A MULTICOVER APPROACH TO NEURAL NETWORKS SAMPLE COMPLEXITY

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

045

Paper under double-blind review

ABSTRACT

Covering numbers are central to estimating sample complexity. Alas, standard techniques for bounding covering numbers fail in estimating the covering numbers of many classes of neural networks. We introduce a generalization of covers, called *multicovers*, which are covers w.r.t. many metrics simultaneously. Contrary to standard covering numbers, multicovering numbers behave better with the layer-wise structure in neural networks. We utilize this property to recover a recent result of Daniely & Granot (2019) who defined a new notion called Approximate Description Length (ADL) to establish tight bounds on the sample complexity of networks with weights of bounded Frobenius norm. We also show that ADL and multicovering numbers are closely related.

022 1 INTRODUCTION

Covering numbers are one of the most basic techniques for bounding the sample complexity of function classes, and can achieve state of the art bounds in various cases. Alas, it is not clear how to estimate covering numbers for function classes of layered architectures, such as neural networks. Indeed, state-of-the-art results still exhibit a polynomial gap between upper and lower sample complexity bounds. This is in contrast to non-layerd function classes, in which the gaps are often logarithmic or even constant.

A major flaw of covering numbers is that it is not clear how to use them inductively on the network's layers. That is, a bound on the covering number for function classes of depth i architectures, is not enough to derive a tight (up to log factors) bound for function classes of depth i + 1 architectures.

In this paper, we present a generalization of covering called *multicover*, which overcomes the above barrier, at least in some cases. This allows the derivation of tight bounds on various families of neural networks. In a nutshell, given a set S of $d \times d$ PSD matrices, an ϵ -multicover of a set $\mathcal{X} \subset \mathbb{R}^d$ is a set $\check{\mathcal{X}}$ that simultaneously forms an ϵ -cover w.r.t. to any metric of the form $d(\mathbf{x}, \mathbf{y}) =$ $\sqrt{(\mathbf{x} - \mathbf{y})^\top R(\mathbf{x} - \mathbf{y})}$ for $R \in S$. The ϵ -multicovering number of \mathcal{X} is the minimal size of an ϵ multicover of \mathcal{X} . We note that if S consists of a single matrix R, multicover is just a standard cover w.r.t. the norm corresponding to R. However, for other sets S we get a notion of covering that is fundamentally different from standard covering w.r.t. a metric.

We present techniques which allow for layerwise induction. Using these techniques and for the case of S being the class of PSD matrices with trace at most 1, we show the following. Given a bound on the multicovering number of \mathcal{X} , a class \mathcal{L} of $d \times d$ matrices, and a non-linearity $\sigma : \mathbb{R} \to \mathbb{R}$, we derive bounds which are often tight on the multicovering number of

$$\mathcal{LX} = \{A\mathbf{x} : A \in \mathcal{L}, \mathbf{x} \in \mathcal{X}\} \text{ and } \sigma(\mathcal{X}) = \{(\sigma(x_1), \dots, \sigma(x_d)) : \mathbf{x} \in \mathcal{X}\}$$

The tools we present allow us to derive nearly tight bounds on the sample complexity of constant depth networks with weights of bounded norm. For instance, assume that the activation is the ReLUlike softplus activation $\sigma(x) = \log(1 + e^x)$ and consider the class \mathcal{N} of networks of depth l, width d, and weight matrices with spectral norm at most O(1) and Frobenius norm at most R. This and similar classes have been studied intensively in recent years, because sample complexity bounds on such classes can potentially be sublinear in the number of network parameters, thus shedding light on a main mystery of modern neural networks. We show that if the input distribution is supported in $[-1, 1]^d$ then the sample complexity of \mathcal{N} is $\tilde{O}(dR^2)$ which is sublinear in number of parameters and is tight up to poly-log factors. As far as we know, despite extensive efforts, such results are not known to be derived via "standard" covering number techniques, or even more generally, via other common techniques such as Radamacher complexity. nevertheless, we note that similar bounds were recently proved (Daniely & Granot, 2019) using a notion called Approximate Description Length (ADL). We show that ADL is closely related to multicover, and in a sense, multicover can be seen as a "dual" approach to ADL.
We hope that having both the ADL technique and the multi-covering technique at our disposal will lead to further progress in the future.

Throughout this paper, absent proofs for theorems, lemmas, and claims appear in full form in the appendix.

063 064 065

076 077

085

1.1 COVERING NUMBERS AND LAYERWISE INDUCTION

We next give a simple example in which layerwise induction fails to establish tight bounds on covering numbers. We emphasize that the goal of this example is to demonstrate the problem with layerwise induction on covering numbers, but it is not a proof that the approach is doomed to fail in general.

We will consider the class \mathcal{H} of linear classifiers of norm ≤ 1 over $B_{\sqrt{d}}^d$. That is, \mathcal{H} consists of all functions $h: B_{\sqrt{d}}^d \to \mathbb{R}$ of the form $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$ for $\mathbf{v} \in B_1^d$. It is well known that the sample complexity of \mathcal{H} is $\tilde{O}\left(\frac{d}{\epsilon^2}\right)$. In order to prove this via covering numbers one can show that for any choice of $\mathbf{x}_1, \ldots, \mathbf{x}_m \in B_{\sqrt{d}}^d$ the covering number of

$$\mathcal{X} = \{(h(\mathbf{x}_1), \dots, h(\mathbf{x}_m)) : h \in \mathcal{H}\}$$

satisfies $\log(N_2(\mathcal{X}_2, \epsilon)) = \tilde{O}\left(\frac{d}{\epsilon^2}\right)$. This can be proved using standard covering number techniques.

For the sake of illustration we will view \mathcal{H} a composition of two function classes, corresponding to a two layer neural network. Fix $\mathbf{u} \in \mathbb{S}^{d-1}$, let \mathcal{H}_1 be the class of all functions $h : B_{\sqrt{d}}^d \to B_{\sqrt{d}}^d$ of the form $h(\mathbf{x}) = A\mathbf{x}$ for $A \in M_{d,d}$ for A with $||A||_F \leq 1$, and let $\mathcal{H}_2 = \{h_{\mathbf{u}}\}$. Note that $\mathcal{H} = \mathcal{H}_2 \circ \mathcal{H}_1$. Now fix $\mathbf{x}_1, \ldots, \mathbf{x}_m \in B_{\sqrt{d}}^d$ and $\mathbf{u} \in \mathbb{S}^{d-1}$. Consider the sets $\mathcal{X}_1 = \{(h(\mathbf{x}_1), \ldots, h(\mathbf{x}_m) : h \in \mathcal{H}_1\}$ and $\mathcal{X}_2 = \{(h(\mathbf{y}_1), \ldots, h(\mathbf{y}_m)) : (\mathbf{y}_1, \ldots, \mathbf{y}_m) \in \mathcal{X}_1, h \in \mathcal{H}_2\}$

As noted above, $\log(N_2(\mathcal{X}_2, \epsilon)) = \tilde{O}\left(\frac{d}{\epsilon^2}\right)$. Suppose now that we want to prove this in an inductive way. Can we guarantee that $\log(N_2(\mathcal{X}_2, \epsilon)) = \tilde{O}\left(\frac{d}{\epsilon^2}\right)$ via a bound on the ℓ^2 covering numbers of \mathcal{X}_1 without specific assumptions on the structure of \mathcal{X}_1 ? (remember that we want an inductive argument that will work for neural networks, in which case it is not clear what further assumptions we can make)? As Claim 1 below shows, without assumptions beyond the fact that $\mathcal{X}_1 \subset \left(\frac{B^d}{\sqrt{d}}\right)^m$, the best bound we can derive is $N_2(\mathcal{X}_2, \epsilon) \leq N_2(\mathcal{X}_1, \epsilon)$. This is not enough as for $\epsilon = 1/4$, Claim 2 below chows that $\log(N_2(\mathcal{X}_1, \epsilon))$ may be as large as $\Omega(d^2)$, thus the best bound we can get is $\log(N_2(\mathcal{X}_2, 1/4)) = O(d^2)$.

Claim 1. There is a set
$$\mathcal{X} \subset \left(B_{\sqrt{d}}^d\right)^m$$
 such that $N_2(\mathcal{X}, \epsilon) = N_2(\mathbf{u}^\top \mathcal{X}, \epsilon)$

097 098

099

100 101 102 *Proof.* (sketch) It is not hard to verify that for $\mathcal{X} = \left\{ \left(a_1 \sqrt{d} \mathbf{u}, \dots, a_m \sqrt{d} \mathbf{u} \right) : a_i \in \{\pm 1\} \right\}$. We have $N_2(\mathcal{X}, \epsilon) = N_2(\{\pm \sqrt{d}\}^m, \epsilon) = N_2(\mathbf{u}^\top \mathcal{X}, \epsilon)$

Claim 2. For $\epsilon \leq 1$, m = d and $\mathbf{x}_i = \sqrt{d} \mathbf{e}_i$ we have $N_2(\mathcal{X}_2 \epsilon) = \Omega(d^2)$

103 104 105

Proof. We have
$$\mathcal{X}_1 = \left\{ (\mathbf{a}_1, \dots, \mathbf{a}_d) : \sum_{i=1}^d \|\mathbf{a}_i\|^2 = d \right\}$$
. Thus, $N_2(\mathcal{X}_1, \epsilon) = N_2\left(B_{\sqrt{d}}^{d^2}, \sqrt{d}\epsilon\right) = N_2\left(B_1^{d^2}, \epsilon\right) \stackrel{\text{Lemma 2.2}}{=} \Omega(d^2)$

108 2 MULTICOVER

110 2.1 NOTATION

We denote by $\mathbf{x}_1 \mathbf{x}_2$ the elementwise product of two vectors $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$. We denote by $\mathbf{e}_1 \dots, \mathbf{e}_d$ the standard basis of \mathbb{R}^d and by $\{E_{ij}\}_{1 \le i,j \le d}$ the standard basis of the space of $d \times d$ matrices. We will use $\|\cdot\|$ to denote the standard Euclidean norm for vector and the spectral norms for matrices. $\|\cdot\|_F$ will be used for the Frobenius norm of matrices. B_r^d will stand for the Euclidean ball of radius r in \mathbb{R}^d . We will use \lesssim to denote inequality up to a constant.

117

118 2.2 BASIC DEFINITIONS

119 Denote by \mathcal{R}^d the convex set of $d \times d$ PSD matrices. Let $\mathcal{S} \subset \mathcal{R}^d$, we say that \mathcal{S} is *nice* if $\forall R \in \mathcal{S}$ 120 and $W \in \mathbb{R}^{d \times d}$ with $||W|| \leq 1$ then $W^{\top} R W \in S$. We denote by $\mathcal{R}_t^d = \{R \in \mathcal{R}^d : \operatorname{Tr}(R) \leq t\}$ the 121 corresponding nice set. We denote the inner product, the norm, and the metric induced by $R \in \mathcal{R}^d$ 122 on \mathbb{R}^d by $\langle \mathbf{x}, \mathbf{y} \rangle_R = \langle \mathbf{x}, R \mathbf{y} \rangle$, $\|\mathbf{x}\|_R = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_R}$ and $d_R(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_R$. Fix $\mathcal{X} \subset \mathbb{R}^d$ and let 123 $\varepsilon > 0$. A set $\check{\mathcal{X}} \subset \mathbb{R}^d$ is an ε -cover of \mathcal{X} w.r.t. a metric d on \mathbb{R}^d if for every $\mathbf{x} \in \mathcal{X}$ there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ 124 such that $d(\mathbf{x}, \check{\mathbf{x}}) \leq \varepsilon$. A set $\check{\mathcal{X}} \subset \mathbb{R}^d$ is an ε -multicover of \mathcal{X} w.r.t. a nice set \mathcal{S} if for any $R \in \mathcal{S}$ and 125 every $\mathbf{x} \in \mathcal{X}$ there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\mathbf{x} - \check{\mathbf{x}}\|_R \leq \sqrt{\mathrm{Tr}(R)\varepsilon}$. Equivalently, for any $R \in \mathcal{S}, \check{\mathcal{X}}$ is an 126 ε -cover of \mathcal{X} w.r.t. d_R . The ε -multicovering-number of \mathcal{X} , w.r.t. \mathcal{S} , and denoted by $M_{\mathcal{S}}(\mathcal{X}, \varepsilon)$ is the 127 minimal size of an ε -multicover of \mathcal{X} w.r.t. \mathcal{S} . Likewise, the ε -covering-number of \mathcal{X} w.r.t. a metric 128 d, denoted by $N_d(\mathcal{X},\varepsilon)$, is the minimal size of an ε -cover of \mathcal{X} w.r.t. d. We will use $N_p(\mathcal{X},\varepsilon)$ when 129 the metric is $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_p$ and $N_R(\mathcal{X}, \varepsilon)$ when the metric is d_R for $R \in \mathcal{R}^d$. 130

Note that if S is a nice set, $A \in S$ and $B \leq A$, then w.l.o.g. we may assume $B \in \mathcal{R}$. This is because for every $x \in \mathbb{R}^d x^\top Bx \leq x^\top Ax$, thereby $||x||_B \leq ||x||_A$. i.e. adding all PSD matrices of lower PSD order to S keeps the ε -multicover w.r.t. S valid. On the other hand, adding matrices only adds constraints and therefore cannot decrease the multicovering number. Overall we get $M_S(\mathcal{X}, \varepsilon) = M_{S \cup \{B\}}(\mathcal{X}, \varepsilon)$ for any \mathcal{X} and ε

We will also use the notion of packing. We say that X̃ ⊂ X is an ε-packing of X w.r.t a metric d
on X if d(x, y) ≥ ε for any pair of points x, y ∈ X̃. We denote by P_d(X, ε) the maximal size of
an ε-packing of X. As with covering, we will use P_p(X, ε) when the metric is d(x, y) = ||x - y||_p
and P_R(X, ε) when the metric is d_R for R ∈ R^d. It is well known (e.g. Vershynin (2018)) that

$$P_d(\mathcal{X}, 2\epsilon) \le N_d(\mathcal{X}, \epsilon) \le P_d(\mathcal{X}, \epsilon) \tag{1}$$

2.3 Some preliminary lemmas

Lemma 2.1. Let X_1, \ldots, X_k be independent r.v. with that that are σ -estimators to μ . Then

$$\Pr\left(|\text{median}(X_1,\ldots,X_k)-\mu|>r\sigma\right)<\left(\frac{2}{r}\right)^k$$

Lemma 2.2 (e.g. Vershynin (2018)). For any $\varepsilon \leq M$, $(M/\varepsilon)^d \leq N_2(B_M^d, \varepsilon) \leq (3M/\varepsilon)^d$

150 Lemma 2.3. $P_2(\{\pm 1\}^d, d) \ge e^{d/8}$

2.4 MULTICOVER AND ESTIMATORS

We say that a random variable $X \in \mathbb{R}^d$ is an ε -estimator of $\mathbf{x} \in \mathbb{R}^d$ if for any $\mathbf{u} \in \mathbb{S}^{d-1}$, $\mathbb{E}\langle \mathbf{u}, X - \mathbf{x} \rangle^2 \leq \varepsilon^2$. Equivalently, for any $R \in \mathcal{R}^d$, $\mathbb{E} ||X - \mathbf{x}||_R^2 \leq \operatorname{Tr}(R)\varepsilon^2$. We say that X is unbiased if $\mathbb{E}X = \mathbf{x}$.

Lemma 2.4. Let $\mathcal{X} \subset \mathbb{R}^d$. A set $\tilde{\mathcal{X}} \subset \mathbb{R}^d$ is an ε -multicover of \mathcal{X} w.r.t. \mathcal{R}^d if and only if for any $\mathbf{x} \in \mathcal{X}$ there is a random vector $X \in \tilde{\mathcal{X}}$ that is an ε -estimator of \mathbf{x} .

160

140 141 142

143 144

145

146 147 148

152

153

161 Proof. Write $\tilde{\mathcal{X}} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$. Suppose that $\tilde{\mathcal{X}}$ is a ε -multicover and let $\mathbf{x} \in \mathcal{X}$. It is enough to show that there is a r.v. X whose range is $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ such that for any $R \in \mathcal{R}_1^d, \mathbb{E} ||X - \mathbf{x}||_R^2 \leq \varepsilon^2$.

¹⁶² Such a r.v. exists if and only if

$$\min_{\lambda \in \Delta^{T-1}} \max_{R \in \mathcal{R}_1^d} \sum_{i=1}^T \lambda_i \| \mathbf{x}_i - \mathbf{x} \|_R^2 \le \varepsilon^2$$

since the objective $\sum_{i=1}^{T} \lambda_i ||\mathbf{x}_i - \mathbf{x}||_R^2 = \sum_{i=1}^{T} \lambda_i (\mathbf{x}_i - \mathbf{x})^\top R(\mathbf{x}_i - \mathbf{x})$ is bi-linear in λ and R, and since Δ^{T-1} and \mathcal{R}_1^d are both convex and compact, we can apply the minmax theorem to conclude that a r.v. X as described above exists if and only if

$$\max_{R \in \mathcal{R}_1^d} \min_{\lambda \in \Delta^{T-1}} \sum_{i=1}^T \lambda_i \|\mathbf{x}_i - \mathbf{x}\|_R^2 \le \varepsilon^2$$

which is equivalent to

$$\max_{R \in \mathcal{R}_1^d} \min_{i \in [T]} \|\mathbf{x}_i - \mathbf{x}\|_R \le \varepsilon$$

Which is indeed the case as $\check{\mathcal{X}}$ is an ε -multicover on \mathcal{X} .

Suppose now that for any $\mathbf{x} \in \mathcal{X}$ there is a r.v. X whose range is $\check{\mathcal{X}}$ such that for any $R \in \mathcal{R}_1^d$, $\mathbb{E} \|X - \mathbf{x}\|_R^2 \leq \varepsilon^2$. This implies that for any $\mathbf{x} \in \mathcal{X}$ and any $R \in \mathcal{R}_1^d$ there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\check{\mathbf{x}} - \mathbf{x}\|_R \leq \varepsilon$. This implies that $\check{\mathcal{X}}$ is an ε -multicover of \mathcal{X} .

2.5 The multicovering-number of an Euclidean Ball

Lemma 2.5. For the ball $B_M^d = \{ \mathbf{x} \in \mathbb{R}^d | \|\mathbf{x}\|_2 \leq M \}$ and $\varepsilon \leq M$ we have

$$2^{\min(d,\lfloor (M/2\varepsilon)^2 \rfloor)} \le M_{\mathcal{R}_1^d}(B_M^d,\varepsilon) \le \min\left((4d^2 \lceil M \rceil + 6d)^{\lceil \frac{2M^2 + \frac{1}{4}}{\varepsilon^2} \rceil}, (3M/\varepsilon)^d \right)$$

The idea behind the proof is constructing a sparse covering set by picking a convex hull that covers the ball, then using averaging to make the sparse cover k-sparse, in the spirit of Maury's lemma (Pisier, 1980-1981).

2.6 MULTICOVER CALCULUS

Lemma 2.6. Let $S \subset \mathbb{R}^d$ be a nice set, then:

- 1. For $\mathcal{X} \subset \mathbb{R}^{d_1}$ and a $d_2 \times d_1$ matrix A we have $M_{\mathcal{S}}(A\mathcal{X}, ||A||_{\mathcal{E}}) \leq M(\mathcal{X}, \varepsilon)$
- 2. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^d$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have $M_{\mathcal{S}}(\sum_{i=1}^n \mathcal{X}_i, \sum_{i=1}^n \varepsilon_i) \leq \prod_{i=1}^n M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon_i)$
- 3. For $\mathcal{X} \subset \mathbb{R}^d$, $\varepsilon > 0$, orthonormal matrix U and $\mathbf{b} \in \mathbb{R}^d$ we have $M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b}, \varepsilon) = M_{\mathcal{S}}(\mathcal{X}, \varepsilon)$

4. For
$$\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^d$$
 and $\varepsilon > 0$ we have $M_{\mathcal{S}}(\cup_{i=1}^n \mathcal{X}_i, \varepsilon) \leq \sum_{i=1}^n M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon)$

5. For
$$S = \mathcal{R}_1^d$$
 and $\mathcal{X}_i \subset [-M_i, M_i]^d$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have

$$M_{\mathcal{R}_1^d}\left(\prod_{i=1}^n \mathcal{X}_i, \prod_{i=1}^n (M_i + \varepsilon_i) - \prod_{i=1}^n M_i\right) \le \prod_{i=1}^n M_{\mathcal{R}_1^d}(\mathcal{X}_i, \varepsilon_i)$$

6. If for any maximal $R \in S$ (w.r.t. PSD order), $Tr(R) \ge 1$. Fix $\mathcal{X} \subset B_M^{d_1}$, and $\mathcal{L} \subset \mathbb{R}^{d_2, d_1}$ matrices with spectral norm $\le r$. Denote $\|A\|_S := \min\{t > 0 : \frac{1}{t}A \in S\}$, then

$$M_{\mathcal{S}}\left(\mathcal{LX}, \varepsilon_2 \sqrt{2r^2 + 2\varepsilon_1^2 \|I_{d_1}\|_{\mathcal{S}}} + \varepsilon_1 M\right) \le M_{\mathcal{R}_1^d}(\mathcal{L}, \varepsilon_1) \cdot M_{\mathcal{S}}(\mathcal{X}, \varepsilon_2)$$

We next prove each item separately.

Proof. (of item 1.) Let $\check{\mathcal{X}}$ be an ε -multicover of \mathcal{X} w.r.t. \mathcal{S} . It is enough to show that $A\check{\mathcal{X}}$ is an $(||A||\varepsilon)$ -multicover of $A\mathcal{X}$. Fix $\mathcal{X} \in \mathcal{X}$ and a PSD matrix $R \in \mathcal{S}$. We need to show that there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $||A\mathbf{x} - A\check{\mathbf{x}}||_R^2 \leq \operatorname{Tr}(R)||A||^2\varepsilon^2$. Now, for any $\check{\mathbf{x}}$ we have

$$\|A\mathbf{x} - A\check{\mathbf{x}}\|_{R}^{2} = \|A\|^{2} (\mathbf{x} - \check{\mathbf{x}})^{\top} \frac{A^{\top}}{\|A\|} R \frac{A}{\|A\|} (\mathbf{x} - \check{\mathbf{x}}) = \|A\|^{2} \|\mathbf{x} - \check{\mathbf{x}}\|_{\frac{A^{\top}}{\|A\|} R \frac{A}{\|A\|}}^{2}$$

Finally, since $\check{\mathcal{X}}$ is an ε -multicover w.r.t. \mathcal{S} , there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\mathbf{x} - \check{\mathbf{x}}\|_{\frac{A^{\top}}{\|A\|}}^2 R_{\frac{A}{\|A\|}} \leq \operatorname{Tr}(\frac{1}{\|A\|}A^{\top}R_{\frac{1}{\|A\|}}A)\varepsilon^2 \leq \operatorname{Tr}(R)\varepsilon^2$. Therefore overall

$$\|A\mathbf{x} - A\check{\mathbf{x}}\|_{R}^{2} = \|A\|^{2} \|\mathbf{x} - \check{\mathbf{x}}\|_{\|A\|}^{2} R_{\frac{A}{\|A\|}} \leq \operatorname{Tr}(R) \|A\|^{2} \varepsilon^{2}$$

Proof. (of item 5.) We first prove the item for n = 2. We will then show that the general case follows by induction. In the proof of this item we will denote by $A \circ B$ the elementwise product of two $d \times d$ matrices, and by $\operatorname{diag}(A)$ the diagonal matrix obtained by zeroing the non-diagonal entries of A.

Let $\check{\mathcal{X}}_i$ be an ε_i -multicover of \mathcal{X}_i w.r.t. \mathcal{R}_1^d . Fix $\mathbf{x}_i \in \mathcal{X}_i$ and a PSD matrix $R \ge 0$ with $\operatorname{Tr}(R) \le 1$. It is enough to show that there is $\check{\mathbf{x}}_i \in \check{\mathcal{X}}_i$ with $\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \le M_1\varepsilon_2 + M_2\varepsilon_1 + \varepsilon_1\varepsilon_2$. We have

$$\begin{aligned} \|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R &\leq \|\mathbf{x}_1\mathbf{x}_2 - \mathbf{x}_1\check{\mathbf{x}}_2\|_R + \|\mathbf{x}_1\check{\mathbf{x}}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \\ &= \|\mathbf{x}_2 - \check{\mathbf{x}}_2\|_{R\circ\mathbf{x}_1\mathbf{x}_1^\top} + \|\mathbf{x}_1 - \check{\mathbf{x}}_1\|_{R\circ\check{\mathbf{x}}_2\check{\mathbf{x}}_2} \end{aligned}$$

Now, $\operatorname{Tr}(R \circ \check{\mathbf{x}}_{2}\check{\mathbf{x}}_{2}^{\top}) = \|\check{\mathbf{x}}_{2}\|^{2}_{\operatorname{diag}(R)}$. Thus, we can choose $\check{\mathbf{x}}_{1}$ such that $\|\mathbf{x}_{1} - \check{\mathbf{x}}_{1}\|_{R \circ \check{\mathbf{x}}_{2}\check{\mathbf{x}}_{2}^{\top}} \leq \|\check{\mathbf{x}}_{2}\|_{\operatorname{diag}(R)} \varepsilon_{1}$. We get for any 0

$$\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \stackrel{(*)}{\leq} \|\mathbf{x}_2 - \check{\mathbf{x}}_2\|_{\frac{1}{p}R \circ \mathbf{x}_1\mathbf{x}_1^\top + \frac{1}{1-p}\operatorname{diag}(\varepsilon_1^2 R)} + M_2\varepsilon_1$$

Where (*) follows from straight-forward calculations that appear fully in the appendix version of this proof. Now, we can choose $\check{\mathbf{x}}_2 \in \check{\mathcal{X}}_2$ with

$$\begin{aligned} \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{\frac{1}{p}R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top} + \frac{1}{1-p}\operatorname{diag}(\varepsilon_{1}^{2}R)} &\leq \varepsilon_{2}\sqrt{\operatorname{Tr}\left(\frac{1}{p}R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top} + \frac{1}{1-p}\operatorname{diag}(\varepsilon_{1}^{2}R)\right)} \\ &\leq \varepsilon_{2}\sqrt{M_{1}^{2}/p + \varepsilon_{1}^{2}/(1-p)} \end{aligned}$$

for $p = M_1/(M_1 + \varepsilon_1)$ we get that

$$\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \le \varepsilon_2(M_1 + \varepsilon_1) + M_2\varepsilon_1 = M_1\varepsilon_2 + M_2\varepsilon_1 + \varepsilon_1\varepsilon_2$$

We next consider n > 2 and conclude the proof by induction. Denote $\mathcal{X}'_2 = \prod_{i=2}^n \mathcal{X}_i, M'_2 = \prod_{i=2}^n M_i$ and $\varepsilon'_2 = \prod_{i=2}^n (M_i + \varepsilon_i) - \prod_{i=2}^n M_i$.

By the induction hypothesis we have $M(\mathcal{X}'_2, \varepsilon'_2) \leq \prod_{i=2}^n M(\mathcal{X}_i, \varepsilon_i)$. By the case n = 2 we have $M(\mathcal{X}_1, \mathcal{X}'_2, M_1 \varepsilon'_2 + M'_2 \varepsilon_1 + \varepsilon_1 \varepsilon'_2) \leq M(\mathcal{X}_1, \varepsilon_1) M(\mathcal{X}'_2, \varepsilon'_2)$

$$= M(\mathcal{X}_1, \varepsilon_1) \prod_{i=2}^n M(\mathcal{X}_i, \varepsilon_i) = \prod_{i=2}^n M(\mathcal{X}_i, \varepsilon_i)$$

this concludes the proof as $\mathcal{X}_1\mathcal{X}_2' = \prod_{i=1}^n \mathcal{X}_i$ and

$$M_{1}\varepsilon_{2}' + M_{2}'\varepsilon_{1} + \varepsilon_{1}\varepsilon_{2}' = (M_{1} + \varepsilon_{1})\left(\prod_{i=2}^{n}(M_{i} + \varepsilon_{i}) - \prod_{i=2}^{n}M_{i}\right) + \varepsilon_{1}\prod_{i=2}^{n}M_{i}$$
$$= \prod_{i=1}^{n}(M_{i} + \varepsilon_{i}) - \prod_{i=1}^{n}M_{i}$$

r		
L		
L		
L		

$$\begin{split} \|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_{R} &\leq \|W\mathbf{x} - W\check{\mathbf{x}}\|_{R} + \|W\check{\mathbf{x}} - \check{W}\check{\mathbf{x}}\|_{R} = \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \|(W - \check{W})\check{\mathbf{x}}\|_{R} \\ &= \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \sqrt{\check{\mathbf{x}}^{\top}(W - \check{W})^{\top}R(W - \check{W})\check{\mathbf{x}}} \end{split}$$

Now, $(W_1, W_2) \mapsto \check{\mathbf{x}}^\top W_1^\top R W_2 \check{\mathbf{x}}$ is a symmetric and positive bi-linear form on the space of $d_2 \times d_1$ matrices of trace

$$\sum_{i=1}^{d_2} \sum_{j=1}^{d_1} \check{\mathbf{x}}^\top E_{ij}^\top R E_{ij} \check{\mathbf{x}} = \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} (\check{x}_j \mathbf{e}_i)^\top R(\check{x}_j \mathbf{e}_i) = \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} \check{x}_j^2 R_{ii} = \operatorname{Tr}(R) \|\check{\mathbf{x}}\|^2$$

Thus, there is $\check{W} \in \check{\mathcal{L}}$ such that $\check{\mathbf{x}}^{\top} (W - \check{W})^{\top} R (W - \check{W}) \check{\mathbf{x}} \leq \operatorname{Tr}(R) \|\check{\mathbf{x}}\|^2 \varepsilon_1^2$. For this \check{W} we have

$$\|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_{R} \stackrel{(*)}{\leq} \sqrt{2}\|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW + \varepsilon_{1}^{2}\mathrm{Tr}(R)I_{d_{1}}} + \varepsilon_{1}\sqrt{\mathrm{Tr}(R)}M$$

291 Where (*) follows from simple calculations, that appear fully in the appendix version of this proof. 292 Thus, it is possible to choose $\check{\mathbf{x}} \in \check{\mathcal{X}}$ s.t. $\|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW + \varepsilon_1^2 I_{d_1}} \le \varepsilon_2 \sqrt{\|W^{\top}RW + \varepsilon_1^2 \mathrm{Tr}(R)I_{d_1}\|_{\mathcal{S}}}$. 293 Finally, $\|W^{\top}RW + \varepsilon_1^2 \mathrm{Tr}(R)I_{d_1}\|_{\mathcal{S}} \le r^2 + \varepsilon_1^2 \mathrm{Tr}(R)\|I_{d_1}\|_{\mathcal{S}} \le Tr(R)r^2 + \varepsilon_1^2 \mathrm{Tr}(R)\|I_{d_1}\|_{\mathcal{S}}$

Lemma 2.7. Fix $\mathcal{X} \subset \mathbb{R}^d$, $\varepsilon > 0$ and r > 2. It holds that $N_{\infty}(\mathcal{X}, r\varepsilon) \leq \left(M_{\mathcal{R}_1^d}(\mathcal{X}, \varepsilon)\right)^{\left\lceil \log_{r/2}(d) \right\rceil}$.

Proof. Let $\check{\mathcal{X}}$ be an ε -multicover of \mathcal{X} of size $M(\mathcal{X}, \varepsilon)$. By lemma 2.4 for any $\mathbf{x} \in \mathcal{X}$ there is a distribution $\mathcal{D}_{\mathbf{x}}$ on $\check{\mathcal{X}}$ such that if $X \sim \mathcal{D}_{\mathbf{x}}$ then X is an ε -estimator of \mathbf{x} . In particular, for any coordinate $i \in [d]$ we have

$$\mathbb{E}_X (X_i - x_i)^2 = \mathbb{E}_X (X - \mathbf{x})^\top E_{ii} (X - \mathbf{x}) \le \varepsilon^2$$

Denote $k = \left| \log_{r/2}(d) \right|$. By the above equation and lemma 2.1 we conclude that if $X^1, \ldots, X^k \sim \mathcal{D}_{\mathbf{x}}$ then for every $i \in [d]$

$$\Pr\left(\exists i \in [d] \text{ s.t. } |\text{median}(X_i^1, \dots, X_i^k) - x_i| > r\varepsilon\right) < d\left(\frac{2}{r}\right)^k \le 1$$

in particular, there exists $\mathbf{x}^1, \ldots, \mathbf{x}^k \in \check{\mathcal{X}}$ such that for any $i \in [d]$, $|\text{median}(x_i^1, \ldots, x_i^k) - x_i| \leq r\varepsilon$. This implies that

is an ε -cover of \mathcal{X} w.r.t. the ℓ^{∞} norm. This concludes the proof as $|\text{median}(\check{\mathcal{X}}^k)| \leq |\check{\mathcal{X}}|^k$

$$\mathrm{median}(\check{\mathcal{X}}^k) := \left\{ \left(\mathrm{median}(x_1^1, \dots, x_1^k), \dots, \mathrm{median}(x_d^1, \dots, x_d^k) \right) : \mathbf{x}^1, \dots, \mathbf{x}^k \in \check{\mathcal{X}} \right\}$$

3 MULTICOVER FOR NEURAL NETWORKS SAMPLE COMPLEXITY

3.1 MULTICOVER FOR SEQUENCE OF VECTORS

We denote by $\mathbb{R}^{d,m}$ the vector space of sequences $\mathbf{x} = (\mathbf{x}^1, \dots, \mathbf{x}^m)$ of m vectors in \mathbb{R}^d . We next extend the notion of multicover, as well as multicover calculus, to subsets $\mathbb{R}^{d,m}$. This extension is useful for sample complexity analysis via multicover.

Denote by $\mathbb{S}^{d,m}$ the collection of all sequences $(\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathbb{R}^{d,m}$ with $\sum_{i=1}^m ||\mathbf{u}_i||^2 = 1$. We also denote by $\mathcal{R}^{d,m}$ the convex set of sequences $R = (R_1, \dots, R_m)$ of $m \ d \times d$ PSD matrices.

We denote $\operatorname{Tr}(R) = \sum_{i=1}^{m} \operatorname{Tr}(R_i)$ and $\mathcal{R}_t^{d,m} = \{R \in \mathcal{R}^{d,m} : \operatorname{Tr}(R) \leq t\}$. We say that a set $\mathcal{S} \subset \mathcal{R}^{d,m}$ is *nice*, if $\forall R \in \mathcal{S}, i \in [m]$ and $W \in \mathbb{R}^{d \times d,m}$ with $||W_i|| \leq 1$ then $W^{\top}RW = (W_1^{\top}R_1W_1, \dots, W_m^{\top}R_mW_m) \in \mathcal{S}$. We denote the inner product, norm and metric induced by $R \in \mathcal{R}^{d,m}$ on $\mathbb{R}^{d,m}$ by $\langle \mathbf{x}, \mathbf{y} \rangle_R = \sum_{i=1}^{m} \langle \mathbf{x}^i, R_i \mathbf{y}^i \rangle$, $||\mathbf{x}||_R = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_R}$ and $d_R(\mathbf{x}, \mathbf{y}) = ||\mathbf{x} - \mathbf{y}||_R$. A set $\check{\mathcal{X}} \subset \mathbb{R}^{d,m}$ is an ε -multicover of \mathcal{X} w.r.t. a nice set \mathcal{S} if for any $R \in \mathcal{S}$ and every $\mathbf{x} \in \mathcal{X}$. 324 325 326 327 328 there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\mathbf{x} - \check{\mathbf{x}}\|_R \leq \sqrt{\operatorname{Tr}(R)}\varepsilon$. Equivalently, for any $R \in \mathcal{R}_1^{d,m} \check{\mathcal{X}}$ is an ε -cover of \mathcal{X} w.r.t. d_R . The ε -multicovering-number of \mathcal{X} , denoted by $M(\mathcal{X}, \varepsilon)$ is the minimal size of an 330 ε -multicover of \mathcal{X} . 332 We will use $N_R(\mathcal{X},\varepsilon)$ for the covering number w.r.t. the metric d_R , $N_{\infty}(\mathcal{X},\varepsilon)$ when the 333 metric is $d(\mathbf{x}, \mathbf{y}) = \max_{j \in [m]} \|\mathbf{x}^j - \mathbf{y}^j\|_{\infty}$ and $N_2(\mathcal{X}, \varepsilon)$ when the metric is $d(\mathbf{x}, \mathbf{y}) =$ 334 $\sqrt{\frac{1}{m}\sum_{j=1}^m \|\mathbf{x}^j - \mathbf{y}^j\|_2^2}.$ 335 336 We say that a random variable $X \in \mathbb{R}^{d,m}$ is an ε -estimator of $\mathbf{x} \in \mathbb{R}^{d,m}$ if for any $\mathbf{u} \in \mathbb{S}^{d,m}$, we 337 have $\sum_{i=1}^{m} \mathbb{E} \langle \mathbf{u}^{j}, X^{j} - \mathbf{x}^{j} \rangle^{2} \leq \varepsilon^{2}$. Equivalently, for any $R \in \mathcal{R}^{d,m}$, $\mathbb{E} \| X - \mathbf{x} \|_{R}^{2} \leq \operatorname{Tr}(R) \varepsilon^{2}$. 338 339 We next generalize Lemmas 2.4, 2.6 and 2.7 to the extended definition of multicover. The proofs 340 of the generalized lemmas are similar to the proofs of the original lemmas and are deffered to the 341 appendix, similarly to the other absent proofs. 342 **Lemma 3.1.** Let $\mathcal{X} \subset \mathbb{R}^{d,m}$. A set $\check{\mathcal{X}} \subset \mathbb{R}^{d,m}$ is an ε -multicover of \mathcal{X} w.r.t. $\mathcal{R}^{d,m}$ if and only if for 343 any $\mathbf{x} \in \mathcal{X}$ there is a random vector $X \in \check{\mathcal{X}}$ that is an ε -estimator of \mathbf{x} . 344 1. For $\mathcal{X} \subset \mathbb{R}^{d_1,m}$ and $A \in \mathbb{R}^{d_2 \times d_1}$ we have $M_{\mathcal{S}}(A\mathcal{X}, ||A||_{\mathcal{E}}) \leq M_{\mathcal{S}}(\mathcal{X}, \varepsilon)$ Lemma 3.2. 345 346 2. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^{d,m}$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have $M_{\mathcal{S}}(\sum_{i=1}^n \mathcal{X}_i, \sum_{i=1}^n \varepsilon_i) \leq \infty$ 347 $\prod_{i=1}^{n} M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon_i)$ 348 3. For $\mathcal{X} \subset \mathbb{R}^{d,m}$, $\varepsilon > 0$ and $\mathbf{b} \in \mathbb{R}^{d,m}$ we have $M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b}, \varepsilon) = M_{\mathcal{S}}(\mathcal{X}, \varepsilon)$ 349 350 4. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^{d,m}$ and $\varepsilon > 0$ we have $M_{\mathcal{S}}(\bigcup_{i=1}^n \mathcal{X}_i, \varepsilon) \leq \sum_{i=1}^n M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon)$ 351 352 5. For $S = \mathcal{R}_1^{d,m} \mathcal{X}_i \subset [-M_i, M_i]^{d,m}$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have 353 354 $M_{\mathcal{R}_1^{d,m}}\left(\prod_{i=1}^n \mathcal{X}_i, \prod_{i=1}^n (M_i + \varepsilon_i) - \prod_{i=1}^n M_i\right) \le \prod_{i=1}^n M_{\mathcal{R}_1^{d,m}}(\mathcal{X}_i, \varepsilon_i)$ 355 356 357 6. For $\mathcal{S} = \mathcal{R}_1^{d,m}$, fix $\mathcal{X} \subset B_M^{d_1,m}$, $\mathcal{L} \subset \mathbb{R}^{d_2,d_1}$ matrices with spectral norm $\leq r$. Then, $M_{\mathcal{S}}\left(\mathcal{LX},\varepsilon_{2}\sqrt{2r^{2}+2\varepsilon_{1}^{2}d_{1}}+\varepsilon_{1}M\right) \leq M_{\mathcal{R}_{1}^{d,m}}(\mathcal{L},\varepsilon_{1})\cdot M_{\mathcal{S}}(\mathcal{X},\varepsilon_{2})$ 360 361 362 **Lemma 3.3.** Fix $\mathcal{X} \subset \mathbb{R}^{d,m}$, $\varepsilon > 0$ and r > 2. 363 Then $N_{\infty}(\mathcal{X}, r\varepsilon) \leq \left(M_{\mathcal{R}_{1}^{d,m}}(\mathcal{X}, \varepsilon)\right)^{\left\lceil \log_{r/2}(dm) \right\rceil}$. 364 365 366 3.1.1 STRONGLY BOUNDED ACTIVATION 367 368 In this section we will develop tools to calculate $M_{\mathcal{R}_1^{d,m}}(\rho(\mathcal{X}),\varepsilon)$ for a smooth enough ρ . For the 369 sake of cleanliness we will denote $M(\cdot, \cdot) := M_{\mathcal{R}^{d,m}_{1}}(\cdot, \cdot)$. The smoothness requirements are given 370 in the following definition. 371 **Definition 3.4.** A function $\rho : \mathbb{R} \to \mathbb{R}$ is B-strongly-bounded if for all $n \ge 1$, $\|\rho^{(n)}\|_{\infty} \le n!B^n$. 372 373 Likewise, ρ is strongly-bounded if it is B-strongly-bounded for some B 374

As shown in Daniely & Granot (2019) the ReLU-like function $\log(1 + e^x)$ is strongly bounded, as well as the sigmoid function $\frac{e^x}{1+e^x}$. It is also shown in Daniely & Granot (2019) that

³⁷⁷

¹This claim can be generalized to a more general S. We present the case $S = \mathcal{R}_1^{d,m}$ for simplicity.

Fact 3.5. If ρ is *B*-strongly-bounded then ρ is analytic and its Taylor coefficients around any point are bounded by B^n for any $n \ge 1$.

We will utilize this fact in order to calculate the effect of a non-linearity on the multicovering number. **Lemma 3.6** (β -Swish Activation Ramachandran et al. (2017)). For a constant $\beta \ge 0$, the function $\frac{x}{1+e^{-\beta x}}$ is strongly-bounded

Lemma 3.7 (Hyperbolic Tangent). The function $\frac{e^{2x}-1}{e^{2x}+1}$ is strongly-bounded

In order to analyze $M(\rho(\mathcal{X}), \varepsilon)$ for a strongly bounded ρ , we first analyze $M(p(\mathcal{X}), \varepsilon)$ for a polynomial p, and then utilize fact 3.5.

Lemma 3.8. Let $p(x) = \sum_{i=0}^{k} a_i X^i$ be a polynomial with $|a_i| \leq B^i$ and suppose that $\mathcal{X} \subset \left[-\frac{1}{8B}, \frac{1}{8B}\right]^{d,m}$. Then, for any Let $0 < \varepsilon \leq 1$, $M\left(p(\mathcal{X}), \varepsilon\right) \leq \left(M\left(\mathcal{X}, \frac{\varepsilon}{8B}\right)\right)^{\frac{k(k+1)}{2}}$

We are now ready to present our main tool for analysing $M(\rho(\mathcal{X}), \varepsilon)$ for strongly bounded ρ . Lemma 3.9. Let $\mathcal{X} \subset \mathbb{R}^{d,m}$. Let $\rho : \mathbb{R} \to \mathbb{R}$ be B strongly bounded. Then for $1 \ge \varepsilon > 0$,

$$M(\rho(\mathcal{X}), \varepsilon + \sqrt{d} 8^{-(k+1)}) \leq \left(M\left(\mathcal{X}, \frac{1}{32B}\right) \right)^{\lceil \log_2(dm) \rceil} \left(M\left(\mathcal{X}, \frac{\varepsilon}{8B}\right) \right)^{\frac{k(k+1)}{2}}$$

Proof. Assume first that $\mathcal{X} \subset \mathbb{R}^{d,m}$ is contained in an ℓ^{∞} ball of radius $\frac{1}{8B}$. Since multicovering numbers are invariant to translations (i.e. $M(\mathcal{X}, \varepsilon) = M(\mathcal{X} + \mathbf{b}, \varepsilon)$ for any $\mathbf{b} \in \mathbb{R}^{m,d}$), we can assume w.l.o.g. that $\mathcal{X} \subset \left[-\frac{1}{8B}, \frac{1}{8B}\right]^{d,m}$. Let p be the Taylor polynomial of ρ around 0 of degree k and let $r = \rho - p$. We have that for any $x \in \left[-\frac{1}{8B}, \frac{1}{8B}\right], |r(x)| \leq B^{k+1}|x|^{k+1} \leq 8^{-(k+1)}$. Thus $\{0\}$ is an $\left(\sqrt{d}8^{-(k+1)}\right)$ -multicover of $r(\mathcal{X})$. Indeed, if $R \in \mathcal{R}_1^{d,m}$ and $\mathbf{x} \in r(\mathcal{X})$ then

$$\|\mathbf{x}\|_{R}^{2} = \sum_{i=1}^{m} \left\langle \mathbf{x}^{i}, R_{i} \mathbf{x}^{i} \right\rangle \leq \sum_{i=1}^{m} \operatorname{Tr}(R_{i}) \|\mathbf{x}^{i}\|^{2}$$

$$\leq \sum_{i=1}^{m} \operatorname{Tr}(R_i) d8^{-2(k+1)} = d8^{-2(k+1)} \operatorname{Tr}(R) \leq d8^{-2(k+1)}$$

In particular $M(r(\mathcal{X}), \sqrt{d}8^{-(k+1)}) = 1$. Now, we have

$$\begin{split} M(\rho(\mathcal{X}), \varepsilon + \sqrt{d} 8^{-(k+1)}) & \stackrel{\rho(\mathcal{X}) \subset p(\mathcal{X}) + r(\mathcal{X})}{\leq} & M(p(\mathcal{X}) + r(\mathcal{X}), \varepsilon + \sqrt{d} 8^{-(k+1)}) \\ & \underset{\leq}{\overset{\text{Lemma 3.2}}{\leq}} & M\left(r(\mathcal{X}), \sqrt{d} 8^{-(k+1)}\right) M(p(\mathcal{X}), \varepsilon) \\ & \underset{\leq}{\overset{\text{Lemma 3.8}}{\leq}} & \left(M\left(\mathcal{X}, \frac{\varepsilon}{8B}\right)\right)^{\frac{k(k+1)}{2}} \end{split}$$

Finally, by lemma 3.3, \mathcal{X} is a union of $\left(M\left(\mathcal{X}, \frac{1}{32B}\right)\right)^{\lceil \log_2(dm) \rceil}$ sets \mathcal{X}_i , such that each \mathcal{X}_i is contained in an ℓ^{∞} ball of radius $\frac{1}{8B}$. Applying the above argument to each \mathcal{X}_i implies the lemma. \Box

422 3.2 NEURAL NETWORK SAMPLE COMPLEXITY VIA MULTICOVER

423 3.2.1 MULTICOVER AND SAMPLE COMPLEXITY

Fix an instance space \mathcal{Z} , a label space \mathcal{Y} and a loss $\ell : \mathbb{R}^d \times \mathcal{Y} \to [0, \infty)$. We say that ℓ has some property p (e.g. boundness, Lipschitzness, etc.) if for any $y \in \mathcal{Y}$, $\ell(\cdot, y)$ has the property p. Fix a class \mathcal{H} from \mathcal{Z} to \mathbb{R}^d . For a distribution \mathcal{D} and a sample $S \in (\mathcal{Z} \times \mathcal{Y})^m$ we define the *representativeness* of S as

$$\operatorname{rep}_{\mathcal{D}}(S,\mathcal{H}) = \sup_{h \in \mathcal{H}} \ell_{\mathcal{D}}(h) - \ell_{S}(h)$$

431 Where $\ell_{\mathcal{D}}(h) = \mathbb{E}_{(x,y)\sim\mathcal{D}}\ell(h(x),y)$ and $\ell_{S}(h) = \frac{1}{m}\sum_{i=1}^{m}\ell(h(x_{i}),y_{i})$. We note that if $\operatorname{rep}_{\mathcal{D}}(S,\mathcal{H}) \leq \varepsilon$ then any algorithm that is guaranteed to return a function $\hat{h} \in \mathcal{H}$ will enjoy a generalization bound $\ell_{\mathcal{D}}(h) \leq \ell_{\mathcal{S}}(h) + \varepsilon$. In particular, the ERM algorithm will return a function whose loss is optimal, up to an additive factor of ε .

We will focus on bounds on rep_D(S, \mathcal{H}) when $S \sim \mathcal{D}^m$. To this end, we will rely on the connection between representativeness and the *covering numbers* of \mathcal{H} . For $x_1, \ldots, x_m \in \mathcal{Z}$ we denote

$$\mathcal{H}(x_1,\ldots,x_m) = \{(h(x_1),\ldots,h(x_m)) : h \in \mathcal{H}\}$$

Given a metric d on $\mathbb{R}^{d,m}$ we denote $N_d(\mathcal{H}, m, \epsilon) = \sup_{x_1, \dots, x_m \in \mathcal{Z}} N_d(\mathcal{H}(x_1, \dots, x_m), m, \epsilon)$. Similarly, we denote $M(\mathcal{H}, m, \epsilon) = \sup_{x_1, \dots, x_m \in \mathbb{Z}} M(\mathcal{H}(x_1, \dots, x_m), m, \epsilon).$

Lemma 3.10. (Shalev-Shwartz & Ben-David, 2014) Let $\ell : \mathbb{R}^d \times \mathcal{Y} \to \mathbb{R}$ be B-bounded. Then for any distribution \mathcal{D} on \mathcal{Z}

$$\mathbb{E}_{S \sim \mathcal{D}^m} \operatorname{rep}_{\mathcal{D}}(S, \mathcal{H}) \le B2^{-M+1} + \frac{12B}{\sqrt{m}} \sum_{k=1}^M 2^{-k} \sqrt{\ln\left(N_2(\ell \circ \mathcal{H}, m, B2^{-k})\right)}$$

We conclude with a special case of the above lemma, which will be useful in this paper.

Lemma 3.11. Let $\ell : \mathbb{R}^d \times \mathcal{Y} \to \mathbb{R}$ be L-Lipschitz w.r.t. $\|\cdot\|_{\infty}$ and B-bounded. Assume that for any $\frac{\sqrt{nB}}{\sqrt{m}8L} \leq \varepsilon \leq 1$, $\ln M(\mathcal{H}, m, \varepsilon) \leq \frac{n}{\varepsilon^2}$. Then for any distribution \mathcal{D} on \mathcal{Z}

$$\mathbb{E}_{S \sim \mathcal{D}^m} \operatorname{rep}_{\mathcal{D}}(S, \mathcal{H}) \lesssim \frac{(L+B)\sqrt{n}}{\sqrt{m}} \sqrt{\log(dm)} \log(m)$$

3.2.2 SAMPLE COMPLEXITY OF NEURAL NETWORKS

Fix the instance space to be the ball of radius $\sqrt{d_0}$ in \mathbb{R}^{d_0} (in particular $[-1, 1]^{d_0} \subset \mathcal{X}$). Fix also a *B*-strongly-bounded activation function ρ . Consider the class

$$\mathcal{N}_{r,R}^{\rho}(d_0,\ldots,d_t) = \left\{ W_t \circ \rho \circ W_{t-1} \circ \rho \ldots \circ \rho \circ W_1 : W_i \in M_{d_{i-1}d_i} \| W_i \| \le r, \| W_i \|_F \le R \right\}$$

and more generally, for matrices $W_i^0 \in M_{d_i, d_{i-1}}, i = 1, \ldots, t$ consider

$$\mathcal{H} = \mathcal{N}_{r,R}^{\rho}(W_1^0, \dots, W_t^0) = \left\{ W_t \circ \rho \circ W_{t-1} \circ \rho \dots \circ \rho \circ W_1 : \|W_i - W_i^0\| \le r, \|W_i - W_i^0\|_F \le R \right\}$$

denote $d = \max(d_0, \ldots, d_t)$. We will assume that t, $||W_i^0||$, r are all bounded by some constant C > 0, and will allow hidden constants to depend C. This is motivated by the fact that in practice, $||W_i^0||, r$ are often bounded by small constants. For instance, if the initial weights are sampled form the standard Xavier initialization then $||W_i^0|| \approx \sqrt{2}$ w.h.p. for resnets we have $||W_i^0|| \approx 1$. We will also allow hidden constant to depend on the activation ρ and the depth t.

Theorem 3.12. Let $\ell : \mathbb{R}^d \times \mathcal{Y} \to \mathbb{R}$ be O(1)-Lipschitz w.r.t. $\|\cdot\|_{\infty}$ and O(1)-bounded. Then for any distribution \mathcal{D} on \mathcal{Z}

$$\mathbb{E}_{S \sim \mathcal{D}^m} \operatorname{rep}_{\mathcal{D}}(S, \mathcal{H}) \lesssim \sqrt{\frac{dR^2}{m}} \log^{t+2}(Rdm)$$

The theorem is implied by the following lemma together with lemma 3.11.

Lemma 3.13. For any $0 < \epsilon \leq 1$, $M(\mathcal{H}, m, \epsilon) \lesssim (\log(dm) + \log^2(d/\epsilon))^t \log(dR) \frac{dR^2}{\epsilon^2}$

The proof follows a peeling argument, applying lemmas 3.2, 3.9 and 2.5 inductively, for each layer.

RELATED WORK

In recent years, there has been active work in the area of the sample complexity of neural networks. For the the remaining of this section, we refer the reader to Table 1 to explain the notation used in different works.



5 FUTURE DIRECTIONS

numbers appears in appendix B

Future direction arising from our work Covering-Packing relations for multicover, the (in)Existence of proper multicover (that is, a multicover in which each point is in the class), and the behaviour of multicover w.r.t. Lipschitz functions (specifically, ReLU).

The resemblance of the results of Daniely & Granot (2019) to ours is not coincidental. A full discus-

sion of the connection between the notion of Approximate Description Length and multicovering

²Equivalent to Daniely & Granot (2019), and Vardi et al. (2022) result for 2 layer networks

REFERENCES

542	Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for
543	neural networks. Advances in neural information processing systems, 30, 2017.

- Amit Daniely and Elad Granot. Generalization bounds for neural networks via approximate description length. arXiv preprint arXiv:1910.05697, 2019.
- Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In Conference On Learning Theory, pp. 297-299. PMLR, 2018.
- Daniel Hsu, Ziwei Ji, Matus Telgarsky, and Lan Wang. Generalization bounds via distillation, 2021.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In Conference on Learning Theory, pp. 1376–1401. PMLR, 2015.
- Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. arXiv preprint arXiv:1707.09564, 2017.
- G. Pisier. Remarques sur un résultat non publié de B. Maurey. Séminaire d'Analyse fonctionnelle (dit "Maurey-Schwartz"), pp. 1-12, 1980-1981. URL http://www.numdam.org/item/ SAF_1980-1981____A5_0/. talk:5.
- Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. arXiv preprint arXiv:1710.05941, 2017.
- Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algo-rithms. Cambridge university press, 2014.
- Gal Vardi, Ohad Shamir, and Nati Srebro. The sample complexity of one-hidden-layer neural net-works. Advances in Neural Information Processing Systems, 35:9139–9150, 2022.
 - Roman Vershynin. High-dimensional probability: An introduction with applications in data science, volume 47. Cambridge university press, 2018.

A OMITTED PROOFS

A.1 PRELIMINARY LEMMAS

Lemma A.1 (2.1). Let X_1, \ldots, X_k be independent r.v. with that that are σ -estimators to μ . Then

$$\Pr\left(|\operatorname{median}(X_1,\ldots,X_k)-\mu|>r\sigma\right)<\left(\frac{2}{r}\right)^k$$

Proof of 2.1. We have that $\Pr(|X_i - \mu| > r\sigma) \le \frac{1}{r^2}$. It follows that the probability that $\ge \frac{k}{2}$ of X_1, \ldots, X_k fall outside of the segment $(\mu - r\sigma, \mu + r\sigma)$ is bounded by

 $\binom{k}{\lceil k/2\rceil} \left(\frac{1}{r^2}\right)^{\lceil k/2\rceil} < 2^k \left(\frac{1}{r^2}\right)^{\lceil k/2\rceil} \le \left(\frac{2}{r}\right)^k$

Lemma A.2 (2.3). $P_2(\{\pm 1\}^d, d) \ge e^{d/8}$

 Proof of 2.3. If $\mathbf{y}, \mathbf{x} \in \{\pm 1\}^d$ are two independent uniform vectors then by Hoeffding's bound we have

$$\Pr\left(\|\mathbf{x} - \mathbf{y}\|^2 \le d\right) = \Pr\left(\sum_{i=1}^d \mathbb{1}[x_i \ne y_i] \le d/4\right) \le e^{-d/8}$$

Thus, there are at least $e^{d/8}$ vectors in $\{\pm 1\}^d$ such that the distance between each pair is more than d.

A.2 MULTICOVER FOR VECTORS

Lemma A.3 (2.5). For the ball $B_M^d = \{ \mathbf{x} \in \mathbb{R}^d | \|\mathbf{x}\|_2 \leq M \}$ and $\varepsilon \leq M$ we have

$$2^{\min(d,\lfloor (M/2\varepsilon)^2 \rfloor)} \le M_{\mathcal{R}_1^d}(B_M^d,\varepsilon) \le \min\left((4d^2 \lceil M \rceil + 6d)^{\lceil \frac{2M^2 + \frac{1}{4}}{\varepsilon^2} \rceil}, (3M/\varepsilon)^d \right)$$

Proof of lemma 2.5. We first show that $M_{\mathcal{R}_1^d}(B_M^d,\varepsilon) \leq (4d^2\lceil M\rceil + 6d)^{\lceil\frac{2M^2+\frac{1}{4}}{\varepsilon^2}\rceil}$. By lemma 2.4, it is enough to show that there is a set $\mathcal{X} \subset \mathbb{R}^d$ of size $(4d^2\lceil M\rceil + 6d)^{\lceil\frac{2M^2+\frac{1}{4}}{\varepsilon^2}\rceil}$ such that for every $\mathbf{x} \in B_M^d$ there is a random vector $X \in \mathcal{X}$ which is an ε -estimator of \mathbf{x} . Define

$$\tilde{\mathcal{X}} = \{k\mathbf{e}_i | i \in [d], k \in [-2dM - 1, 2dM + 1] \cap \mathbb{Z}\}$$

 $Let \mathbf{x} \in B_M^d, \text{ we next define a } \sqrt{2M^2 + \frac{1}{4}} \text{ estimator } X \in \tilde{\mathcal{X}} \text{ for } \mathbf{x} \text{: First sample a coordinate } i \text{ w.p.}$ 647 $p_i = \frac{\mathbf{x}_i^2}{2\|\mathbf{x}\|^2} + \frac{1}{2d}, \text{ and let } \check{\mathbf{x}} = \left(\left\lfloor \frac{\mathbf{x}_i}{p_i} \right\rfloor + b \right) e_i \text{ where } b \sim Ber\left(\left\langle \frac{\mathbf{x}_i}{p_i} \right\rangle \right) \text{ and } \left\langle \frac{\mathbf{x}_i}{p_i} \right\rangle \coloneqq \frac{\mathbf{x}_i}{p_i} - \left\lfloor \frac{\mathbf{x}_i}{p_i} \right\rfloor.$

678 679

680

681 682

683 684

685

686 687

688 689

690

691 692 693

649	Note that $\mathbb{E}X = \mathbf{x}$. Fix	$\mathbf{u} \in$	\mathbb{S}^{a-1} . We need to show that $\mathbb{E}\langle \mathbf{u}, X - \mathbf{x} \rangle^2 \leq 2M^2 + \frac{1}{4}$. Indeed,	
650	$\mathbb{E}\langle \mathbf{u}, X - \mathbf{x} \rangle^2$	<	$\mathbb{E}\langle \mathbf{u}, X \rangle^2$	
651		_	(1) (1) (1) (2) (1) (1) (2)	
652		=	$\sum p_i \left(\left\langle \frac{\mathbf{x}_i}{\mathbf{x}_i} \right\rangle \left(\left \frac{\mathbf{x}_i}{\mathbf{x}_i} \right + 1 \right)^2 + \left(1 - \left\langle \frac{\mathbf{x}_i}{\mathbf{x}_i} \right\rangle \right) \left \frac{\mathbf{x}_i}{\mathbf{x}_i} \right ^2 \right) \mathbf{u}_i^2$	
653			$\sum_{i} \sum_{i} \sum_{j} \sum_{i} \left(\left p_{i} \right \right) \left p_{i} \right \right) = \left(\left p_{i} \right \right) \left p_{i} \right \right) = \left(\left p_{i} \right \right) \left p_{i} \right \right) = \left(\left p_{i} \right \right) = \left(\left p_{i} \right \right) \left p_{i} \right \right) = \left(\left p_{i} \right \right) = $	
654			$- \left(\langle \mathbf{x}_{1} \rangle \mathbf{x}_{2} \langle \mathbf{x}_{2} \rangle \mathbf{x}_{2} \rangle \right)$	
655		=	$\sum p_i \left(2 \left\langle \frac{\mathbf{x}_i}{n} \right\rangle \left \frac{\mathbf{x}_i}{n} \right + \left\langle \frac{\mathbf{x}_i}{n} \right\rangle + \left \frac{\mathbf{x}_i}{n} \right \right) \mathbf{u}_i^2$	
656			\overline{i} (\Pi / \Pi] \Pi / \Pi]	
657			$\sum \left(\left(\left \mathbf{x}_i \right\rangle - \left \mathbf{x}_i \right \right)^2 - \left \mathbf{x}_i \right\rangle \left(\left \mathbf{x}_i \right\rangle \right) \right) = 2$	
658		=	$\sum p_i \left(\left(\left\langle \frac{1}{p_i} \right\rangle + \left \frac{1}{p_i} \right \right) + \left\langle \frac{1}{p_i} \right\rangle \left(1 - \left\langle \frac{1}{p_i} \right\rangle \right) \right) \mathbf{u}_i^2$	
659			$i \left($	
660		_	$\sum_{n} \left(\left(\mathbf{x}_{i} \right)^{2} + \left/ \mathbf{x}_{i} \right\rangle \left(\left(\left(\mathbf{x}_{i} \right)^{2} + \left(\mathbf{x}_{i} \right)^{2} \right) \right)^{2} \right)^{2} \right)$	
661		_	$\sum_{i} p_{i} \left(\left(\frac{\overline{p_{i}}}{p_{i}} \right)^{-+} \left(\frac{\overline{p_{i}}}{p_{i}} \right)^{-+} \left(\frac{\overline{p_{i}}}{p_{i}} \right)^{} \right)^{$	
662			$\begin{pmatrix} & & \\ & & \\ & & \end{pmatrix}$	
663		\leq	$\sum p_i \left(\left(\frac{\mathbf{x}_i}{-} \right) + \frac{1}{t} \right) \mathbf{u}_i^2$	(2)
665		_	$\sum_{i} \left(\left(p_i \right) - 4 \right)^{-i}$	
666		,	1_{\parallel} \mathbf{x}_{i}^{2}	
667		\leq	$\frac{1}{4} \ \mathbf{u}\ _{\infty} + \sum_{i} \frac{1}{p_{i}} \mathbf{u}_{i}^{*}$	
668			<i>i i</i> 1	
669		\leq	$\frac{1}{4} + \sum 2 \ \mathbf{x}\ ^2 \mathbf{u}_i^2$	(3)
670			$\frac{4}{i}$	
671		_	$2\ \mathbf{x}\ ^2 + \frac{1}{2}$	
672		_		
673		<	$2M^2 + \frac{1}{2}$	
674			4	
675	Where equation 2 is true	e sinc	$x = x(1-x) \le 1/4$ for any $0 \le x \le 1$ and equation 3 is true by plus	gging
676	in the definition of $n =$	\mathbf{x}_i^2	$\frac{1}{2} + \frac{1}{2}$ and by the fact that u is a unit vector	
677	In the domittion of $p_i =$	$2\ \mathbf{x}\ $	$\ ^2 + 2d$, and by the fact that a is a difference.	

We next construct an ε -estimator by averaging independent copies of X. Let $\tilde{X} = \frac{1}{k} \sum_{i=1}^{k} X_i$ where every X_i is sampled i.i.d. like X. We claim that \tilde{X} is $\sqrt{\frac{2M^2 + \frac{1}{4}}{k}}$ -estimator. Let $\mathbf{u} \in \mathbb{S}^{d-1}$. We have that $\langle \mathbf{u}, \tilde{X} - \mathbf{x} \rangle = \frac{1}{k} \sum_{i=1}^{k} \langle \mathbf{u}, X_i - \mathbf{x} \rangle$ is an average of k i.i.d. r.v. with mean 0 and variance bounded by $2M^2 + \frac{1}{4}$. Thus, $\mathbb{E} \langle \mathbf{u}, \tilde{X} - \mathbf{x} \rangle^2 \leq \frac{2M^2 + \frac{1}{4}}{k}$. Plugging $k = \lceil \frac{2M^2 + \frac{1}{4}}{\varepsilon^2} \rceil$, we get that $\mathbb{E} \langle \mathbf{u}, \tilde{X} - \mathbf{x} \rangle^2 \leq \varepsilon^2$. Note that \tilde{X} gets values in $\mathcal{X} = \left\{ \frac{1}{k} \sum_{i=1}^{k} \mathbf{x}_i : \mathbf{x}_i \in \tilde{\mathcal{X}} \right\}$. By lemma 2.4 we have that \mathcal{X} is a multicover. Finally, $|\mathcal{X}| \leq (4d^2 \lceil M \rceil + 6d)^{\lceil \frac{2M^2 + \frac{1}{4}}{\varepsilon^2} \rceil}$, implying that $M_{\mathcal{R}_1^d}(B_M^d, \varepsilon) \leq (4d^2 \lceil M \rceil + 6d)^{\lceil \frac{2M^2 + \frac{1}{4}}{\varepsilon^2} \rceil}$.

We next show that $M_{\mathcal{R}_1^d}(B_M^d,\varepsilon) \leq (3M/\varepsilon)^d$. Given Lemma 2.2, it is enough to show that $M_{\mathcal{R}_1^d}(B_M^d,\varepsilon) \leq N_2(B_M^d,\varepsilon)$. Indeed, fix $R \in \mathcal{R}_1^d$. We have for any $\mathbf{x}, \check{\mathbf{x}} \in \mathbb{R}^d$

$$\|\mathbf{x} - \check{\mathbf{x}}\|_R^2 = (\mathbf{x} - \check{\mathbf{x}})^\top R(\mathbf{x} - \check{\mathbf{x}}) \le \|R\| \cdot \|\mathbf{x} - \check{\mathbf{x}}\|_2^2 \le \operatorname{Tr}(R) \|\mathbf{x} - \check{\mathbf{x}}\|_2^2 \le \|\mathbf{x} - \check{\mathbf{x}}\|_2^2$$

⁶⁹⁴ Thus, any ε -cover w.r.t. the Euclidean norm is an ε -cover w.r.t. d_R . Since this is true for any $R \in \mathcal{R}_1^d$, ⁶⁹⁵ we have that any ε -cover w.r.t. the Euclidean norm is an ε -multicover. Thus, $M_{\mathcal{R}_1^d}(B_M^d, \varepsilon) \leq N_2(B_M^d, \varepsilon)$

we can use standard upper bounds for covering numbers of sets using volume (Vershynin, 2018) to upper bound the cover of B_M^d with $O\left((3M/\varepsilon)^d\right)$ balls of radius ε in ℓ_2 . This is an upper bound for the multicover of B_M^d as well, considering that $\|\mathbf{x} - \check{\mathbf{x}}\|_R^2 = (\mathbf{x} - \check{\mathbf{x}})^\top R(\mathbf{x} - \check{\mathbf{x}}) \leq Tr(R)(\mathbf{x} - \check{\mathbf{x}})^\top (\mathbf{x} - \check{\mathbf{x}}) = Tr(R) \|\mathbf{x} - \check{\mathbf{x}}\|_2^2$ where the inequality is by cauchy-schwarz and the fact that $\|x\|_2 \leq \|x\|_1$. Therefore we have shown the second upper bound.

For the lower bound, let $d' = \min(|(M/2\varepsilon)^2|, d)$, and let $R = \frac{1}{d'}\tilde{I}_{d'}$ where \tilde{I}_k is a diagonal matrix s.t. $\tilde{I}_{ii} = 0$ for i > k and $\tilde{I}_{ii} = 1$ for $i \le k$. We have $M_{\mathcal{R}^d}(B^d_M,\varepsilon) > N_R(B^d_M,\varepsilon) = N_2(B^{d'}_M,\sqrt{d'}\varepsilon) > N_2(B^{d'}_M,M/2) \stackrel{\text{Lemma 2.2}}{>} 2^{d'}$ **Lemma A.4** (2.6). Let $S \subset \mathbb{R}^d$ be a nice set, then: 1. For $\mathcal{X} \subset \mathbb{R}^{d_1}$ and a $d_2 \times d_1$ matrix A we have $M_{\mathcal{S}}(A\mathcal{X}, ||A||\varepsilon) \leq M(\mathcal{X}, \varepsilon)$ 2. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^d$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have $M_{\mathcal{S}}(\sum_{i=1}^n \mathcal{X}_i, \sum_{i=1}^n \varepsilon_i)$ \leq $\prod_{i=1}^{n} M_{\mathcal{S}}(\mathcal{X}_{i}, \varepsilon_{i})$ 3. For $\mathcal{X} \subset \mathbb{R}^d, \varepsilon > 0$, orthonormal matrix U and $\mathbf{b} \in \mathbb{R}^d$ we have $M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b}, \varepsilon) =$ $M_{\mathcal{S}}(\mathcal{X},\varepsilon)$ 4. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^d$ and $\varepsilon > 0$ we have $M_{\mathcal{S}}(\bigcup_{i=1}^n \mathcal{X}_i, \varepsilon) \leq \sum_{i=1}^n M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon)$ 5. For $S = \mathcal{R}_1^d$ and $\mathcal{X}_i \subset [-M_i, M_i]^d$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have $M_{\mathcal{R}_1^d}\left(\prod_{i=1}^n \mathcal{X}_i, \prod_{i=1}^n (M_i + \varepsilon_i) - \prod_{i=1}^n M_i\right) \le \prod_{i=1}^n M_{\mathcal{R}_1^d}(\mathcal{X}_i, \varepsilon_i)$ 6. If for any maximal $R \in S$ (w.r.t. PSD order), $Tr(R) \ge 1$. Fix $\mathcal{X} \subset B_M^{d_1}$, and $\mathcal{L} \subset \mathbb{R}^{d_2,d_1}$ matrices with spectral norm $\le r$. Denote $||A||_{\mathcal{S}} := \min\{t > 0 : \frac{1}{t}A \in S\}$, then $M_{\mathcal{S}}\left(\mathcal{LX}, \varepsilon_2 \sqrt{2r^2 + 2\varepsilon_1^2 \|I_{d_1}\|_{\mathcal{S}}} + \varepsilon_1 M\right) \le M_{\mathcal{R}_1^d}(\mathcal{L}, \varepsilon_1) \cdot M_{\mathcal{S}}(\mathcal{X}, \varepsilon_2)$ *Proof.* We next prove each item separately. *Proof.* (of item 1.) Let $\check{\mathcal{X}}$ be an ε -multicover of \mathcal{X} w.r.t. \mathcal{S} . It is enough to show that $A\check{\mathcal{X}}$ is an $(||A||\varepsilon)$ -multicover of $A\mathcal{X}$. Fix $\mathcal{X} \in \mathcal{X}$ and a PSD matrix $R \in \mathcal{S}$. We need to show that there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|A\mathbf{x} - A\check{\mathbf{x}}\|_{R}^{2} \leq \operatorname{Tr}(R) \|A\|^{2} \varepsilon^{2}$ Now, for any $\check{\mathbf{x}}$ we have $\|A\mathbf{x} - A\check{\mathbf{x}}\|_R^2 = \|A\|^2 (\mathbf{x} - \check{\mathbf{x}})^\top \frac{A^\top}{\|A\|} R \frac{A}{\|A\|} (\mathbf{x} - \check{\mathbf{x}}) = \|A\|^2 \|\mathbf{x} - \check{\mathbf{x}}\|_{\frac{A^\top}{\|A\|} R \frac{A}{\|A\|}}^2$ Finally, since $\check{\mathcal{X}}$ is an ε -multicover w.r.t. \mathcal{S} , there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\mathbf{x} - \check{\mathbf{x}}\|_{\frac{A^{\top}}{1+\pi}R_{\frac{A^{\top}}{1+\pi}}}^2 \leq \varepsilon$

 Proof. (of item 2.) Let $\check{\mathcal{X}}_i$ be an ε_i -multicover of \mathcal{X}_i w.r.t. \mathcal{S} . It is not hard to verify that $\sum_{i=1}^n \check{\mathcal{X}}_i$ is an $(\sum_{i=1}^n \varepsilon_i)$ -multicover of $\sum_{i=1}^n \mathcal{X}_i$, which establishes the proof.

 $\|A\mathbf{x} - A\check{\mathbf{x}}\|_{R}^{2} = \|A\|^{2} \|\mathbf{x} - \check{\mathbf{x}}\|_{\frac{A^{\top}}{\|A^{\top}\|}R\frac{A}{\|A^{\top}\|}}^{2} \leq \operatorname{Tr}(R) \|A\|^{2} \varepsilon^{2}$

Proof. (of item 3.) We have

 $\operatorname{Tr}(\frac{1}{\|A\|}A^{\top}R\frac{1}{\|A\|}A)\varepsilon^2 \leq \operatorname{Tr}(R)\varepsilon^2$. Therefore overall

$$M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b}, \varepsilon) \xrightarrow{Item 2} M_{\mathcal{S}}(U\mathcal{X}, \varepsilon) \cdot M_{\mathcal{S}}(\{\mathbf{b}\}, 0)$$

$$\begin{split} M_{\mathcal{S}}(\{\mathbf{b}\},0) &= 1, \|U\| = 1 \\ &= M_{\mathcal{S}}(U\mathcal{X}, \|U\|\varepsilon) \\ &\stackrel{Item \ 1}{\leq} M_{\mathcal{S}}(\mathcal{X},\varepsilon) \\ \text{Similarly } M_{\mathcal{S}}(\mathcal{X},\varepsilon) &= M_{\mathcal{S}}(U^{-1}(U\mathcal{X} + \mathbf{b}) - U^{-1}\mathbf{b},\varepsilon) \leq M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b},\varepsilon) \text{ implying that} \\ M_{\mathcal{S}}(\mathcal{X},\varepsilon) &= M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b},\varepsilon) \end{split}$$

Proof. (of item 4.) Let $\check{\mathcal{X}}_i$ be an ε -multicover of \mathcal{X}_i w.r.t. \mathcal{S} . It is not hard to verify that $\bigcup_{i=1}^n \check{\mathcal{X}}_i$ is an ε -multicover of $\bigcup_{i=1}^n \mathcal{X}_i$ w.r.t. \mathcal{S} , which establishes the proof.

Proof. (of item 5.) We first prove the item for n = 2. We will then show that the general case follows by induction. In the proof of this item we will denote by $A \circ B$ the elementwise product of two $d \times d$ matrices, and by diag(A) the diagonal matrix obtained by zeroing the non-diagonal entries of A.

⁷⁶³ Let $\check{\mathcal{X}}_i$ be an ε_i -multicover of \mathcal{X}_i w.r.t. \mathcal{R}_1^d . Fix $\mathbf{x}_i \in \mathcal{X}_i$ and a PSD matrix $R \ge 0$ with $\operatorname{Tr}(R) \le 1$. It is enough to show that there is $\check{\mathbf{x}}_i \in \check{\mathcal{X}}_i$ with $\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \le M_1\varepsilon_2 + M_2\varepsilon_1 + \varepsilon_1\varepsilon_2$. We have

$$\begin{aligned} \|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R &\leq \|\mathbf{x}_1\mathbf{x}_2 - \mathbf{x}_1\check{\mathbf{x}}_2\|_R + \|\mathbf{x}_1\check{\mathbf{x}}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \\ &= \|\mathbf{x}_2 - \check{\mathbf{x}}_2\|_{R\circ\mathbf{x}_1\mathbf{x}_1^\top} + \|\mathbf{x}_1 - \check{\mathbf{x}}_1\|_{R\circ\check{\mathbf{x}}_2\check{\mathbf{x}}_2^\top} \end{aligned}$$

769 Now, $\operatorname{Tr}(R \circ \check{\mathbf{x}}_{2} \check{\mathbf{x}}_{2}^{\top}) = \|\check{\mathbf{x}}_{2}\|_{\operatorname{diag}(R)}^{2}$. Thus, we can choose $\check{\mathbf{x}}_{1}$ such that $\|\mathbf{x}_{1} - \check{\mathbf{x}}_{1}\|_{R \circ \check{\mathbf{x}}_{2} \check{\mathbf{x}}_{2}^{\top}} \leq \|\check{\mathbf{x}}_{2}\|_{\operatorname{diag}(R)} \varepsilon_{1}$. We get for any 0

$$\begin{aligned} \|\mathbf{x}_{1}\mathbf{x}_{2} - \check{\mathbf{x}}_{1}\check{\mathbf{x}}_{2}\|_{R} &\leq \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2}\|_{\operatorname{diag}(R)}\varepsilon_{1} \\ &\leq \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{\operatorname{diag}(R)}\varepsilon_{1} + \|\mathbf{x}_{2}\|_{\operatorname{diag}(R)}\varepsilon_{1} \\ &\leq \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{\operatorname{diag}(R)}\varepsilon_{1} + M_{2}\varepsilon_{1} \\ &= \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{\operatorname{diag}(\varepsilon_{1}^{2}R)} + M_{2}\varepsilon_{1} \\ &\leq \sqrt{\frac{\|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top}}}{p} + \frac{\|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{\operatorname{diag}(\varepsilon_{1}^{2}R)}^{2}}{1 - p} + M_{2}\varepsilon_{1} \\ &= \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{\frac{1}{p}R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top}} + \frac{1}{1 - p}\operatorname{diag}(\varepsilon_{1}^{2}R)}{1 - p} + M_{2}\varepsilon_{1} \end{aligned}$$

Where (*) follows from the fact that $a + b \le \sqrt{a^2/p + b^2/(1-p)}$ Now, we can choose $\check{\mathbf{x}}_2 \in \check{\mathcal{X}}_2$ with

$$\begin{aligned} \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{\frac{1}{p}R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top} + \frac{1}{1-p}\operatorname{diag}(\varepsilon_{1}^{2}R)} &\leq \varepsilon_{2}\sqrt{\operatorname{Tr}\left(\frac{1}{p}R \circ \mathbf{x}_{1}\mathbf{x}_{1}^{\top} + \frac{1}{1-p}\operatorname{diag}(\varepsilon_{1}^{2}R)\right)} \\ &\leq \varepsilon_{2}\sqrt{M_{1}^{2}/p + \varepsilon_{1}^{2}/(1-p)} \end{aligned}$$

for $p = M_1/(M_1 + \varepsilon_1)$ we get that

$$\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \le \varepsilon_2(M_1 + \varepsilon_1) + M_2\varepsilon_1 = M_1\varepsilon_2 + M_2\varepsilon_1 + \varepsilon_1\varepsilon_2$$

We next consider n > 2 and conclude the proof by induction. Denote $\mathcal{X}'_2 = \prod_{i=2}^n \mathcal{X}_i$, $M'_2 = \prod_{i=2}^n M_i$ and $\varepsilon'_2 = \prod_{i=2}^n (M_i + \varepsilon_i) - \prod_{i=2}^n M_i$. By the induction hypothesis we have

$$M\left(\mathcal{X}_{2}^{\prime},\varepsilon_{2}^{\prime}\right) \leq \prod_{i=2}^{n} M\left(\mathcal{X}_{i},\varepsilon_{i}\right)$$

By the case n = 2 we have

$$M\left(\mathcal{X}_{1}\mathcal{X}_{2}^{\prime}, M_{1}\varepsilon_{2}^{\prime} + M_{2}^{\prime}\varepsilon_{1} + \varepsilon_{1}\varepsilon_{2}^{\prime}\right) \leq M\left(\mathcal{X}_{1}, \varepsilon_{1}\right)M\left(\mathcal{X}_{2}^{\prime}, \varepsilon_{2}^{\prime}\right)$$
$$\leq M\left(\mathcal{X}_{1}, \varepsilon_{1}\right)\prod_{i=2}^{n}M(\mathcal{X}_{i}, \varepsilon_{i}) = \prod_{i=2}^{n}M(\mathcal{X}_{i}, \varepsilon_{i})$$

this concludes the proof as $\mathcal{X}_1\mathcal{X}_2' = \prod_{i=1}^n \mathcal{X}_i$ and

$$M_{1}\varepsilon_{2}' + M_{2}'\varepsilon_{1} + \varepsilon_{1}\varepsilon_{2}' = (M_{1} + \varepsilon_{1})\left(\prod_{i=2}^{n} (M_{i} + \varepsilon_{i}) - \prod_{i=2}^{n} M_{i}\right) + \varepsilon_{1}\prod_{i=2}^{n} M_{i}$$
$$= \prod_{i=1}^{n} (M_{i} + \varepsilon_{i}) - \prod_{i=1}^{n} M_{i}$$

г		

Proof. (of item 6) For a nice set $S \subset \mathbb{R}^d$, and PSD matrix $R \in \mathbb{R}^{d \times d}$, define $||R||_S = \min\{t > t\}$ $0: \frac{1}{4}R \in S$. Note that this is almost a norm - the triangle inequality, positive definiteness, and homogeneity for positive scalars apply - but do not apply for negative scalars. Let $\check{\mathcal{L}}$ be an ε_1 -multicover of \mathcal{L} w.r.t. \mathcal{R}_1^d and let \mathcal{X} be an ε_2 -multicover of \mathcal{X} w.r.t. \mathcal{S} . We will show that $\mathcal{L}\mathcal{X}$ is an $\varepsilon_2 \sqrt{2r^2 + 2\varepsilon_1} \|I_d\|_{\mathcal{S}} + \varepsilon_1 M$ -multicover of \mathcal{LX} w.r.t. S. Fix $R \in \mathcal{S}$. W.l.o.g we may assume that it is maximal w.r.t. PSD order. Let $W \in \mathcal{L}$ and $\mathbf{x} \in \mathcal{X}$. We need to show that there are $\check{W} \in \check{\mathcal{L}}$ and $\check{\mathbf{x}} \in \dot{\mathcal{X}}$ with $\|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_R \le \varepsilon_2 \sqrt{2\mathrm{Tr}(R)r^2 + 2\varepsilon_1\mathrm{Tr}(R)}\|I_{d_1}\|_{\mathcal{S}} + \varepsilon_1 \sqrt{\mathrm{Tr}(R)}M$. We have

$$\begin{split} \|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_{R} &\leq \|W\mathbf{x} - W\check{\mathbf{x}}\|_{R} + \|W\check{\mathbf{x}} - \check{W}\check{\mathbf{x}}\|_{R} \\ &= \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \|(W - \check{W})\check{\mathbf{x}}\|_{R} \\ &= \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \sqrt{\check{\mathbf{x}}^{\top}(W - \check{W})^{\top}R(W - \check{W})\check{\mathbf{x}}} \end{split}$$

Now, $(W_1, W_2) \mapsto \check{\mathbf{x}}^\top W_1^\top R W_2 \check{\mathbf{x}}$ is a symmetric and positive bi-linear form on the space of $d_2 \times d_1$ matrices of trace

$$\sum_{i=1}^{d_2} \sum_{j=1}^{d_1} \check{\mathbf{x}}^\top E_{ij}^\top R E_{ij} \check{\mathbf{x}} = \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} (\check{x}_j \mathbf{e}_i)^\top R(\check{x}_j \mathbf{e}_i) = \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} \check{x}_j^2 R_{ii} = \operatorname{Tr}(R) \|\check{\mathbf{x}}\|^2$$

Thus, there is $\check{W} \in \check{\mathcal{L}}$ such that $\check{\mathbf{x}}^{\top} (W - \check{W})^{\top} R(W - \check{W}) \check{\mathbf{x}} \leq \operatorname{Tr}(R) \|\check{\mathbf{x}}\|^2 \varepsilon_1^2$. For this \check{W} we have

$$\begin{split} \|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_{R} &\leq \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \varepsilon_{1}\sqrt{\operatorname{Tr}(R)}\|\check{\mathbf{x}}\| \\ &\leq \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \varepsilon_{1}\sqrt{\operatorname{Tr}(R)}(\|\check{\mathbf{x}} - \mathbf{x}\| + \|\mathbf{x}\|) \\ &\leq \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \|\check{\mathbf{x}} - \mathbf{x}\|_{\varepsilon_{1}^{2}\operatorname{Tr}(R)I_{d_{1}}} + \varepsilon_{1}\sqrt{\operatorname{Tr}(R)}M \\ &\leq \sqrt{2}\sqrt{\|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW}^{2} + \|\check{\mathbf{x}} - \mathbf{x}\|_{\varepsilon_{1}^{2}\operatorname{Tr}(R)I_{d_{1}}}^{2}} + \varepsilon_{1}\sqrt{\operatorname{Tr}(R)}M \\ &= \sqrt{2}\|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW + \varepsilon_{1}^{2}\operatorname{Tr}(R)I_{d_{1}}} + \varepsilon_{1}\sqrt{\operatorname{Tr}(R)}M \end{split}$$

Thus, it is possible to choose $\check{\mathbf{x}} \in \check{\mathcal{X}}$ s.t. $\|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW + \varepsilon_1^2 I_{d_1}} \le \varepsilon_2 \sqrt{\|W^{\top}RW + \varepsilon_1^2 \mathrm{Tr}(R)I_{d_1}\|_{\mathcal{S}}}$. Finally, $\|W^{\top}RW + \varepsilon_1^2 \operatorname{Tr}(R)I_{d_1}\|_{\mathcal{S}} \leq r^2 + \varepsilon_1^2 \operatorname{Tr}(R)\|I_{d_1}\|_{\mathcal{S}} \leq Tr(R)r^2 + \varepsilon_1^2 \operatorname{Tr}(R)\|I_{d_1}\|_{\mathcal{S}}$

A.3 MULTICOVER FOR SEQUENCES OF VECTORS

Proof. (of Lemma 3.1) Write $\check{\mathcal{X}} = \{\mathbf{x}^1, \dots, \mathbf{x}^T\}$. Suppose that $\check{\mathcal{X}}$ is a ε -multicover w.r.t. $\mathcal{R}^{d,m}$ and let $\mathbf{x} \in \mathcal{X}$. It is enough to show that there is a r.v. X whose range is $\{\mathbf{x}^1, \dots, \mathbf{x}^T\}$ such that for any $R \in \mathcal{R}_1^{d,m}, \mathbb{E} \| X - \mathbf{x} \|_R^2 \leq \varepsilon^2$. Such a r.v. exists if and only if

$$\min_{\lambda \in \Delta^{T-1}} \max_{R \in \mathcal{R}_1^{d,m}} \sum_{i=1}^T \lambda_i \| \mathbf{x}^i - \mathbf{x} \|_R^2 \le \varepsilon^2$$

since the objective $\sum_{i=1}^{T} \lambda_i \|\mathbf{x}^i - \mathbf{x}\|_R^2 = \sum_{i=1}^{T} \sum_{m=1}^{m} \lambda_i (\mathbf{x}_j^i - \mathbf{x}_j)^\top R(\mathbf{x}_j^i - \mathbf{x}_j)$ is bi-linear in λ and R, and since Δ^{T-1} and $\mathcal{R}_1^{d,m}$ are both convex and compact, we can apply the minmax theorem to conclude that a r.v. X as described above exists if and only if

$$\max_{R \in \mathcal{R}_1^{d,m}} \min_{\lambda \in \Delta^{T-1}} \sum_{i=1}^T \lambda_i \| \mathbf{x}^i - \mathbf{x} \|_R^2 \le \varepsilon^2$$

which is equivalent to $\max_{R \in \mathcal{R}^d_+} \min_{i \in [T]} \|\mathbf{x}^i - \mathbf{x}\|_R \leq \varepsilon$. Which is indeed the case as $\hat{\mathcal{X}}$ is an ε -multicover on \mathcal{X} w.r.t. $\mathcal{R}^{d,m}$.

Suppose now that for any $\mathbf{x} \in \mathcal{X}$ there is a r.v. X whose range is \mathcal{X} such that for any $R \in \mathcal{R}_1^{d,m}$, $\mathbb{E} \| X - \mathbf{x} \|_{R}^{2} \leq \varepsilon^{2}$. This implies that for any $\mathbf{x} \in \mathcal{X}$ and any $R \in \mathcal{R}_{1}^{d}$ there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\check{\mathbf{x}} - \mathbf{x}\|_R \leq \varepsilon$. This implies that $\check{\mathcal{X}}$ is an ε -multicover of \mathcal{X} w.r.t. $\mathcal{R}^{d,m}$. 1. For $\mathcal{X} \subset \mathbb{R}^{d_1,m}$ and a $d_2 \times d_1$ matrix A we have $M_{\mathcal{S}}(A\mathcal{X}, ||A||\varepsilon) \leq 1$ Lemma A.5 (3.2). $M_{\mathcal{S}}(\mathcal{X},\varepsilon)$ 2. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^{d,m}$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have $M_{\mathcal{S}}(\sum_{i=1}^n \mathcal{X}_i, \sum_{i=1}^n \varepsilon_i) \leq 0$ $\prod_{i=1}^{n} M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon_i)$ 3. For $\mathcal{X} \subset \mathbb{R}^{d,m}$, $\varepsilon > 0$ and $\mathbf{b} \in \mathbb{R}^{d,m}$ we have $M_{\mathcal{S}}(U\mathcal{X} + \mathbf{b}, \varepsilon) = M_{\mathcal{S}}(\mathcal{X}, \varepsilon)$ 4. For $\mathcal{X}_1, \ldots, \mathcal{X}_n \subset \mathbb{R}^{d,m}$ and $\varepsilon > 0$ we have $M_{\mathcal{S}}(\bigcup_{i=1}^n \mathcal{X}_i, \varepsilon) \leq \sum_{i=1}^n M_{\mathcal{S}}(\mathcal{X}_i, \varepsilon)$ 5. For $S = \mathcal{R}_1^{d,m} \mathcal{X}_i \subset [-M_i, M_i]^{d,m}$ and $\varepsilon_1, \ldots, \varepsilon_n > 0$ we have

$$M_{\mathcal{R}_{1}^{d,m}}\left(\prod_{i=1}^{n}\mathcal{X}_{i},\prod_{i=1}^{n}(M_{i}+\varepsilon_{i})-\prod_{i=1}^{n}M_{i}\right)\leq\prod_{i=1}^{n}M_{\mathcal{R}_{1}^{d,m}}(\mathcal{X}_{i},\varepsilon_{i})$$

6. For $\mathcal{S} = \mathcal{R}_1^{d,m}$, fix $\mathcal{X} \subset B_M^{d_1,m}$ a set $\mathcal{L} \subset \mathbb{R}^{d_2,d_1}$ of matrices with spectral norm at most r. We have

$$M_{\mathcal{S}}\left(\mathcal{LX},\varepsilon_{2}\sqrt{2r^{2}+2\varepsilon_{1}^{2}d_{1}}+\varepsilon_{1}M\right) \leq M_{\mathcal{R}_{1}^{d,m}}(\mathcal{L},\varepsilon_{1})\cdot M_{\mathcal{S}}(\mathcal{X},\varepsilon_{2})$$

Proof. (of Lemma 3.2) Fix $R, S \in \mathbb{R}^{d,m}$, $\mathbf{x} \in \mathbb{R}^{d,m}$, $A \in \mathbb{R}^{d_1 \times d}$ and $B \in \mathbb{R}^{d \times d_2}$. In this proof we will denote $R \circ S = (R^1 \circ S^1, \ldots, R^m \circ S^m)$ where $R^j \circ S^j$ the elementwise product of R^j and S^j . We will also denote diag $(R) = (\text{diag}(R^1), \ldots, \text{diag}(R^m))$, $AR = (AR^1, \ldots, AR^m)$, $RB = (R^1B, \ldots, R^mB)$ and $\mathbf{xx}^\top = (\mathbf{x}^1(\mathbf{x}^1)^\top, \ldots, \mathbf{x}^m(\mathbf{x}^m)^\top)$.

We next prove each item separately.

Proof. (of item 1.) Let $\tilde{\mathcal{X}}$ be an ε -multicover of \mathcal{X} w.r.t. \mathcal{S} . It is enough to show that $A\tilde{\mathcal{X}}$ is an $(||A||\varepsilon)$ -multicover of $A\mathcal{X}$ w.r.t. \mathcal{S} . Fix $\mathbf{x} \in \mathcal{X}$ and $R \in \mathcal{R}^{d,m}$. We need to show that there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that

$$\|A\mathbf{x} - A\check{\mathbf{x}}\|_{R}^{2} \leq \operatorname{Tr}(R) \|A\|^{2} \varepsilon^{2}$$

Now, for any $\check{\mathbf{x}}$ we have

$$\|A\mathbf{x} - A\check{\mathbf{x}}\|_R^2 = \sum_{i=1}^n (\mathbf{x}_i - \check{\mathbf{x}}_i)^\top A^\top R_i A(\mathbf{x}_i - \check{\mathbf{x}}_i) = \|A\|^2 \|\mathbf{x} - \check{\mathbf{x}}\|_{\frac{A^\top}{\|A\|} R \frac{A}{\|A\|}}^2$$

Finally, since $\check{\mathcal{X}}$ is an ε -multicover, there is $\check{\mathbf{x}} \in \check{\mathcal{X}}$ such that $\|\mathbf{x} - \check{\mathbf{x}}\|_{\|A\|}^2 \leq \operatorname{Tr}(R)\varepsilon^2$. Plugging in $\|A\mathbf{x} - A\check{\mathbf{x}}\|_R^2 = \|A\|^2 \|\mathbf{x} - \check{\mathbf{x}}\|_{\|A\|}^2 \operatorname{Re}_{\|A\|}^A$ we establish the proof.

Proof. (of item 2.) Let $\check{\mathcal{X}}_i$ be an ε_i -multicover of \mathcal{X}_i w.r.t. \mathcal{S} . It is not hard to verify that $\sum_{i=1}^n \check{\mathcal{X}}_i$ is an $(\sum_{i=1}^n \varepsilon_i)$ -multicover of $\sum_{i=1}^n \mathcal{X}_i$ w.r.t. \mathcal{S} , which establishes the proof.

Proof. (of item 3.) We have

$$M_{\mathcal{S}}(\mathcal{X} + \mathbf{b}, \varepsilon) \stackrel{Item 2}{\leq} M_{\mathcal{S}}(\mathcal{X}, \varepsilon) \cdot M_{\mathcal{S}}(\{\mathbf{b}\}, 0) \\ = M_{\mathcal{S}}(\mathcal{X}, \varepsilon)$$

912 Similarly $M_{\mathcal{S}}(\mathcal{X},\varepsilon) = M_{\mathcal{S}}((\mathcal{X}+\mathbf{b})-\mathbf{b},\varepsilon) \le M_{\mathcal{S}}(\mathcal{X}+\mathbf{b},\varepsilon)$ implying that $M_{\mathcal{S}}(\mathcal{X},\varepsilon) = M_{\mathcal{S}}(\mathcal{X}+\mathbf{b},\varepsilon)$ 914 \square

Proof. (of item 4.) Let $\check{\mathcal{X}}_i$ be an ε -multicover of \mathcal{X}_i w.r.t. \mathcal{S} . It is not hard to verify that $\bigcup_{i=1}^n \check{\mathcal{X}}_i$ is an ε -multicover of $\bigcup_{i=1}^n \mathcal{X}_i$ w.r.t. \mathcal{S} , which establishes the proof.

³This claim can be generalized to a more general S. We present the case $S = \mathcal{R}_1^{d,m}$ for simplicity.

Proof. (of item 5.) We first prove the item for n = 2. We will then show that the general case follows by induction. Let \check{X}_i be an ε_i -multicover of \mathcal{X}_i . Fix $\mathbf{x}_i \in \mathcal{X}_i$ and a $R \in \mathcal{R}_1^{d,m}$. It is enough to show that there is $\check{\mathbf{x}}_i \in \check{\mathcal{X}}_i$ with $\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \le M_1\varepsilon_2 + M_2\varepsilon_1 + \varepsilon_1\varepsilon_2$. We have

$$\begin{aligned} \|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R &\leq \|\mathbf{x}_1\mathbf{x}_2 - \mathbf{x}_1\check{\mathbf{x}}_2\|_R + \|\mathbf{x}_1\check{\mathbf{x}}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \\ &= \|\mathbf{x}_2 - \check{\mathbf{x}}_2\|_{R \circ \mathbf{x}_1\mathbf{x}_1^\top} + \|\mathbf{x}_1 - \check{\mathbf{x}}_1\|_{R \circ \check{\mathbf{x}}_2\check{\mathbf{x}}_2^\top} \end{aligned}$$

925 Now, $\operatorname{Tr}(R \circ \check{\mathbf{x}}_{2}\check{\mathbf{x}}_{2}^{\top}) = \|\check{\mathbf{x}}_{2}\|_{\operatorname{diag}(R)}^{2}$. Thus, we can choose $\check{\mathbf{x}}_{1}$ such that $\|\mathbf{x}_{1} - \check{\mathbf{x}}_{1}\|_{R \circ \check{\mathbf{x}}_{2}\check{\mathbf{x}}_{2}^{\top}} \leq$ 926 $\|\check{\mathbf{x}}_{2}\|_{\operatorname{diag}(R)}\varepsilon_{1}$. We get for any 0

$$\begin{aligned} \|\mathbf{x}_{1}\mathbf{x}_{2} - \check{\mathbf{x}}_{1}\check{\mathbf{x}}_{2}\|_{R} &\leq & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{Rox_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2}\|_{diag(R)}\varepsilon_{1} \\ &\leq & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{Rox_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{diag(R)}\varepsilon_{1} + \|\mathbf{x}_{2}\|_{diag(R)}\varepsilon_{1} \\ &\leq & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{Rox_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{diag(R)}\varepsilon_{1} + M_{2}\varepsilon_{1} \\ &= & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{Rox_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{diag(\varepsilon_{1}^{2}R)} + M_{2}\varepsilon_{1} \\ &= & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{Rox_{1}\mathbf{x}_{1}^{\top}} + \|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{diag(\varepsilon_{1}^{2}R)}^{2} + M_{2}\varepsilon_{1} \\ &= & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{Rox_{1}\mathbf{x}_{1}^{\top}}^{2} + \frac{\|\check{\mathbf{x}}_{2} - \mathbf{x}_{2}\|_{diag(\varepsilon_{1}^{2}R)}^{2}}{1 - p} + M_{2}\varepsilon_{1} \\ &= & \|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{\frac{1}{p}Ro\mathbf{x}_{1}\mathbf{x}_{1}^{\top}} + \frac{1}{1 - p}\mathrm{diag}(\varepsilon_{1}^{2}R)} + M_{2}\varepsilon_{1} \end{aligned}$$

Now, we can choose $\check{\mathbf{x}}_2 \in \check{\mathcal{X}}_2$ with

$$\|\mathbf{x}_{2} - \check{\mathbf{x}}_{2}\|_{\frac{1}{p}R \circ \mathbf{x}_{1} \mathbf{x}_{1}^{\top} + \frac{1}{1-p} \operatorname{diag}(\varepsilon_{1}^{2}R)} \leq \varepsilon_{2} \sqrt{\operatorname{Tr}\left(\frac{1}{p}R \circ \mathbf{x}_{1} \mathbf{x}_{1}^{\top} + \frac{1}{1-p} \operatorname{diag}(\varepsilon_{1}^{2}R)\right)} \leq \varepsilon_{2} \sqrt{M_{1}^{2}/p + \varepsilon_{1}^{2}/(1-p)}$$

for $p = M_1/(M_1 + \varepsilon_1)$ we get that

$$\|\mathbf{x}_1\mathbf{x}_2 - \check{\mathbf{x}}_1\check{\mathbf{x}}_2\|_R \le \varepsilon_2(M_1 + \varepsilon_1) + M_2\varepsilon_1 = M_1\varepsilon_2 + M_2\varepsilon_1 + \varepsilon_1\varepsilon_2$$

We next consider n > 2 and conclude the proof by induction. Denote $\mathcal{X}'_2 = \prod_{i=2}^n \mathcal{X}_i$, $M'_2 = \prod_{i=2}^n M_i$ and $\varepsilon'_2 = \prod_{i=2}^n (M_i + \varepsilon_i) - \prod_{i=2}^n M_i$. By the induction hypothesis we have

$$M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{2}^{\prime},\varepsilon_{2}^{\prime}\right) \leq \prod_{i=2}^{n} M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{i},\varepsilon_{i}\right)$$

By the case n = 2 we have

$$M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{1}\mathcal{X}_{2}^{\prime}, M_{1}\varepsilon_{2}^{\prime} + M_{2}^{\prime}\varepsilon_{1} + \varepsilon_{1}\varepsilon_{2}^{\prime}\right) \leq M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{1}, \varepsilon_{1}\right) M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{2}^{\prime}, \varepsilon_{2}^{\prime}\right)$$
$$\leq M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{1}, \varepsilon_{1}\right) \prod_{i=2}^{n} M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{i}, \varepsilon_{i}\right) = \prod_{i=2}^{n} M_{\mathcal{R}_{1}^{d,m}}\left(\mathcal{X}_{i}, \varepsilon_{i}\right)$$

this concludes the proof as $\mathcal{X}_1 \mathcal{X}'_2 = \prod_{i=1}^n \mathcal{X}_i$ and

$$M_{1}\varepsilon_{2}' + M_{2}'\varepsilon_{1} + \varepsilon_{1}\varepsilon_{2}' = (M_{1} + \varepsilon_{1})\left(\prod_{i=2}^{n} (M_{i} + \varepsilon_{i}) - \prod_{i=2}^{n} M_{i}\right) + \varepsilon_{1}\prod_{i=2}^{n} M_{i}$$
$$= \prod_{i=1}^{n} (M_{i} + \varepsilon_{i}) - \prod_{i=1}^{n} M_{i}$$

Proof. (of item 6) Let $S = \mathcal{R}_1^{d,m}$, let $\check{\mathcal{L}}$ be an ε_1 -multicover of \mathcal{L} w.r.t. $\mathcal{R}_1^{d,m}$ and let $\check{\mathcal{X}}$ be an ε_2 -multicover of \mathcal{X} w.r.t. S. We will show that $\check{\mathcal{L}}\check{\mathcal{X}}$ is an $(\varepsilon_2\sqrt{2r^2+2\varepsilon_1^2d_1}+\varepsilon_1M)$ -multicover of

973 \mathcal{LX} w.r.t. \mathcal{S} . Fix $R \in \mathcal{R}_1^{d,m}$. Let $W \in \mathcal{L}$ and $\mathbf{x} \in \mathcal{X}$. We need to show that there are $\check{W} \in \check{\mathcal{L}}$ and 974 $\check{\mathbf{x}} \in \check{\mathcal{X}}$ with $||W\mathbf{x} - \check{W}\check{\mathbf{x}}||_R \le \varepsilon_2 \sqrt{2r^2 + 2\varepsilon_1^2 d_1} + \varepsilon_1 M$. We have

$$\begin{aligned} \|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_{R} &\leq \|W\mathbf{x} - W\check{\mathbf{x}}\|_{R} + \|W\check{\mathbf{x}} - \check{W}\check{\mathbf{x}}\|_{R} \\ &= \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \|(W - \check{W})\check{\mathbf{x}}\|_{R} \end{aligned}$$

$$= \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \sqrt{\sum_{i=1}^{m} (\check{\mathbf{x}}^{i})^{\top} (W - \check{W})^{\top} R^{i} (W - \check{W}) \check{\mathbf{x}}^{i}}$$

Now, $(W_1, W_2) \mapsto \sum_{i=1}^m (\check{\mathbf{x}}^i)^\top W_1^\top R^i W_2 \check{\mathbf{x}}^i$ is a symmetric and positive bi-linear form on the space of $d_2 \times d_1$ matrices of trace

$$\sum_{k=1}^{m} \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} (\check{\mathbf{x}}^k)^\top E_{ij}^\top R^k E_{ij} \check{\mathbf{x}}^k = \sum_{k=1}^{m} \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} (\check{x}_j^k \mathbf{e}_i)^\top R^k (\check{x}_j^k \mathbf{e}_i) = \sum_{k=1}^{m} \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} (\check{x}_j^k)^2 R_{ii}^k = \sum_{k=1}^{m} \operatorname{Tr}(R^k) \|\check{\mathbf{x}}^k\|^2 \le \max_k \|\check{\mathbf{x}}^k\|^2$$

Denote $max := \arg \max_k \|\check{\mathbf{x}}^k\|^2$. By the last inequality there is $\check{W} \in \check{\mathcal{L}}$ such that

$$\sum_{i=1}^{m} (\check{\mathbf{x}}^{i})^{\top} (W - \check{W})^{\top} R^{i} (W - \check{W}) \check{\mathbf{x}}^{i} \leq \varepsilon_{1}^{2} \|\check{\mathbf{x}}^{max}\|^{2}$$

We have

$$\begin{split} \|W\mathbf{x} - \check{W}\check{\mathbf{x}}\|_{R} &\leq \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \varepsilon_{1}\|\check{\mathbf{x}}^{max}\| \\ &\leq \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \varepsilon_{1}(\|\check{\mathbf{x}}^{max} - \mathbf{x}^{max}\| + \|\mathbf{x}^{max}\|)) \\ &\leq \|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW} + \|\check{\mathbf{x}}^{max} - \mathbf{x}^{max}\|_{\varepsilon_{1}^{2}I}^{2} + \varepsilon_{1}M \\ &\leq \sqrt{2}\sqrt{\|\mathbf{x} - \check{\mathbf{x}}\|_{W^{\top}RW}^{2}} + \|\check{\mathbf{x}}^{max} - \mathbf{x}^{max}\|_{\varepsilon_{1}^{2}I}^{2}} + \varepsilon_{1}M \\ &\leq \sqrt{2}\sqrt{r^{2}\varepsilon_{2}^{2}} + d_{1}\varepsilon_{1}^{2}\varepsilon_{2}^{2}} + \varepsilon_{1}M \end{split}$$

1004 Proof. (of Lemma 3.3) Let $\check{\mathcal{X}}$ be an ε -multicover of \mathcal{X} of size $M(\mathcal{X}, \varepsilon)$. By lemma 3.1 for any 1005 $\mathbf{x} \in \mathcal{X}$ there is a distribution $\mathcal{D}_{\mathbf{x}}$ on $\check{\mathcal{X}}$ such that if $X \sim \mathcal{D}_{\mathbf{x}}$ then X is an ε -estimator of \mathbf{x} . In 1006 particular, for any coordinate $i \in [d]$ and $j \in [m]$ we have

$$\mathbb{E}_X (X_i^j - x_i^j)^2 = \mathbb{E}_X (X^j - \mathbf{x}^j)^\top E_{ii} (X^j - \mathbf{x}^j) \le \varepsilon^2$$

1010 Denote $k = \left\lceil \log_{r/2}(d) \right\rceil$. By the above equation and lemma 2.1 we conclude that if 1011 $X(1), \ldots, X(k) \sim \mathcal{D}_{\mathbf{x}}$ then for every $i \in [d]$ and $j \in [m]$

$$\Pr\left(\exists i \in [d] \text{ s.t. } |\text{median}(X(1)_i^j, \dots, X(k)_i^j) - x_i^j| > r\varepsilon\right) < d\left(\frac{2}{r}\right)^k \le 1$$

1015 in particular, there exists $\mathbf{x}(1), \ldots, \mathbf{x}(k) \in \check{\mathcal{X}}$ such that for any $i \in [d]$ and $j \in [m]$, 1016 $|\text{median}(x(1)_i^j, \ldots, x(k)_i^j) - x_i^j| \leq r\varepsilon$. This implies that $\text{median}(\check{\mathcal{X}}^k)$ is an ε -cover of \mathcal{X} w.r.t. 1017 the ℓ^{∞} norm. This concludes the proof as $|\text{median}(\check{\mathcal{X}}^k)| \leq |\check{\mathcal{X}}|^k$ \Box

1019 A.4 STRONGLY BOUNDED ACTIVATION FUNCTIONS

1021 Lemma A.6 (β -Swish Activation Ramachandran et al. (2017)). For a constant $\beta \ge 0$, the function $\frac{x}{1+e^{-\beta x}}$ is strongly-bounded

- 1024 It is shown in Daniely & Granot (2019) that
- **1025** Fact A.7. If ρ is *B*-strongly-bounded then ρ is analytic and its Taylor coefficients around any point are bounded by B^n for any $n \ge 1$.

Proof. For the case of $\beta = 0$ the swish becomes a linear function, and the claim is trivial. For $\beta > 1$, consider the complex function $f(z) = \frac{z}{1+e^{-\beta z}}$. It is defined in the strip $\{z = x + iy : |y| < \frac{1}{\beta}\pi\}$. By Cauchy integral formula, for any $r < \frac{\pi}{\beta}$, $a \in \mathbb{R}$ and $n \ge 0$,

$$f^{(n)}(a) = \frac{n!}{2\pi i} \int_{|z-a|=r} \frac{f(z)}{(z-a)^{n+1}}$$

It follows that

$$\left| f^{(n)}(a) \right| \le \frac{n!}{r^n} \max_{|z-a|=r} |f(z)| \le \frac{n!}{r^n} \max_{x+iy:|y|< r} |f(x+iy)|$$

Now, if $|y| < r < \frac{\pi}{2\beta}$, we have

$$|f(x+iy)| = \frac{|x+iy|}{|1+e^{-i\beta y}e^{-\beta x}|} \le \frac{r}{|1+\cos(-\beta y)e^{-\beta x}|} \le \frac{r}{|1+\cos(\beta r)e^{-\beta x}|} \le r$$

1042 This implies that $\frac{x}{1+e^{-\beta x}}$ is strongly bounded.

Lemma A.8 (Hyperbolic Tangent). The function $\frac{e^{2x}-1}{e^{2x}+1}$ is strongly-bounded

Proof. Consider the complex function $f(z) = \frac{e^{2z}-1}{e^{2z}+1}$. It is defined in the strip $\{z = x + iy : |y| < \frac{1}{2}\pi\}$. By Cauchy integral formula, for any $r < \frac{\pi}{2}\pi$, $a \in \mathbb{R}$ and $n \ge 0$,

$$f^{(n)}(a) = \frac{n!}{2\pi i} \int_{|z-a|=r} \frac{f(z)}{(z-a)^{n+1}}$$

¹⁰⁵⁵ It follows that

$$\left| f^{(n)}(a) \right| \le \frac{n!}{r^n} \max_{|z-a|=r} |f(z)| \le \frac{n!}{r^n} \max_{x+iy:|y|< r} |f(x+iy)|$$

Now, if $|y| < r < \frac{\pi}{8}$ we have that

$$|f(x+iy)| = \frac{|e^{2x}e^{2iy} - 1|}{|e^{2x}e^{2iy} + 1|} \le \frac{2\max\{e^{2x}, 1\}}{|e^{2x}e^{2iy} + 1|} \le \frac{2\max\{e^{2x}, 1\}}{|e^{2x}\cos(2y) + 1|} \le \frac{2\max\{e^{2x}, 1\}}{|e^{2x}\cos(2r) + 1|} \le \frac{2\max\{e^{$$

This implies that $\frac{e^{2x}-1}{e^{2x}+1}$ is strongly bounded.

Lemma A.9 (3.8). Let $p(x) = \sum_{i=0}^{k} a_i X^i$ be a polynomial with $|a_i| \leq B^i$ and suppose that $\mathcal{X} \subset \left[-\frac{1}{8B}, \frac{1}{8B}\right]^{d,m}$. Then, for any Let $0 < \varepsilon \leq 1$,

$$M\left(p(\mathcal{X}),\varepsilon\right) \leq \left(M\left(\mathcal{X},\frac{\varepsilon}{8B}\right)\right)^{\frac{k(k+1)}{2}}$$

Proof. As $M(p(\mathcal{X}), \varepsilon) = M(p(\mathcal{X}) - a_0, \varepsilon)$ we can assume w.l.o.g. that $a_0 = 0$. Denote $a = \frac{1}{8B}$ and $\varepsilon_i = i2^{-2i-1}\varepsilon$. Note that since for -1 < x < 1, $\frac{1}{(1-x)^2} = \sum_{i=1}^{\infty} ix^{i-1}$ we have that k = 1

$$\sum_{i=1}^{k} \varepsilon_i \le \frac{\varepsilon}{4} \sum_{i=1}^{\infty} i(1/2)^{i-1} = \frac{\varepsilon}{4} \frac{1}{(1/2)^2} = \varepsilon$$
(4)

1080	Hence, we have that	
1082	$p(\mathcal{X}) \subset \sum^{k} a_{i} \mathcal{X}^{i}$	$\begin{pmatrix} k \end{pmatrix}$
1083	$M\left(p(\mathcal{X}),\varepsilon\right) \overset{P(\mathcal{X})=\sum_{i=1}^{r} w_{i}\leq$	$M\left(\sum a_i \mathcal{X}^i, \varepsilon\right)$
1084		$\left(\frac{1}{i=1}\right)$
1085	$\sum_{i=1}^k arepsilon_i \leq \varepsilon$	$\begin{pmatrix} k & k \end{pmatrix}$
1086	\leq	$M\left(\sum a_i \mathcal{X}^i, \sum \varepsilon_i\right)$
1087		i=1 $i=1$
1088	<i>Lem</i> .3.2	k
1089	\leq	$\prod M(a_i \mathcal{X}^i, \varepsilon_i)$
1090		<i>i</i> =1
1091	Lem.3.2	$\prod_{k=1}^{k} M\left(\mathcal{Y}^{i} \subset \alpha \right)$
1092	2	$\prod_{i=1}^{M} \left(\mathcal{X}_{i}, \mathcal{E}_{i} / \mathcal{U}_{i} \right)$
1093		k = 1
1094	$\stackrel{ a_i \leq B}{<}$	$\prod M(\mathcal{X}^i, \varepsilon_i/B^i)$
1095	—	$\prod_{i=1}^{n} (i + i) (i + j)$
1096	Claim?	$k \left(\left(\begin{array}{c} c \end{array} \right)^{i} \right)$
1097	\leq	$\prod M\left(\mathcal{X}^{i}, \left(a + \frac{\varepsilon_{i}}{(\alpha_{i})^{i-1}D^{i}}\right) - a^{i}\right)$
1098		$\prod_{i=1} \left(\left(1(2a)^{i-1}B^{i}\right) \right)$
11099	Lem.3.2	$\frac{k}{k} \left(\left(\varepsilon_{i} \right) \right)^{i}$
1100	\leq	$\prod \left(M \left(\mathcal{X}, \frac{\varepsilon_i}{i(2a)^{i-1}B^i} \right) \right)$
1102		$i=1$ ((2α) D))
1103		$\prod_{k=1}^{k} \left(\sum_{i=1}^{k} i2^{-2i-1} \varepsilon \right)^{i}$
1104	=	$\prod \left(M\left(\mathcal{X}, \frac{i(2a)^{i-1}B^i}{i(2a)^{i-1}B^i} \right) \right)$
1105		
1106	=	$\prod_{i=1}^{n} \left(M\left(\chi_{i} - \frac{\varepsilon}{1-\varepsilon}\right) \right)^{i}$
1107		$\prod_{i=1}^{I} \left(\frac{a}{(8aB)^{i-1}8B} \right) $
1108		k (()) i
1109	aaB=1	$\prod \left(M\left(\mathcal{X}, \frac{c}{2R}\right) \right)^{*}$
1110		i=1
1111	_	$\left(M\left(\gamma \varepsilon \right)\right)^{\frac{k(k+1)}{2}}$
1112	—	$\left(M\left(\mathcal{X}, \overline{\mathbf{8B}}\right)\right)$
1114	$($ $)^{i}$	
1115	Claim 3. $\left(a + \frac{\varepsilon_i}{i(2a)^{i-1}B^i}\right) - a^i \leq \frac{\varepsilon_i}{B^i}$	
1116		
1117	<i>Proof.</i> Denote $f(x) = x^i$. Since f is conver	x on \mathbb{R}_+ we have
1118		
1119	$f\left(a+\frac{\varepsilon_i}{1-\varepsilon_i}\right) - f(a)$	$ \psi_i \le f' \left(a + \frac{\varepsilon_i}{1 + $
1120	$i(2a)^{i-1}B^i$	$y = y (1 + i(2a)^{i-1}B^i) i(2a)^{i-1}B^i$
1121	Now, $\frac{\varepsilon_i}{(aB)^i} \leq a \Leftrightarrow \varepsilon_i \leq i2^{i-1}(aB)^i \in$	$\Rightarrow i2^{-2i-1}\varepsilon \leq i2^{i-1}(aB)^i \Leftrightarrow \varepsilon \leq 8^i(aB)^i = 1$. Hence
1122	$\frac{\varepsilon_i}{\varepsilon_i} \leq a \text{Since } f' \text{ is monotone on } \mathbb{R}$	we have
1123	$i(2a)^{i-1}B^i - a$. Since f is monotone on a	
1124	$\varepsilon \left(\begin{array}{c} \varepsilon_i \end{array} \right)$	ε_{i}
1125	$f\left(a+\frac{1}{i(2a)^{i-1}B^{i}}\right)$	$\int -f(a) \le f'(2a) \frac{1}{i(2a)^{i-1}B^i}$
1120		
1128	i his translate to	
1129	$\left(\begin{array}{c} & \varepsilon_i \end{array} \right)^i$	$\varepsilon^i < i (\Omega_{\alpha})^{i-1} \varepsilon_i \varepsilon_i$
1130	$\left(a + \frac{1}{i(2a)^{i-1}B^i}\right) =$	$a \ge i (2a) \qquad \frac{1}{i(2a)^{i-1}B^i} = \frac{1}{B^i}$
1131		· /
1132		
1133		_

A.5 BOUNDING THE MULTICOVERING NUMBER OF NEURAL NETWORKS

Proof of 3.13. As $M(\mathcal{H}, m, \epsilon)$ is monotonically decreasing with ϵ , and the inequality is up to constant, it is enough to prove the lemma for $\epsilon \leq \frac{1}{32B}$. Denote

$$\mathcal{L}_i = \{ W : \| W - W_i^0 \| \le r, \ \| W - W_i^0 \|_F \le R \}$$

Fix examples $\mathbf{x}^1, \dots, \mathbf{x}^m \in B^d_{\sqrt{d}}$. Denote $\mathcal{X}_0 = \{(\mathbf{x}^1, \dots, \mathbf{x}^m)\} \subset \mathbb{R}^{d,m}$. For $1 \leq i \leq t$ denote $\mathcal{X}_i = \rho(\mathcal{L}_i \mathcal{X}_{i-1})$. We need to show that

$$M(\mathcal{X}_t,\epsilon) \lesssim (\log(dm) + \log^2(d/\epsilon))^t \log(dR) \frac{dR^2}{\epsilon^2}$$

where the hidden constant does not depend on the choice of $\mathbf{x}^1, \ldots, \mathbf{x}^m$.

Note that $\mathcal{X}_i \subset (B^d_M)^m$ for $M \leq \sqrt{D}$. Let $k = \lfloor \log_8(\sqrt{d}/\epsilon) \rfloor$ and choose $\epsilon_2 > 0$ such that for $\epsilon_1 = \frac{\epsilon_2}{\sqrt{d+M}}$ we have

$$\epsilon = 8B\epsilon_2\sqrt{2r^2 + 2\epsilon_1^2d} + 8B\epsilon_1M + \sqrt{d}8^{-(k+1)}$$

Note that

$$\epsilon \le 8B\epsilon_2\sqrt{2r^2+2} + 8B\epsilon_2 + \epsilon/8 \Rightarrow \epsilon \le \frac{8}{7}8B(\sqrt{2r^2+2}+2)\epsilon_2 =: C\epsilon_2$$

We have

$$\begin{array}{rcl} \mathbf{1156} & M\left(\mathcal{X}_{t},\epsilon\right) & = & M\left(\rho(\mathcal{L}_{t}\mathcal{X}_{t-1}),\epsilon\right) \\ \mathbf{1158} & & \mathbf{Lemma 3.9} \\ \mathbf{1159} & \leq & \left(M\left(\mathcal{L}_{t}\mathcal{X}_{t-1},\frac{1}{32B}\right)\right)^{\lceil \log_{2}(dm)\rceil} \left(M\left(\mathcal{L}_{t}\mathcal{X}_{t-1},\epsilon_{2}\sqrt{2r^{2}+2\epsilon_{1}^{2}d}+\epsilon_{1}M\right)\right)^{\frac{k(k+1)}{2}} \\ \mathbf{1160} & \\ \mathbf{1161} & \leq & \left(M\left(\mathcal{L}_{t}\mathcal{X}_{t-1},\epsilon_{2}\sqrt{2r^{2}+2\epsilon_{1}^{2}d}+\epsilon_{1}M\right)\right)^{\lceil \log_{2}(dm)\rceil+\frac{k(k+1)}{2}} \\ \mathbf{1163} & & \mathbf{Lemma 3.2} \\ \mathbf{1164} & \leq & \left(M(\mathcal{L}_{t},\epsilon_{1})M\left(\mathcal{X}_{t-1},\epsilon_{2}\right)\right)^{\lceil \log_{2}(dm)\rceil+\frac{k(k+1)}{2}} \\ \mathbf{1165} & & \epsilon^{/C\leq\epsilon_{2}} \\ \leq & \left(M(\mathcal{L}_{t},\epsilon_{1})M\left(\mathcal{X}_{t-1},\epsilon/C\right)\right)^{\lceil \log_{2}(dm)\rceil+\frac{k(k+1)}{2}} \end{array}$$

By lemma 2.5 and since $\epsilon_1 = \frac{\epsilon_2}{\sqrt{d}+M}$ we have

$$\log(M(\mathcal{L}_t, \epsilon_1)) \le \left\lceil 2(d+M^2) \frac{2R^2 + 1/4}{\epsilon_2^2} \right\rceil \log(4d^4 \lceil R \rceil + 6d^2)$$

Thus we get

$$\log \left(M\left(\mathcal{X}_{t},\epsilon\right) \right) = \log \left(M\left(\rho(\mathcal{L}_{t}\mathcal{X}_{t-1}),\epsilon\right) \right)$$
$$\lesssim \left(\log(dm) + \log^{2}(d/\epsilon)\right) \left(\frac{dR^{2}}{\epsilon^{2}}\log(dR) + M\left(\mathcal{X}_{t-1},\epsilon/C\right)\right)$$

Inductively, we get that

$$\log\left(M\left(\mathcal{X}_{t},\epsilon\right)\right) \lesssim (\log(dm) + \log^{2}(d/\epsilon))^{i} \log(dR) \frac{dR^{2}}{\epsilon^{2}}$$

A.6 BOUNDING REPRESENTATIVENESS WITH MULTICOVER

Lemma A.10 (3.11). Let $\ell : \mathbb{R}^d \times \mathcal{Y} \to \mathbb{R}$ be L-Lipschitz w.r.t. $\|\cdot\|_{\infty}$ and B-bounded. Assume that for any $\frac{\sqrt{nB}}{\sqrt{m}BL} \leq \varepsilon \leq 1$, $\ln M(\mathcal{H}, m, \varepsilon) \leq \frac{n}{\varepsilon^2}$. Then for any distribution \mathcal{D} on \mathcal{Z}

$$\mathbb{E}_{S \sim \mathcal{D}^m} \operatorname{rep}_{\mathcal{D}}(S, \mathcal{H}) \lesssim \frac{(L+B)\sqrt{n}}{\sqrt{m}} \sqrt{\log(dm)} \log(m)$$

 Choosing $M = \left[\log_2 \left(\sqrt{\frac{m}{n}} \right) \right]$ we get

$$A \leq B \sqrt{\frac{n}{m}} + \frac{12B\sqrt{n \left\lceil \log_2(dm) \right\rceil}}{\sqrt{m}} \left(\frac{4L \log(m)}{B} + 1\right)$$

В APPROXIMATE DESCRIPTION LENGTH AND MULTICOVER

In this section we show that multicover is closely related to the notion of approximate description length (ADL) as defined by Daniely & Granot (2019). We start with a definition that is slightly different from the definition used in Daniely & Granot (2019). We say that \mathcal{X} has ε -ADL of n if there is a protocol between two entities, Alice and Bob with the following properties. Upon seeing $\mathbf{x} \in \mathcal{X}$, Alice, that is allowed to use randomness, sends a message $s \in \{0, 1\}^n$ to Bob. Upon seeing s, Bob generates a vector $\hat{\mathbf{x}}$ that is an ε -estimator to \mathbf{x} . Formally, there is a probability space (Ω, P) (representing Alice's randomness) and functions $A: \mathcal{X} \times \Omega \to \{0,1\}^n$ and $B: \{0,1\}^n \to \mathbb{R}^{d,m}$ such that for any $\mathbf{x} \in \mathcal{X}$ the random variable $\omega \mapsto B(A(\mathbf{x}, \omega))$ is an ε -estimator of \mathbf{x} . We denote by $ADL(\mathcal{X}, \varepsilon)$ the minimal k for which \mathcal{X} has an ε -ADL of k.

The following lemma shows that ADL is closely related to multicover.

1228 Lemma B.1.
$$ADL(\mathcal{X}, \varepsilon) = \lfloor \log_2 (M(\mathcal{X}, \varepsilon)) \rfloor$$

Proof. Observe that \mathcal{X} has ε -ADL of n if and only if there is a set \mathcal{X} such that for any $\mathbf{x} \in \mathcal{X}$ there is a random vector $X \in \mathcal{X}$ that is a ε -estimator of x. By lemma 3.1 this is valid if and only if \mathcal{X} is an ε -multicover. It follows that ADL(\mathcal{X}, ε) is the minimal k for which \mathcal{X} has an ε -multicover of size 2^k . In other words, $ADL(\mathcal{X}, \varepsilon) = |\log_2(M(\mathcal{X}, \varepsilon))|$.

We next turn to the definition used in Daniely & Granot (2019). We define *unbiased* ε -ADL, by making two modification to the definition of ε -ADL. First, we require that $\mathbb{E}_{\omega \sim P} B(A(\mathbf{x}, \omega)) = \mathbf{x}$. Second, we allow sending messages of unbounded length (i.e. a message in $\{0,1\}^*$), and just require that the expected number of sent bits will be at most n. We denote by $uADL(\mathcal{X},\varepsilon)$ the minimal k for which \mathcal{X} has an unbiased ε -ADL of k. We note that Daniely & Granot (2019) defined the ADL of \mathcal{X} to be $uADL(\mathcal{X}, 1)$. The following lemma connects unbiased ADL and ADL by showing that ignoring poly-logarithmic factors, $uADL(\mathcal{X}, 1) \leq k$ if and only if $ADL(\mathcal{X}, \varepsilon) \leq \frac{k}{\varepsilon^2}$. By lemma B.1 this happens if and only if $\log_2(M(\mathcal{X}, \varepsilon)) \leq \frac{k}{\varepsilon^2}$.

1242 **Lemma B.2.** *Fix* $\mathcal{X} \subset \mathbb{R}^{d,m}$ *. We have* 1243 • $\forall 0 < \epsilon \leq 1, \text{ ADL}(\mathcal{X}, \varepsilon) \leq O\left(\frac{u \text{ADL}(\mathcal{X}, 1)}{\varepsilon^2}\right)$ 1244 1245 1246 • If $ADL(\mathcal{X}, \varepsilon) \leq \frac{k}{\varepsilon^2}$ for any $0 < \epsilon \leq 1$ then $uADL(\mathcal{X}, 1) = O\left(\log^2(dm)k\right)$ 1247 Where the constant in the big-O notation are universal. 1248 1249 *Proof.* (sketch) Denote $k = uADL(\mathcal{X}, 1)$. Given $\mathbf{x} \in \mathcal{X}$ and using $O\left(\frac{k}{\epsilon^2}\right)$ expected bits Alice can 1250 send to Bob $\begin{bmatrix} \frac{1}{\epsilon^2} \end{bmatrix}$ independent and unbiased 1-estimators of x. If Bob averages these estimators, 1251 he gets an ϵ -estimator of **x**. This implies that $uADL(\mathcal{X}, \varepsilon) \leq O\left(\frac{uADL(\mathcal{X}, 1)}{\varepsilon^2}\right)$. It is therefore 1253 enough to show that $ADL(\mathcal{X}, \varepsilon) \leq O\left(uADL(\mathcal{X}, \varepsilon/\sqrt{2})\right)$. By Lemma B.1 it is enough to show that 1254 $\log(M(\mathcal{X},\varepsilon)) \le O\left(u \text{ADL}(\mathcal{X},\varepsilon/\sqrt{2})\right).$ 1255 1256 Denote $k = u \text{ADL}(\mathcal{X}, \varepsilon/\sqrt{2})$ and fix a probability space (Ω, P) and functions $A : \mathcal{X} \times \Omega \to \{0, 1\}^*$ 1257 and $B: \{0,1\}^* \to \mathbb{R}^{d,m}$ such that for any $\mathbf{x} \in \mathcal{X}$ the random variable $\omega \mapsto B(A(\mathbf{x},\omega))$ is an unbiased $(\varepsilon/\sqrt{2})$ -estimator of \mathbf{x} , and $\mathbb{E}_{\omega} \operatorname{len}(A(\mathbf{x},\omega)) \leq k$. Fix $R \in \mathcal{R}_{1}^{d,m}$. We have $\mathbb{E}_{\omega} \|\mathbf{x} - B(A(\mathbf{x},\omega))\|_{R}^{2} \leq \epsilon^{2}/2$ By Markov inequality, there exists ω such that $\|\mathbf{x} - B(A(\mathbf{x},\omega))\|_{R}^{2} \leq \epsilon^{2}$ 1258 1259 and $len(A(\mathbf{x}, \omega)) \leq 2k$. This implies that $\mathcal{X} := \{B(s) : s \in \{0, 1\}^*, \ len(s) \leq 2k\}$ is a ϵ -cover 1261 of \mathcal{X} w.r.t. R. This is true for any $R \in \mathcal{R}_1^{d,m}$, and therefore $\check{\mathcal{X}}$ is a ϵ -multicover of \mathcal{X} . This implies 1262 that $\log(M(\mathcal{X},\varepsilon)) \leq 4k$ 1263

For the second item, let X_n and \bar{X}_n be $\frac{\epsilon}{\sqrt{2^n}}$ -estimators of **x**, which can be encoded using $\frac{k2^n}{\epsilon^2}$ bits each. Let Z_n be a r.v. that is 2^n w.p. 2^{-n} and 0 otherwise. Assume that all these random variables are independent. Consider now the estimator

$$X^{N} = X_{1} + \sum_{n=1}^{N} Z_{n} (X_{n+1} - \bar{X}_{n})$$

1270 We first claim that Bob can generate such an estimator using $O\left(\frac{kN}{\epsilon^2}\right)$ expected bits set from Alice. 1271 Indeed, Alice can first sample the Z_n 's. Then, for any n, if $Z_n \neq 0$, send the index n using $O(\log(n))$ bits as well as X_n and \bar{X}_n using $\frac{k2^n}{\epsilon^2}$. The expected number of sent bits is $O\left(2^{-n}\frac{k2^n}{\epsilon^2}\right) = O\left(\frac{k}{\epsilon^2}\right)$. The total expected number of sent bits is therefore $O\left(\frac{kN}{\epsilon^2}\right)$ bits. We next show that X^N is an $(\epsilon\sqrt{1+4N})$ -estimator of x. Indeed, for any unit vector u we have

$$\operatorname{Var}\left(\left\langle \mathbf{u}, X^{N} \right\rangle\right) = \operatorname{Var}\left(\left\langle \mathbf{u}, X_{1} \right\rangle\right) + \sum_{n=1}^{N} \operatorname{Var}\left(Z_{n} \left\langle \mathbf{u}, X_{n+1} - X_{n} \right\rangle\right)$$

1279
1280
$$\leq \epsilon^2 + \sum_{n=1}^N \mathbb{E} \left(Z_n \langle \mathbf{u}, X_{n+1} - X_n \rangle \right)^2$$
1281

1268 1269

1276 1277 1278

$$\epsilon^{2} + \sum_{n=1}^{N} \mathbb{E}Z_{n}^{2} \mathbb{E}\left(\langle \mathbf{u}, X_{n+1} - X_{n} \rangle\right)^{2}$$

$$\sum_{n=1}^{N} \mathbb{E}Z_{n}^{2} \mathbb{E}\left(\langle \mathbf{u}, X_{n+1} - X_{n} \rangle\right)^{2}$$

1285
1286
1287

$$= \epsilon^{2} + \sum_{n=1}^{N} 2^{n} \mathbb{E} \langle \mathbf{u}, X_{n+1} - X_{n} \rangle^{2}$$
N

1288
1289
1290
$$\leq \epsilon^{2} + 2\sum_{n=1}^{\infty} 2^{n} \left(\mathbb{E} \langle \mathbf{u}, X_{n+1} - \mathbf{x} \rangle^{2} + \mathbb{E} \langle \mathbf{u}, X_{n} - \mathbf{x} \rangle^{2} \right)$$
N

1291
1292
1293

$$\leq \epsilon^2 + \epsilon^2 2 \sum_{n=1}^{N} 2^n \left(2^{-(n+1)} + 2^{-n} \right)$$

 $\leq \epsilon^2 (1+4N)$

Let Y^N be an unbiased (1/2)-estimator of $\mathbf{y} := \mathbf{x} - \mathbb{E}X^N = \mathbf{x} - \mathbb{E}X_{N+1}$. Since X_{N+1} is $\left(\frac{1}{2(N+1)/2}\right)$ -estimator of \mathbf{x} , we have that the absolute value of each coordinate of \mathbf{y} is at most

1296 1297 1298 1299 1300 1301	$\frac{\epsilon}{2^{(N+1)/2}}$. Thus, Alice can send $\frac{1}{2}$ estimator of y as follows: for any $i \in [m]$ and $j \in [d]$, w.p. $ y_j^i $ send $\operatorname{sign}(y_j^i)$ and the indices i and j . If the pair (i, j) were sent, Bob will define $Y_j^i = \operatorname{sign}(y_j^i)$. Otherwise, he will define $Y_j^i = 0$. It is not hard to verify (see Daniely & Granot (2019) for details) that Y is an unbiased $\frac{\epsilon}{2^{(N+1)/2}}$ -estimator of Y . Likewise, the expected number of sent bits per coordinate is $O\left(\frac{\epsilon \log(md)}{2^{(N+1)/2}}\right)$, resulting with a total cost of $O\left(\frac{\epsilon md \log(md)}{2^{(N+1)/2}}\right) = O\left(\frac{md \log(md)}{2^{(N+1)/2}}\right)$ bits.
1302 1303 1304	Finally, $X = X^N + Y^N$ is an unbiased $\sqrt{1/2 + \epsilon^2(1+4N)}$ -estimator of x which costs $O\left(\frac{md \log(md)}{2(N+1)/2} + \frac{kN}{\epsilon^2}\right)$ bits to encode. Choosing $\epsilon = \sqrt{1/(2+4N)}$ and $N = 2\log_2(md)$ we
1305	gen an unbiased 1-estimator of x which costs $O(k \log^2(md))$ bits to encode
1306	gen an anomased i estimator of \mathbf{x} which costs o $(n \log (n \omega))$ ons to encode.
1307	
1308	
1309	
1310	
1311	
1312	
1314	
1315	
1316	
1317	
1318	
1319	
1320	
1321	
1322	
1323	
1324	
1325	
1320	
1328	
1329	
1330	
1331	
1332	
1333	
1334	
1335	
1336	
1337	
1338	
1339	
1340	
1342	
1343	
1344	
1345	
1346	
1347	
1348	
1349	