_{c=1}^n \mathbb{1}(c \not\sim r \wedge c \neq r) + \beta_3^2 \right) \\
&= n \left( \beta_1^2 2(v-2) + \beta_2^2 \left( \frac{v^2 - 3v + 6}{2} \right) + \beta_3^2 \right) \\
&\leq n(\beta_2^2 v^2 + 2\beta_1^2 v + \beta_3^2) \\
&\leq \beta_2^2 v^4 + 2\beta_1^2 v^3 + \beta_3^2 v^2.
\end{aligned}$$

□

We now present a low norm construction for memorizing  $k$ -cliques: in particular, Lemma 4.8 illustrates that the  $k$ -clique graph isomorphism class can be memorized using a parameter whose norm is  $O(\sqrt{n})$ . This is in contrast to the general construction used in the proof of Lemma 4.7 whose norm grew as  $\Theta(n)$ .

**Lemma 4.8.** *If  $\boldsymbol{\beta} = [-10/k, 38/k^2, 0] \in \mathbb{R}^3$ ,  $\boldsymbol{\theta} = F(\boldsymbol{\beta})$ , and  $k \geq 5$ , then the following hold.*

1.  $E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \geq 1$  for all  $\mathbf{x} \in \mathcal{C}_{v,k}$  and all  $j \in [n]$ .
2. If  $k = cv$  for some constant  $c \in (0, 1]$ , then there exists a constant  $C > 0$  such that  $\|\boldsymbol{\theta}\|^2 \leq Cv$ .

*Proof.* Using Lemma B.1, and substituting  $\beta_1 = -10/k$ ,  $\beta_2 = 38/k^2$  and  $\beta_3 = 0$ , we obtain

$$\mathbf{u}_j(\mathbf{x})^T \boldsymbol{\omega} = \begin{cases} 19 \left(1 - \frac{1}{k}\right), & r = 0, \\ 9 - \frac{47}{k} + \frac{38}{k^2}, & r = 1, \\ 1 + \frac{55}{k} - \frac{114}{k^2}, & r = 2. \end{cases}$$

If  $k \geq 5$ , then

$$19 \left(1 - \frac{1}{k}\right) \geq 1,$$

while

$$9 - \frac{47}{k} + \frac{38}{k^2} \geq 9 - \frac{47}{5} + \frac{38}{25} > 1,$$

and

$$1 + \frac{55}{k} - \frac{114}{k^2} \geq 1.$$

Hence  $E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \geq 1$  for all  $\mathbf{x} \in \mathcal{C}_{v,k}$  and all  $j \in [n]$ .

It remains to bound the norm. Since  $F(\boldsymbol{\beta})$  assigns weight  $\beta_1$  to adjacent edge pairs, weight  $\beta_2$  to non-adjacent edge pairs, and has zero bias, we have

$$\|\boldsymbol{\theta}\|^2 = \|\mathbf{W}\|_F^2 \leq n2(v-2)\beta_1^2 + n(n-1)\beta_2^2.$$

Since  $n = \binom{v}{2}$ ,  $\beta_1^2 = 100/k^2$ , and  $\beta_2^2 = 38^2/k^4$ , if  $k = cv$  then

$$n2(v-2)\beta_1^2 = O(v), \quad n(n-1)\beta_2^2 = O(1).$$

Therefore there exists a constant  $C > 0$ , depending only on  $c$ , such that  $\|\boldsymbol{\theta}\|^2 \leq Cv$ . □

### B.3.1 Invariance of the full orbit HSVM solution

The following lemma states a well known result that the average orbit action of a parameter is invariant to the action of the underlying group.

**Lemma B.3.** *Let  $\theta = (\mathbf{W}, \mathbf{b}) \in \Theta$  and  $\Gamma_n$  denote a subgroup of  $\mathcal{P}_n$ . Then  $\text{Proj}_{\Gamma_n}(\theta) = \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\theta \in \Psi(\Gamma_n)$ .*

*Proof.* For typographical ease let  $\theta' = \text{Proj}_{\Gamma_n}(\theta)$ . Then

$$\theta' = (\mathbf{W}', \mathbf{b}') := \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\theta = \left( \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{W} \mathbf{Q}, \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{b} \right).$$

Given  $\Gamma_n$  is a subgroup, then for any  $\mathbf{Q}' \in \Gamma_n$

$$\begin{aligned} \mathbf{Q}'\theta' &= (\mathbf{Q}'^T \mathbf{W}' \mathbf{Q}', \mathbf{Q}'^T \mathbf{b}') \\ &= \left( \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} (\mathbf{Q}\mathbf{Q}')^T \mathbf{W} (\mathbf{Q}\mathbf{Q}'), \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} (\mathbf{Q}\mathbf{Q}')^T \mathbf{b} \right) \\ &= \left( \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{W} \mathbf{Q}, \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{b} \right) \\ &= \theta', \end{aligned}$$

therefore  $\theta' \in \Psi(\Gamma_n)$ . □

Using the previous lemma, we show that the full orbit HSVM solution lies in the invariant set.

**Lemma 4.9.** *Let  $\mathbf{x}_0 \in \{0, 1\}^n$  and  $\Gamma_n$  denote a subgroup of  $\mathcal{P}_n$  and assume  $\text{Orb}(\mathbf{x}_0, \Gamma_n)$  can be strictly memorized. If  $\theta^* = \mathbf{V}\omega^*$  where  $\omega^* = \text{HSVM}_{\Theta}(\text{Orb}(\mathbf{x}_0, \Gamma_n))$  then  $\theta^* \in \Psi(\Gamma_n)$ .*

*Proof.* We use a symmetrization argument. To this end, with  $\theta^* = (\mathbf{W}^*, \mathbf{b}^*)$  let

$$\theta = (\mathbf{W}, \mathbf{b}) := \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\theta^* = \left( \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{W}^* \mathbf{Q}, \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{b}^* \right).$$

By Lemma B.3 we know that  $\theta \in \Psi(\Gamma_n)$ . By the definition of  $\theta^*$  we also have

$$E(\mathbf{x}^{(j)}, \theta^*) - E(\mathbf{x}, \theta^*) \geq 1.$$

Therefore, using both the intertwining property (9) and Lemma A.1, for any  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  and  $j \in [n]$ , and with  $\mathbf{z}(\mathbf{x}) = [\mathbf{x}, 1]$ , we have

$$\begin{aligned} E(\mathbf{x}^{(j)}, \theta) - E(\mathbf{x}, \theta) &= (2x_j - 1)\mathbf{z}(\mathbf{x})^T \theta_j \\ &= \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} (2x_j - 1)\mathbf{z}(\mathbf{x})^T \mathbf{Q}\theta_j^* \\ &= \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} E(\mathbf{x}^{(j)}, \mathbf{Q}\theta^*) - E(\mathbf{x}, \mathbf{Q}\theta^*) \\ &= \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} E(\mathbf{x}^{(j)}, \theta^*) - E(\mathbf{x}, \theta^*) \\ &= E(\mathbf{x}^{(j)}, \theta^*) - E(\mathbf{x}, \theta^*) \\ &\geq 1. \end{aligned}$$

As a result,  $\theta$  is a feasible point of the HSVM problem (7) defined on the full orbit dataset  $\text{Orb}(\mathbf{x}_0, \Gamma_n)$ . Therefore, by the definition of  $\theta^*$  it must follow that  $\|\theta^*\| \leq \|\theta\|$ . On the other hand, using the triangle inequality and the fact that  $Q \in \Gamma_n$  is a permutation, we have

$$\|\theta\| = \left\| \frac{1}{|\Gamma_n|} \sum_{Q \in \Gamma_n} Q\theta^* \right\| \leq \frac{1}{|\Gamma_n|} \sum_{Q \in \Gamma_n} \|Q\theta^*\| = \|\theta^*\|.$$

This implies  $\frac{1}{2}\|\theta^*\|^2 = \frac{1}{2}\|\theta\|^2$ , as this objective is 1-strongly convex this in turn implies  $\theta^* = \theta \in \Psi(\Gamma_n)$ .  $\square$

## B.4 Approximately invariant parameters

### B.4.1 Approximate invariance is sufficient for generalization

The following lemma states a sufficient condition for strict memorization of an orbit dataset based on proximity to the relevant invariant space. In particular, given a graph  $G \in \mathcal{G}_v$  and letting  $\mathbf{x}_0 = \mathcal{E}_{rep}(G)$ , if  $\omega^* = \text{HSVM}(\mathcal{S})$ , where  $\mathcal{S} \subset \text{Orb}(\mathbf{x}_0, \Phi_n)$ , and  $\omega^*$  is sufficiently close to the subspace  $\mathcal{Q}_n$ , then  $\theta^* = \mathbf{E}\omega^*$  will strictly memorize all graphs isomorphic to  $G$ .

**Lemma B.4.** *Let  $\mathbf{x}_0 \in \{0, 1\}^n$  satisfy  $\|\mathbf{x}_0\| = m \in \mathbb{N}_{\geq 2}$ ,  $\theta = \mathbf{E}\omega \in \Theta_n$  satisfy  $E(\mathbf{x}_0^{(j)}; \theta) - E(\mathbf{x}_0; \theta) \geq 1$  for all  $j \in [n]$ , and  $\theta' = \mathbf{E}\omega' \in \Psi(\Gamma_n)$  be such that  $\|\omega - \omega'\| \leq \frac{1}{4\sqrt{m}}$ . Then  $E(\mathbf{x}^{(j)}; \theta') - E(\mathbf{x}; \theta) \geq \frac{1}{2}$  for all  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  and  $j \in [n]$ .*

*Proof.* Inspecting Lemma A.7, if  $\|\mathbf{x}_0\| = m \in \mathbb{N}_{\geq 2}$  then  $\|\mathbf{u}_j(\mathbf{x})\| \leq \sqrt{m}$  for all  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  and  $j \in [n]$ . By assumption  $\mathbf{u}_j(\mathbf{x}_0)^T \omega \geq 1$  for all  $j \in [n]$ , therefore

$$\begin{aligned} E(\mathbf{x}_0^{(j)}; \theta') - E(\mathbf{x}_0; \theta') &= y_j(\mathbf{x}_0)z(\mathbf{x}_0)^T \theta'_j \\ &= \mathbf{u}_j(\mathbf{x}_0)^T \omega' \\ &= \mathbf{u}_j(\mathbf{x}_0)^T \omega - \mathbf{u}_j(\mathbf{x}_0)^T (\omega - \omega') \\ &\geq 1 - \|\mathbf{u}_j(\mathbf{x}_0)\| \|\omega - \omega'\| \\ &\geq \frac{3}{4} \end{aligned}$$

for all  $j \in [n]$ . As  $\theta' \in \Psi(\Gamma_n)$ , then Lemma 4.5 implies for any  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  that

$$E(\mathbf{x}^{(j)}; \theta') - E(\mathbf{x}; \theta') = \mathbf{u}_j(\mathbf{x})^T \omega' \geq \frac{3}{4}.$$

Moreover, and using the same trick as before, we also observe for all  $\mathbf{x} \in \text{Orb}(\mathbf{x}_0, \Gamma_n)$  that

$$\begin{aligned} E(\mathbf{x}^{(j)}; \theta) - E(\mathbf{x}; \theta) &= y_j(\mathbf{x})z(\mathbf{x})^T \theta_j \\ &= \mathbf{u}_j(\mathbf{x})^T \omega \\ &= \mathbf{u}_j(\mathbf{x})^T \omega' - \mathbf{u}_j(\mathbf{x})^T (\omega' - \omega) \\ &\geq \frac{3}{4} - \|\mathbf{u}_j(\mathbf{x})\| \|\omega' - \omega\| \\ &\geq \frac{1}{2} \end{aligned}$$

for all  $j \in [n]$ .  $\square$

### B.4.2 Proximity of AHSVM solution to the invariant set

**Lemma B.5.** *Let  $\mathcal{S} \subset \{0, 1\}^n$ ,  $\mu_{\mathcal{S}} = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \bar{\mathbf{u}}(\mathbf{x})$  and assume  $\omega^* = \text{AHSVM}(\mathcal{S})$  is feasible. Then  $\omega^* = \frac{\mu_{\mathcal{S}}}{\|\mu_{\mathcal{S}}\|^2}$ .*

*Proof.* Note by the feasibility assumption  $\frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \bar{\mathbf{u}}(\mathbf{x}) \neq \mathbf{0}_q$ . Forming the Lagrangian with Lagrange variable  $\alpha$  we have

$$\mathcal{L}(\omega, \alpha) = \frac{1}{2}\|\omega\|^2 + \alpha \left( 1 - \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \langle \bar{\mathbf{u}}(\mathbf{x}), \omega \rangle \right).$$

Clearly this is a strongly convex objective with a unique minimizer  $\boldsymbol{\omega}^*$ . Zeroing the gradient with respect to  $\boldsymbol{\omega}$  and rearranging gives the identity

$$\boldsymbol{\omega}^* = \alpha^* \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \bar{\mathbf{u}}(\mathbf{x})$$

on the solution pair  $(\boldsymbol{\omega}^*, \alpha^*)$ . In addition, as there is only a single constraint and the problem is feasible then the constraint must be active, meaning

$$\frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \langle \bar{\mathbf{u}}(\mathbf{x}), \boldsymbol{\omega}^* \rangle = 1.$$

As a result

$$\|\boldsymbol{\omega}^*\|^2 = \alpha^* \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \langle \bar{\mathbf{u}}(\mathbf{x}), \boldsymbol{\omega}^* \rangle = \alpha^*. \quad (11)$$

Therefore

$$\frac{\boldsymbol{\omega}^*}{\|\boldsymbol{\omega}^*\|^2} = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \bar{\mathbf{u}}(\mathbf{x}) = \boldsymbol{\mu}_{\mathcal{S}}.$$

As  $\boldsymbol{\mu}_{\mathcal{S}} \neq \mathbf{0}_q$  then

$$\|\boldsymbol{\mu}_{\mathcal{S}}\|^2 = \left\| \frac{\boldsymbol{\omega}^*}{\|\boldsymbol{\omega}^*\|^2} \right\|^2 = \frac{1}{\|\boldsymbol{\omega}^*\|^2},$$

giving

$$\boldsymbol{\omega}^* = \frac{\boldsymbol{\mu}_{\mathcal{S}}}{\|\boldsymbol{\mu}_{\mathcal{S}}\|^2}$$

as claimed.  $\square$

Similar to Lemma 4.9, we now show that the full orbit AHSVM solution lies on the relevant invariance space.

**Lemma B.6.** *Let  $\mathbf{x}_0 \in \{0, 1\}^n$ ,  $\Gamma_n$  denote a subgroup of  $\mathcal{P}_n$ , and assume  $\mathcal{O} = \text{Orb}(\mathbf{x}_0, \Gamma_n)$  satisfies  $|\mathcal{O}| = |\Gamma_n|$ . Let  $\boldsymbol{\omega}^* = \text{AHSVM}_{\Theta}(\text{Orb}(\mathbf{x}_0, \Gamma_n))$  be feasible and define  $\boldsymbol{\theta}^* = \mathbf{V}\boldsymbol{\omega}^*$ , then  $\boldsymbol{\theta}^* \in \Psi(\Gamma_n)$ .*

*Proof.* Again we use a symmetrization argument. To this end, with  $\boldsymbol{\theta}^* = (\mathbf{W}^*, \mathbf{b}^*)$  let

$$\boldsymbol{\theta} = (\mathbf{W}, \mathbf{b}) := \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\boldsymbol{\theta}^* = \left( \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{W}^* \mathbf{Q}, \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}^T \mathbf{b}^* \right).$$

By Lemma B.3 we know that  $\boldsymbol{\theta} \in \Psi(\Gamma_n)$ . Let  $\mathbf{x} \in \mathcal{O}$  and  $\mathbf{Q} \in \Gamma_n$  be such that  $\mathbf{Q}\mathbf{x} = \mathbf{x}_0$ , and let  $\pi$  denote the permutation associated with  $\mathbf{Q}$ . As  $\boldsymbol{\theta} \in \Psi(\Gamma_n)$  then using the intertwining property (9)

$$\begin{aligned} \bar{\mathbf{u}}(\mathbf{x})^T \boldsymbol{\omega} &= \frac{1}{n} \sum_{j=1}^n E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{1}{n} \sum_{j=1}^n E(\mathbf{Q}\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{Q}\mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{1}{n} \sum_{j=1}^n E((\mathbf{Q}\mathbf{x})^{(\pi(j))}; \boldsymbol{\theta}) - E(\mathbf{Q}\mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{1}{n} \sum_{l=1}^n E(\mathbf{x}_0^{(l)}; \boldsymbol{\theta}) - E(\mathbf{x}_0; \boldsymbol{\theta}) \\ &= \bar{\mathbf{u}}(\mathbf{x}_0)^T \boldsymbol{\omega}. \end{aligned}$$

As the AHSVM problem is feasible and has a single constraint then  $\frac{1}{|\mathcal{O}|} \sum_{\mathbf{x} \in \mathcal{O}} \bar{\mathbf{u}}(\mathbf{x})^T \boldsymbol{\omega}^* = 1$ . Therefore

$$\begin{aligned}
\frac{1}{|\mathcal{O}|} \sum_{\mathbf{x} \in \mathcal{O}} \bar{\mathbf{u}}(\mathbf{x})^T \boldsymbol{\omega} &= \bar{\mathbf{u}}(\mathbf{x}_0)^T \boldsymbol{\omega} \\
&= \frac{1}{n} \sum_{l=1}^n E(\mathbf{x}_0^{(l)}; \boldsymbol{\theta}) - E(\mathbf{x}_0; \boldsymbol{\theta}) \\
&= \frac{1}{n} \sum_{l=1}^n E\left(\mathbf{x}_0^{(l)}; \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\boldsymbol{\theta}^*\right) - E\left(\mathbf{x}_0; \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\boldsymbol{\theta}^*\right) \\
&= \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \frac{1}{n} \sum_{l=1}^n E(\mathbf{Q}\mathbf{x}_0^{(l)}; \boldsymbol{\theta}^*) - E(\mathbf{Q}\mathbf{x}_0; \boldsymbol{\theta}^*) \\
&= \frac{1}{|\mathcal{O}|} \sum_{\mathbf{x} \in \mathcal{O}} \frac{1}{n} \sum_{j=1}^n E(\mathbf{x}^{(j)}; \boldsymbol{\theta}) - E(\mathbf{x}; \boldsymbol{\theta}) \\
&= \frac{1}{|\mathcal{O}|} \sum_{\mathbf{x} \in \mathcal{O}} \bar{\mathbf{u}}(\mathbf{x})^T \boldsymbol{\omega}^* \\
&= 1.
\end{aligned}$$

As a result,  $\boldsymbol{\theta}$  is a feasible point of the AHSVM problem defined on the full orbit dataset  $\text{Orb}(\mathbf{x}_0, \Gamma_n)$ . Therefore, by the definition of  $\boldsymbol{\theta}^*$  it must follow that  $\|\boldsymbol{\theta}^*\| \leq \|\boldsymbol{\theta}\|$ . On the other hand, using the triangle inequality and the fact that  $\mathbf{Q} \in \Gamma_n$  is a permutation, we have

$$\|\boldsymbol{\theta}\| = \left\| \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \mathbf{Q}\boldsymbol{\theta}^* \right\| \leq \frac{1}{|\Gamma_n|} \sum_{\mathbf{Q} \in \Gamma_n} \|\mathbf{Q}\boldsymbol{\theta}^*\| = \|\boldsymbol{\theta}^*\|.$$

This implies  $\frac{1}{2}\|\boldsymbol{\theta}^*\|^2 = \frac{1}{2}\|\boldsymbol{\theta}\|^2$ , as this objective is 1-strongly convex this in turn implies  $\boldsymbol{\theta}^* = \boldsymbol{\theta} \in \Psi(\Gamma_n)$ .  $\square$

The lemma below bounds the difference between the sample AHSVM solution and the population AHSVM solution, leveraging only the boundedness of the data.

**Lemma 4.10.** *Let  $\mathcal{O} \subseteq \{0, 1\}^n$  satisfy  $\|\mathbf{x}\|_0 \leq m \in \mathbb{N}_{\geq 2}$  for all  $\mathbf{x} \in \mathcal{O}$  and assume  $\boldsymbol{\omega}^* = \text{HSVM}(\mathcal{O})$  is feasible. Consider a random sample  $\mathcal{S} = (\mathbf{x}_i)_{i=1}^N$  where  $\mathbf{x}_i \sim U(\mathcal{O})$  are mutually i.i.d. and define  $\boldsymbol{\omega}_{\mathcal{O}} = \text{AHSVM}(\mathcal{O})$  and  $\boldsymbol{\omega}_{\mathcal{S}} = \text{AHSVM}(\mathcal{S})$ . For  $\delta \in (0, 1]$  and  $\epsilon \in \mathbb{R}_{>0}$ , if  $N \gtrsim \epsilon^{-2} \|\boldsymbol{\omega}^*\|^4 m \log(1/\delta)$  then  $\|\boldsymbol{\omega}_{\mathcal{S}} - \boldsymbol{\omega}_{\mathcal{O}}\| \leq \epsilon$  with probability at least  $1 - \delta$ .*

*Proof.* Let  $\boldsymbol{\mu} = \mathbb{E}[\bar{\mathbf{u}}(\mathbf{x})]$  where  $\mathbf{x} \sim U(\mathcal{O})$ , then  $\boldsymbol{\mu} = \frac{1}{|\mathcal{O}|} \sum_{\mathbf{x} \in \mathcal{O}} \bar{\mathbf{u}}(\mathbf{x})$ . In addition, let  $\hat{\boldsymbol{\mu}}_{\mathcal{S}} = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \bar{\mathbf{u}}(\mathbf{x})$ . Then using Lemma B.5 we have

$$\boldsymbol{\omega}_{\mathcal{O}} = \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|^2}, \quad \boldsymbol{\omega}_{\mathcal{S}} = \frac{\hat{\boldsymbol{\mu}}_{\mathcal{S}}}{\|\hat{\boldsymbol{\mu}}_{\mathcal{S}}\|^2}.$$

Taking norms this clearly also implies

$$\|\boldsymbol{\mu}\| = \frac{1}{\|\boldsymbol{\omega}_{\mathcal{O}}\|}, \quad \|\hat{\boldsymbol{\mu}}_{\mathcal{S}}\| = \frac{1}{\|\boldsymbol{\omega}_{\mathcal{S}}\|},$$

Observe by definition that  $\boldsymbol{\omega}^*$  satisfies  $E(\mathbf{x}^{(j)}; \boldsymbol{\theta}^*) - E(\mathbf{x}; \boldsymbol{\theta}^*) = \langle \mathbf{u}_j(\mathbf{x}), \boldsymbol{\omega}^* \rangle \geq 1$ , therefore for any  $S \subseteq \mathcal{O}$  we have

$$\frac{1}{|S|} \sum_{\mathbf{x} \in S} \langle \bar{\mathbf{u}}(\mathbf{x}), \boldsymbol{\omega}^* \rangle = \frac{1}{n|S|} \sum_{\mathbf{x} \in S} \sum_{j=1}^n \langle \mathbf{u}_j(\mathbf{x}), \boldsymbol{\omega}^* \rangle \geq 1.$$

As a result we have both  $\boldsymbol{\omega}^* \in \mathcal{F}_A(\mathcal{S})$  and  $\boldsymbol{\omega}^* \in \mathcal{F}_A(\mathcal{O})$ , which in turn implies  $\|\boldsymbol{\omega}^*\| \geq \|\boldsymbol{\omega}_{\mathcal{S}}\|$  and  $\|\boldsymbol{\omega}^*\| \geq \|\boldsymbol{\omega}_{\mathcal{O}}\|$ . Defining  $f(\mathbf{x}) = \mathbf{x}/\|\mathbf{x}\|^2$  for any  $\mathbf{x} \in \mathbb{R}^q$ , then applying Lemma A.8 this gives

$$\|\boldsymbol{\omega}_{\mathcal{S}} - \boldsymbol{\omega}_{\mathcal{O}}\| = \|f(\hat{\boldsymbol{\mu}}_{\mathcal{S}}) - f(\boldsymbol{\mu})\| \leq \frac{\|\hat{\boldsymbol{\mu}}_{\mathcal{S}} - \boldsymbol{\mu}\|}{\min\{\|\boldsymbol{\omega}_{\mathcal{O}}\|^{-2}, \|\boldsymbol{\omega}_{\mathcal{S}}\|^{-2}\}} \leq 3\|\boldsymbol{\omega}^*\|^2 \|\hat{\boldsymbol{\mu}}_{\mathcal{S}} - \boldsymbol{\mu}\|.$$

Observe

$$\|\hat{\boldsymbol{\mu}}_S - \boldsymbol{\mu}\| = \left\| \frac{1}{N} \sum_{i=1}^N (\bar{\mathbf{u}}(\mathbf{x}) - \boldsymbol{\mu}) \right\|,$$

clearly  $\bar{\mathbf{u}}(\mathbf{x}) - \boldsymbol{\mu}$  is a centered random vector, moreover for any  $\mathbf{x} \in \{0, 1\}$

$$\begin{aligned} \|\bar{\mathbf{u}}(\mathbf{x}) - \boldsymbol{\mu}\| &\leq \|\bar{\mathbf{u}}(\mathbf{x})\| + \|\boldsymbol{\mu}\| \\ &= \|\bar{\mathbf{u}}(\mathbf{x})\| + \left\| \frac{1}{|\mathcal{O}|} \sum_{\mathbf{x}' \in \mathcal{O}} \bar{\mathbf{u}}(\mathbf{x}') \right\| \\ &\leq \|\bar{\mathbf{u}}(\mathbf{x})\| + \frac{1}{|\mathcal{O}|} \sum_{\mathbf{x}' \in \mathcal{O}} \|\bar{\mathbf{u}}(\mathbf{x}')\| \\ &\leq 2\sqrt{m}, \end{aligned}$$

where the last inequality follows from Lemma A.7. We now deploy Lemma A.9, a specialization of (Pinelis, 1994, Th, 3.5). In particular, given some  $\epsilon \in \mathbb{R}_{\geq 0}$  and letting  $S_N = \sum_{i=1}^N (\bar{\mathbf{u}}(\mathbf{x}) - \boldsymbol{\mu})$  then

$$\mathbb{P}(\|\hat{\boldsymbol{\mu}}_S - \boldsymbol{\mu}\| \geq \epsilon) = \mathbb{P}(\|S_N\| \geq N\epsilon) \leq \exp\left(-\frac{N\epsilon^2}{4m}\right)$$

As a result, for  $\delta \in (0, 1]$ , if  $N \geq \frac{m}{4\epsilon^2} \log(1/\delta)$  then

$$\|\boldsymbol{\omega}_S - \boldsymbol{\omega}_{\mathcal{O}}\| \leq 3\|\boldsymbol{\omega}^*\|^2 \epsilon$$

with probability at least  $1 - \delta$ . To arrive at the claimed result we substitute  $\epsilon$  for  $\frac{\epsilon}{3\|\boldsymbol{\omega}^*\|^2}$ .  $\square$

## C Additional Experiments and Further Preliminary Results

### C.1 Further training and test accuracy curves

Fig. 4 shows test accuracy versus training sample size, with mean and min-max over 10 trials, for Hopfield networks trained by MEF, Perceptron, and Delta (the latter two used only as baselines; see Appendix A.2). For small graphs ( $v = 8$ ) we enumerate the full isomorphism class and report the true accuracy, i.e., the fraction of the class memorized. For larger graphs ( $v = 20$ ), accuracy is estimated on an independent random sample of 1000 graphs. Note, within our hyperparameter range, the Delta rule using Adam failed on the  $k$ -clique class, whereas MEF learned all classes and was insensitive to optimizer choice.

### C.2 Further parameter heatmaps

Fig. 5 shows the weight matrices for MEF and Delta on clique and Paley graph data. We observe that the Delta rule also returns solutions which approach the invariant space as the sample size  $N$  increases.

### C.3 Sample Complexity for Robust Exponential Memory

Robust exponential memory in Hopfield networks can be equivalently viewed through the lenses of error-correcting codes and the recovery of latent combinatorial structure under noise. Prior work has shown, specifically for  $k$ -cliques, that Hopfield networks can store entire isomorphism classes of patterns with large basins of attraction, achieving asymptotically optimal noise tolerance. As discussed in Section 4.3, Theorem 3.2, Lemma 4.7, and Lemma 4.8 together imply that a polynomial number of randomly sampled representatives from an isomorphism class suffices for MEF-trained Hopfield networks to memorize nearly all elements of that class. This provides theoretical evidence that MEF learning can recover invariant structure from sparse data, supporting and extending earlier conjectures on robust generalization in structured memory models. An important direction for future work is to determine whether polynomial sample complexity suffices to guarantee memorization of all elements of an isomorphism class, and whether such learned networks

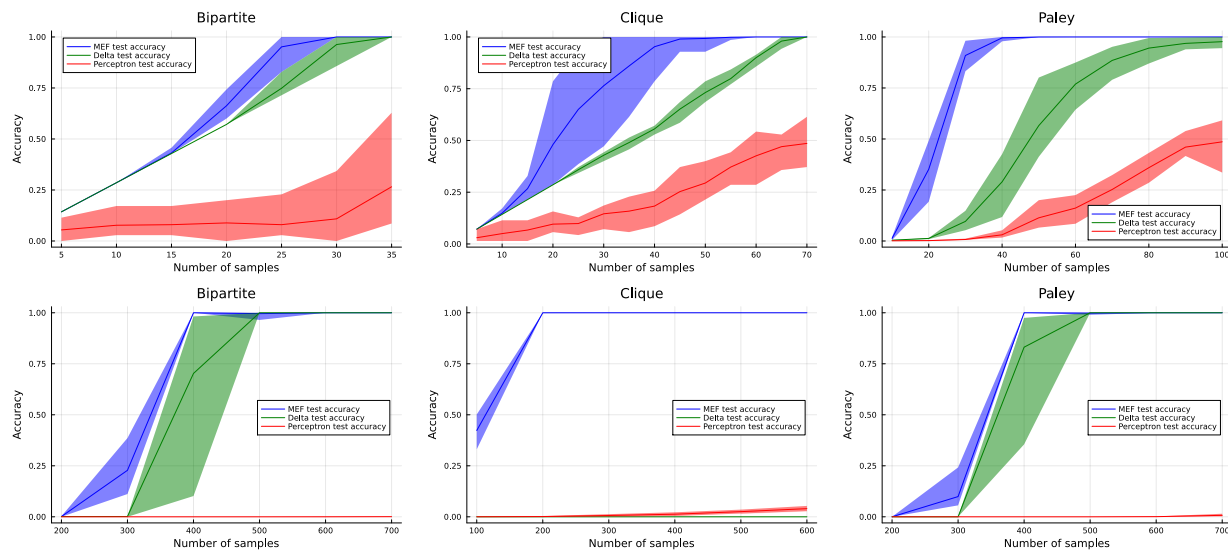


Figure 4: **Test accuracy vs. training sample size for isomorphism classes at two scales.** Top row:  $v = 8$  (isomorphism class sizes: bipartite 35, Paley 2520). Bottom row:  $v = 20$  (for reference class size for bipartite is 92,378). Curves show mean and min-max over 10 trials. Networks are trained with Perceptron, Delta (MSE), and MEF learning rules.

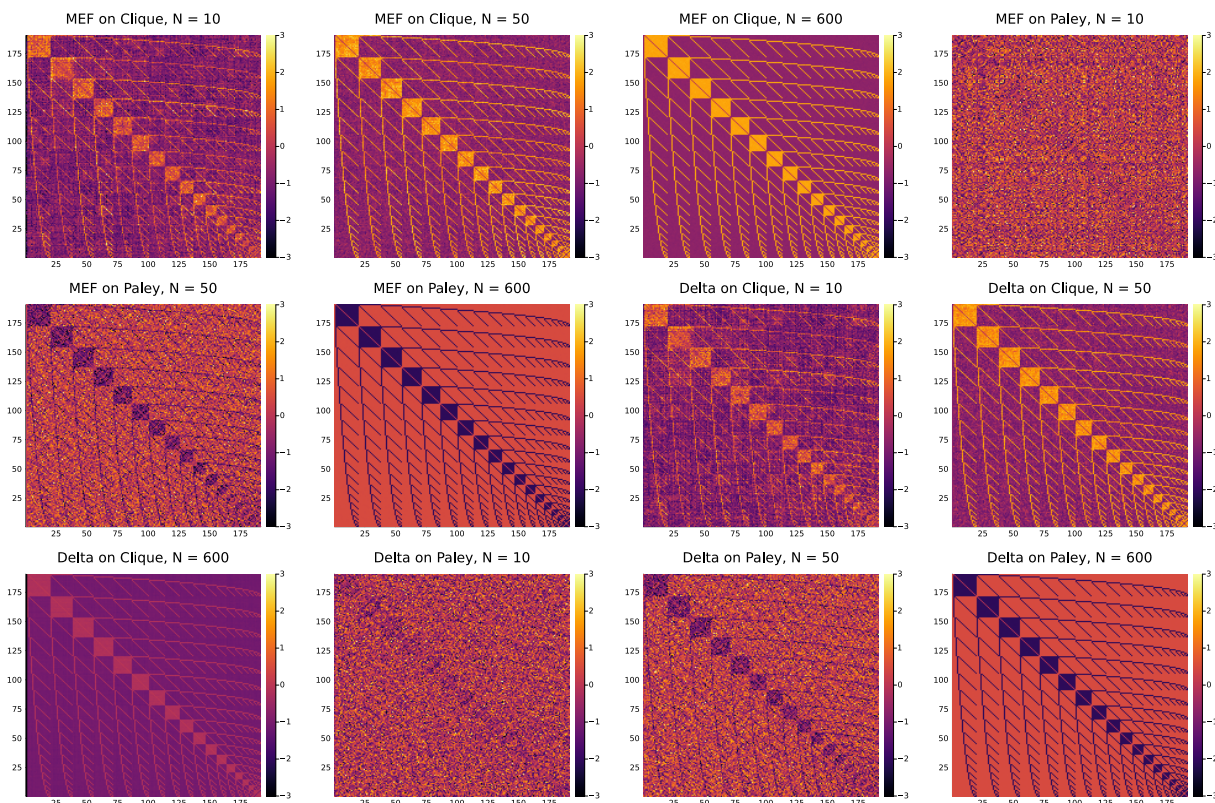


Figure 5: **Weights found by MEF and Delta on clique and Paley graph data while varying  $N$ :** networks where trained on samples from isomorphism class of 10-cliques and Paley graphs on  $v = 20$  vertices with sample size ranging from  $N = 10$  to 600.

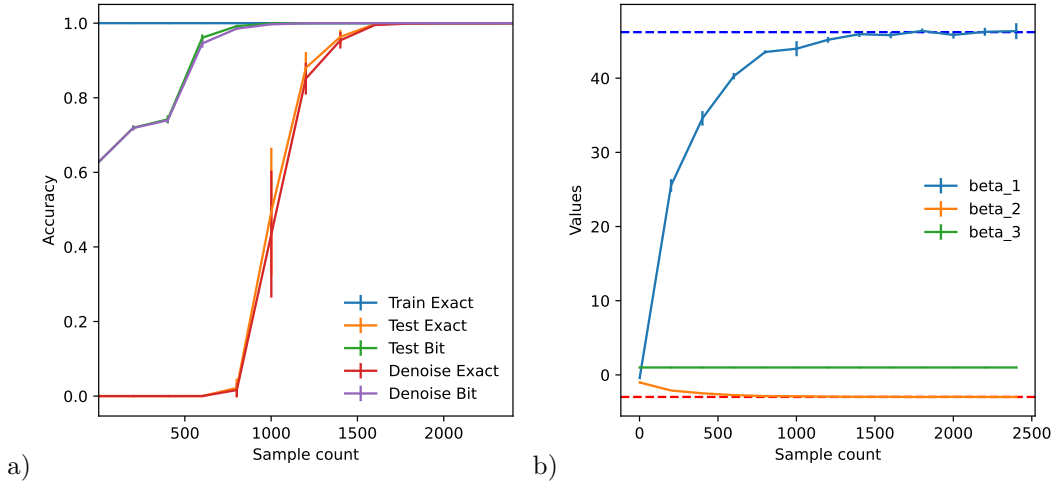


Figure 6: **Generalization on the Hidden Clique Problem.** a) Generalization and denoising accuracy (exact / average bits) are plotted as a function of number of 50-clique samples (in 100-vertex graphs;  $n = 4950$  bit networks) for MEF-trained HNs. Accuracy for generalization was computed using 10000 novel graphs as Test set. Denoising accuracy was computed by corrupting 5% bits in these 10000 and evaluating correctness of the attractor when dynamics is initialized at the noisy patterns (5 trials, standard deviation error bars). b) For each type of learned parameter (weights indexed by adjacent/non-adjacent edges, or thresholds), we plot their average over number of training samples (normalized so that thresholds have mean absolute value 1). Dotted horizontal lines are parameter type averages from training on the largest number of samples.

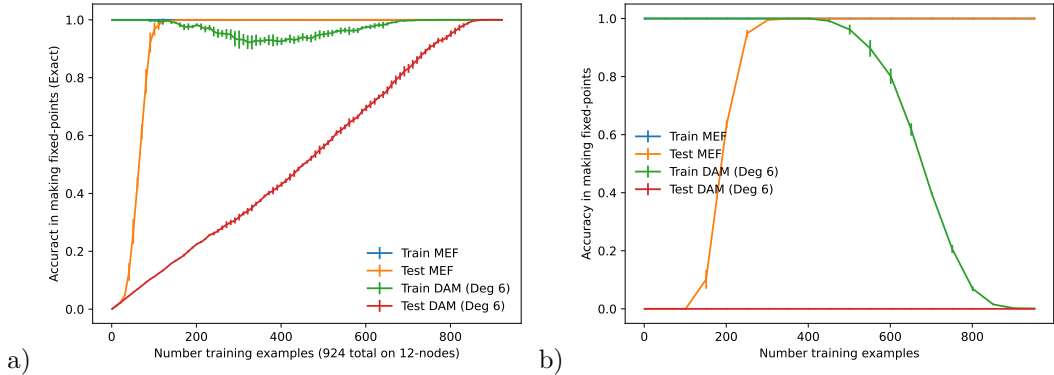


Figure 7: **DAM models trained on cliques.** We compare the generalization performance between DAMs and MEF-trained HNs for a) the 6-clique problem on graphs with  $v = 12$  vertices and b) the 16-clique problem on graphs with  $v = 32$  vertices (496-bit networks). The Train accuracy is the percentage of clique training samples that are neural network attractors. For a, the Test accuracy is the percentage of all 6-cliques that are fixed-points, and for b, it is that for a random set of 10000 16-cliques. Scores are averaged over four trials, with standard deviation error bars.

can reliably perform structured recovery tasks associated with that class under noise. The clique family and the Hidden Clique Problem serve as a concrete illustration of this phenomenon, but the underlying mechanism applies more broadly to invariant families of combinatorial patterns. For an example of this robust generalization, see Fig. 6, which plots over sample count both the generalization and denoising accuracy of very large HNs trained with MEF ( $v = 100$ ;  $n = 4950$  bits).

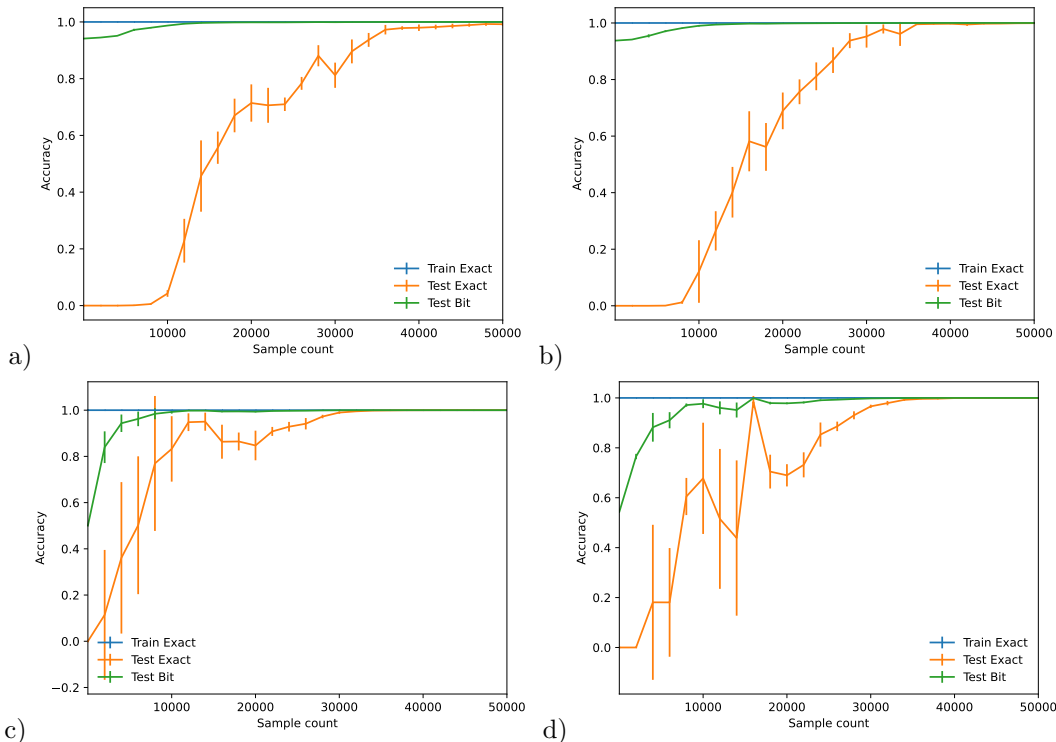


Figure 8: **Generalization for other graph classes.** Accuracy of MEF-trained HN networks generalizing as a function of number of training samples for: a) Chain graph with  $v = 32, k = 16$  (496-bit), b) Cycle graph with  $v = 32, k = 16$  (496-bit), c) Bipartite graph with  $v = 32, k = 16$  (496-bit), and d) Johnson graph  $J(7, 3)$  on  $v = 35$  vertices (595-bit). Averaged over 10 trials, with standard deviation error bars. Test sets consisted of 10000 novel graphs chosen randomly from the isomorphism class.

### C.4 Clique generalization in DAMs

We conducted experiments comparing the generalization performance of DAMs with that of MEF when presented with increasing numbers of cliques as training data. We used the architecture and algorithm described in Krotov & Hopfield (2016) for polynomial degree 6 in the energy function. The results in Fig. 7 show very different generalization behaviors. These preliminary investigations suggest that DAMs do not generalize well in the setting of cliques, but much more work is needed to understand their behavior.

### C.5 Graph Families

We also see generalization occurring for several other families of graphs. Fig. 8 shows Train and test generalization (bit and exact) accuracy for generalization in the four graph classes: chain, cycle, bipartite, and Johnson (see also Fig. 3). In these cases the integer  $k$  corresponds to number of vertices in the corresponding chain, cycle, or bipartite subset. Although most graph classes appear to have monotonic increases in accuracy as a function of number of training samples, it is interesting that for some settings there is a “double ascent/descent” property, which can be seen in Fig. 8d for the Johnson graph family. Another example of this phenomenon for the case of cliques appears in Section C.6. We also list several experiments with circulant graphs in Fig. 9.

### C.6 Double Descent Phenomenon

In our experiments we observed the interesting phenomenon that sometimes the Test error decreased with number of training samples, but then increased for a time, only to decrease again to perfect accuracy (zero

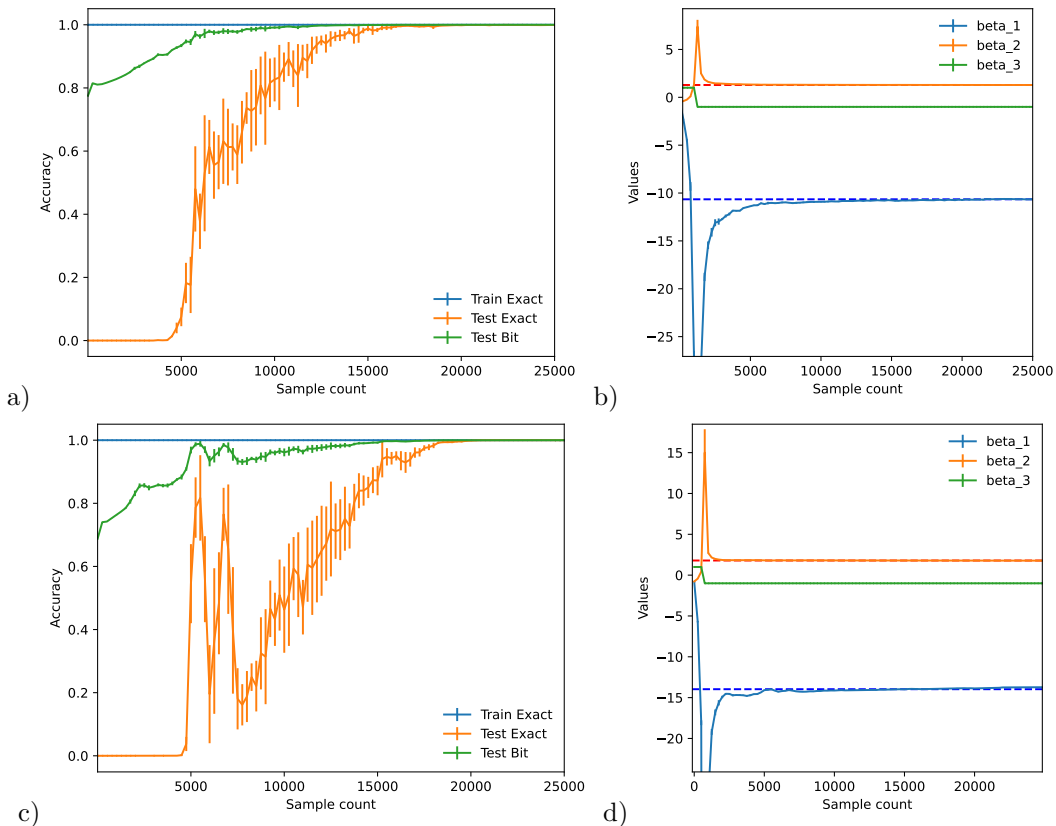


Figure 9: **Generalization for Circulant graphs.** Accuracy of MEF-trained HN networks generalizing as a function of number of training samples for: a) Circulant graph with  $v = 32$ , jump sequence  $[2, 4]$ , b) For each type of learned parameter (weights indexed by adjacent/non-adjacent edges, or thresholds) in data from a, we plot their average over number of training samples (normalized so that thresholds have mean absolute value 1; dotted horizontal lines are parameter type averages from training on the largest number of samples), c) Circulant graph with  $v = 32$ , jump sequence  $[2, 3, 5]$ , and d) Correspondingly as in b for graph in c. Plots are averaged over 10 trials, with standard deviation error bars. Test sets consisted of 10000 novel graphs chosen randomly from the isomorphism class.

error). We give two examples in Fig. 10 for bit error between attractors and samples for the case of cliques. For cliques, we note that this phenomenon seems to occur with smaller  $k$  relative to  $v$ . It is also interesting to see that when  $k$  is small, it takes many more samples to learn the full isomorphism class. In other words, it seems that the easiest  $k$ -clique isomorphism learning problem is when  $k$  is near  $v/2$ .

More examples of multiple (double and triple) ascent in accuracy can be found in Fig. 8cd and Fig. 9c for exact error (percentage of samples that are fixed points of the dynamics) in the case of Johnson and Circulant graphs. We postpone theoretical analysis of these findings for future work.

### C.7 The Hopfield Network Not Graph Isomorphic Check (HNNGIC)

Lemma 4.7 implies any isomorphism class of a graph can be stored in a Hopfield network. This prompts investigation into the potential for using Hopfield networks to check for graph isomorphisms, a fundamental and important problem in computer science. To this end we propose Algorithm 1, which we refer to as the *Hopfield Network Not Graph Isomorphic Check* (HNNGIC). As the name suggests, this algorithm provides a check if two graphs are not graph isomorphic, returning true in certain cases when they are not graph isomorphic and unknown, otherwise.

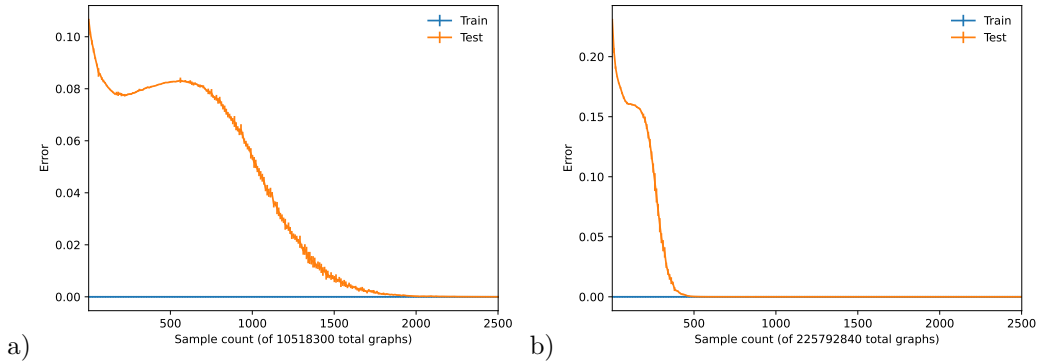


Figure 10: **Double Descent Phenomenon.** a) We trained Hopfield networks using MEF with increasing numbers of training samples for the case of learning cliques with  $v = 32, k = 8$  and find that average test error (on 10000 random cliques) decreases with number of training samples, then mysteriously increases, only to decrease again to zero. (5 trials for choices of randomly samples train clique data, standard deviation error bars). b) Same as in a, but for graphs with  $v = 32, k = 12$ . Note that the double descent phenomenon is less apparent for this larger choice of  $k$ .

---

**Algorithm 1:** Hopfield Network Not Graph Isomorphic Check (HNNGIC)

---

**Input:** two graphs  $\mathbf{x}_1, \mathbf{x}_2 \in \{0, 1\}^n$  and computational budget  $B$

**Output:** True or Unknown

**Step 1:** minimize  $L(F(\boldsymbol{\beta}); \mathbf{x}_1)$  within computational budget  $B$ , return  $\boldsymbol{\beta}^* \in \mathbb{R}^3$ ;

**Step 2:** **if**  $H(\mathbf{x}_1; F(\boldsymbol{\beta}^*)) = \mathbf{x}_1$  **then**

**if**  $H(\mathbf{x}_2; F(\boldsymbol{\beta}^*)) \neq \mathbf{x}_2$  **then**

**return** True

**end**

**else**

**return** Unknown

**end**

---

The idea behind this algorithm is simple: given two graphs, we pick one, i.e.,  $\mathbf{x}_1$ , arbitrarily at random. We then attempt to train the Hopfield network by minimizing the energy flow defined on this single graph, but restrict the parameters to lie on the edge adjacency invariant subspace  $\Psi(\mathcal{Q}_n)$ . If the resulting invariant parameters  $F(\boldsymbol{\beta}^*)$  strictly memorize  $\mathbf{x}_1$  then this implies every point in the orbit of  $\mathbf{x}_1$  under graph isomorphism is also strictly memorized. Therefore, if  $\mathbf{x}_2$  is graph isomorphic to  $\mathbf{x}_1$ , then it must be a fixed point. As a result,  $\mathbf{x}_2$  and  $\mathbf{x}_1$  cannot be graph isomorphic if  $\mathbf{x}_2$  is not a fixed point of an invariant Hopfield network which stores  $\mathbf{x}_1$ . Note, that if  $\mathbf{x}_2$  is a fixed point, it does not follow that  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are isomorphic. Indeed, there may be other fixed points, i.e., “spurious states”, in the landscape, not related to the orbit of  $\mathbf{x}_1$ . Also note that Lemmas 4.7 and 4.5 imply that there is always a 3-parameter network storing any graph; in particular, given enough computation, we are guaranteed to find an approximation of  $\boldsymbol{\beta}^*$  that is sufficient to store  $\mathbf{x}_1$ .