

474 **Supplementary Material**

475 **A Random Matrices**

476 **A.1 Multiplication of Random Matrices**

477 In this section we present and prove some statistical properties of general random matrices and their
 478 multiplications. Let $\{Q_n \in \mathbb{R}^{m_n \times m_{n-1}}\}_{n=1}^N$ denote a set of random matrix whose elements are
 479 sampled iid from a distribution with mean 0 and variance σ_n^2 , with bounded kurtosis. Let

$$\begin{aligned} \mathbf{Q}^l &= \prod_{n=l-1}^1 Q_n = Q_{l-1} \cdots Q_1, & \mathbf{B}^l &= \mathbf{Q}^{l\top} \mathbf{Q}^l \in \mathbb{R}^{m_0 \times m_0} \quad l \in [2 \dots N] \\ \mathbf{Q}^l &= \prod_{n=N}^{l+1} Q_n = Q_N \cdots Q_{l+1}, & \mathbf{A}^l &= \mathbf{Q}^l \mathbf{Q}^{l\top} \in \mathbb{R}^{m_N \times m_N} \quad l \in [1 \dots N-1] \end{aligned} \quad (6)$$

480 **Theorem 3.** $\forall l$

$$\mathbb{E}(\mathbf{B}^l) = \beta_l I \quad \beta_l = \prod_{n=1}^{l-1} m_n \sigma_n^2 \quad (7)$$

$$\mathbb{E}(\mathbf{A}^l) = \alpha_l I \quad \alpha_l = \prod_{n=l+1}^N m_{n-1} \sigma_n^2 \quad (8)$$

481 *Proof.* We only prove (7), as the proof of (8) is similar. To simplify the presentation, we use the
 482 following auxiliary notations: $V = Q_1, U = \prod_{n=l-1}^2 Q_n \implies \mathbf{Q}^l = UV$.

483 Proof proceeds by induction on l .

484 • $l = 2$:

$$\begin{aligned} \mathbb{E}[B_{ij}^l] &= \mathbb{E}\left[\sum_{k=1}^{m_1} V_{ki} V_{kj}\right] \stackrel{i \neq j}{=} \sum_{k=1}^{m_1} \mathbb{E}[V_{ki}] \mathbb{E}[V_{kj}] \\ \mathbb{E}[B_{ii}^l] &= \mathbb{E}\left[\sum_{k=1}^{m_1} V_{ki} V_{ki}\right] = \sum_{k=1}^{m_1} \mathbb{E}[V_{ki}^2] \end{aligned}$$

485 Thus

$$\mathbb{E}[B_{ij}^l] = \begin{cases} 0 & i \neq j \quad (\text{off diagonal}) \\ m_1 \sigma_1^2 & i = j \quad (\text{diagonal}) \end{cases}$$

486 • Assume that (7) holds for $l-1$.

$$B_{ij}^l = \sum_k \mathbf{Q}_{ki}^l \mathbf{Q}_{kj}^l = \sum_k \sum_{\nu} U_{k\nu} V_{\nu i} \sum_{\rho} U_{k\rho} V_{\rho j}$$

487 and therefore

$$\mathbb{E}[B_{ij}^l] = \sum_k \sum_{\nu} \sum_{\rho} \mathbb{E}[U_{k\nu} V_{\nu i} U_{k\rho} V_{\rho j}] = \sum_{\nu} \sum_{\rho} \mathbb{E}[V_{\nu i} V_{\rho j}] \sum_k \mathbb{E}[U_{k\nu} U_{k\rho}]$$

488 where the last transition follows from the independence of U and V . Once again, we
 489 consider the diagonal and off-diagonal elements separately. If $i \neq j$:

$$\mathbb{E}[B_{ij}^l] = \sum_{\nu} \sum_{\rho} \mathbb{E}[V_{\nu i}] \mathbb{E}[V_{\rho j}] \sum_k \mathbb{E}[U_{k\nu} U_{k\rho}] = 0$$

490 If $i = j$:

$$\mathbb{E}[B_{ii}^l] = \sum_{\nu} \sum_{\rho} \mathbb{E}[V_{\nu i} V_{\rho j}] \mathbb{E}[(U^\top U)_{\nu\rho}] = \sum_{\nu} \mathbb{E}[V_{\nu i}^2] \mathbb{E}[(U^\top U)_{\nu\nu}]$$

491

Using the induction assumption

$$\mathbb{E}[B_{ij}^l] = \begin{cases} 0 & i \neq j \quad (\text{off diagonal}) \\ m_1 \sigma_1^2 \prod_{n=2}^{l-1} m_n \sigma_n^2 & i = j \quad (\text{diagonal}) \end{cases}$$

492 from which (7) follows. \square

493 Let m denote the width of the smallest hidden layer, $m = \min(m_1, \dots, m_{N-1})$, and assume that
 494 $\max(m_1, \dots, m_{N-1}) - \min(m_1, \dots, m_{N-1})$ is bounded by some M_b as $m \rightarrow \infty$. Assume the
 495 following initialization scheme

496 **Definition 4.** The elements of $\{Q_n\}_{n=1}^N$ are chosen iid from a distribution with mean 0 and variance
 497 σ_n^2 , where

$$\sigma_n^2 = \frac{2}{m_{n-1} + m_n} \quad 1 < n < N, \quad \sigma_1^2 = \frac{1}{m_1}, \quad \sigma_N^2 = \frac{1}{m_{N-1}}$$

498 For large m , it follows that

$$\begin{aligned} m_n \sigma_n^2 &= 1 + O\left(\frac{1}{m}\right) & n \in [1 \dots N-1] \\ m_{n-1} \sigma_n^2 &= 1 + O\left(\frac{1}{m}\right) & n \in [2 \dots N] \end{aligned}$$

499 **Corollary 3.1.** With initialization as in Def. 4, $\forall l$

$$\mathbb{E}(B^l) = [1 + O\left(\frac{1}{m}\right)]I, \quad \mathbb{E}(A^l) = [1 + O\left(\frac{1}{m}\right)]I$$

500 **Theorem 4.** With initialization as in Def. 4, $\forall l$

$$\text{var}(B^l) = O\left(\frac{1}{m}\right), \quad \text{var}(A^l) = O\left(\frac{1}{m}\right)$$

501 *Proof.* We prove by induction on l that:

$$\mathbb{E}[(B_{ij}^l)^2] = \begin{cases} O\left(\frac{1}{m}\right) & i \neq j \quad (\text{off diagonal}) \\ 1 + O\left(\frac{1}{m}\right) & i = j \quad (\text{diagonal}) \end{cases}, \quad \mathbb{E}[B_{ii}^l B_{jj}^l] = 1 + O\left(\frac{1}{m}\right) \quad (9)$$

502 For $l = 2$, (9) follows from Lemma 2 and Corr 3.1. We now assume that (9) holds for $l - 1$ and prove
 503 for l , using notations as above: $V = Q_1$, $U = \prod_{l=1}^2 Q_n$, $Q^l = UV$.

$$\mathbb{E}[(B_{ij}^l)^2] = \sum_{\nu, \rho} \sum_{\alpha, \beta} \mathbb{E}[V_{\nu i} V_{\rho j} V_{\alpha i} V_{\beta j} \sum_{k, n} U_{k\nu} U_{k\rho} U_{n\alpha} U_{n\beta}]$$

504 Let $B' = U^\top U$. Using the induction assumption

$$\mathbb{E}[(B_{ij}^l)^2] \stackrel{i \neq j}{=} \sum_{\nu, \rho} \mathbb{E}[V_{\nu i}^2 V_{\nu j}^2] \mathbb{E}[(B'_{\nu\rho})^2] = O\left(\frac{1}{m}\right)$$

505 When $i = j$, there are 3 cases where the terms in the sum above do not equal 0: (i) $\nu = \alpha$, $\rho =$
 506 β , $\nu \neq \rho$ or $\nu = \beta$, $\rho = \alpha$, $\nu \neq \rho$; (ii) $\nu = \rho$, $\alpha = \beta$, $\nu \neq \alpha$; (iii) $\nu = \rho = \alpha = \beta$. Case (i) is
 507 similar to the above, and we therefore only expand cases (ii) and (iii) next:

$$(ii) \sum_{\nu, \alpha} \mathbb{E}[V_{\nu i}^2 V_{\alpha i}^2] \mathbb{E}[B'_{\nu\nu} B'_{\alpha\alpha}] = 1 + O\left(\frac{1}{m}\right)$$

$$(iii) \sum_{\nu} \mathbb{E}[V_{\nu i}^4] \mathbb{E}[(B'_{\nu\nu})^2] = O\left(\frac{1}{m}\right)$$

508 In the derivation of (iii) we exploit the assumption that the kurtosis of the distribution used to sample
 509 Q_n is fixed at G and cannot depend on m , indicating that $\mathbb{E}[V_{\nu i}^4] = G\sigma_1^4$.

510 A similar argument would show that $\mathbb{E}[B_{ii}^l B_{jj}^l] = 1 + O\left(\frac{1}{m}\right)$. \square

511 **Theorem 5.** Let $\{X(\mathfrak{m})\}$ denote a sequence of random matrices where $\mathbb{E}[X(\mathfrak{m})] = [1 + O(\frac{1}{\mathfrak{m}})]I$
512 and $\text{var}[X(\mathfrak{m})] = O(\frac{1}{\mathfrak{m}})$. Then $X(\mathfrak{m}) \xrightarrow{p} I$, where \xrightarrow{p} denotes convergence in probability.

513 *Proof.* We need to show that $\forall \epsilon, \delta > 0 \exists \mathfrak{m}' \in \mathbb{N}$, such that $\forall \mathfrak{m} > \mathfrak{m}'$

$$P(|X(\mathfrak{m}) - I| > \epsilon) < \delta$$

514 Henceforth we use X as shorthand for $X(\mathfrak{m})$. Since $\mathbb{E}(X) = [1 + O(\frac{1}{\mathfrak{m}})]I$, it follows that $\forall \epsilon > 0$
515 $\exists \mathfrak{m}_1 \in \mathbb{N}$ such that $\forall \mathfrak{m} > \mathfrak{m}_1$, the following holds element-wise:

$$|\mathbb{E}(X) - I| < \frac{\epsilon}{2}$$

516 Thus

$$P(|X - I| > \epsilon) \leq P(|X - \mathbb{E}(X)| > \frac{\epsilon}{2})$$

517 Since $\text{var}(X) = O(\frac{1}{\mathfrak{m}})$, it follows that $\forall \epsilon, \delta > 0, \exists \mathfrak{m}_2 \in \mathbb{N} \ni \forall \mathfrak{m} > \mathfrak{m}_2$

$$\text{var}(X) < \frac{\epsilon^2}{4} \delta$$

518 From the above, and using Chebyshev inequality

$$P(|X - I| > \epsilon) < \frac{4\text{var}(X)}{\epsilon^2} < \delta$$

519 $\forall \mathfrak{m} > \mathfrak{m}'$, where $\mathfrak{m}' = \max\{\mathfrak{m}_1, \mathfrak{m}_2\}$.

520

□

521 Let $A^l(m)$ and $B^l(\mathfrak{m})$ denote a sequence of random matrices as defined in (6), corresponding to
522 models for which $m = \min(m_1, \dots, m_{L-1})$.

Corollary 5.1.

$$B^l(m) \xrightarrow{p} I \quad \forall l \in [2 \dots N] \quad A^l(m) \xrightarrow{p} I \quad \forall l \in [1 \dots N-1]$$

523 The proof follows from Corr 3.1, Thm 4 and Thm 5.

524 A.2 Dynamics of Random Matrices

525 Consider a dynamical process, where the random matrices defined above are changed as $Q_j \rightarrow$
526 $Q_j - \Delta Q_j \forall j$, and specifically from (21):

$$\Delta Q_j = \mu \left(\prod_{n=N}^{j+1} Q_n \right)^\top E_r \left(\prod_{n=j-1}^1 Q_n \right)^\top, \quad E_r = \left(\prod_{n=N}^1 Q_n \right) \Sigma_{XX} - \Sigma_{YX} \quad (10)$$

527 Denoting $Q^l \rightarrow Q^l - \Delta Q^l$ and applying the product rule

$$\Delta Q^l = \sum_{j=1}^{l-1} \left(\prod_{n=l-1}^{j+1} Q_n \right) \Delta Q_j \left(\prod_{n=j-1}^1 Q_n \right) \quad (11)$$

528 For $B^l = Q^l{}^\top Q^l$ and denoting $B^l \rightarrow B^l - \Delta B^l$:

$$\Delta B^l = [\Delta Q^l{}^\top Q^l + Q^l{}^\top \Delta Q^l] \quad (12)$$

529 Before proceeding to analyze ΔB^l , we note that

$$m_N \sigma_N^2 = \frac{K}{m_{N-1}} = \frac{K}{\mathfrak{m}} [1 + O(\frac{1}{\mathfrak{m}})]$$

530 and therefore from Thm 3

$$\mathbb{E}[(Q^N{}^\top)Q^N] = [\frac{K}{\mathfrak{m}} + O(\frac{1}{\mathfrak{m}})]I \quad (13)$$

531 **Theorem 6.** For sequence $B^l(m)$ defined as above, if

$$B^l(m) \xrightarrow{p} I, \quad \text{var}[B^l(m)] = O\left(\frac{1}{m}\right)$$

532 then

$$\Delta B^l(m) \xrightarrow{p} 0, \quad \text{var}[\Delta B^l(m)] = O\left(\frac{1}{m}\right)$$

533 *Proof.* $B^l(m) \xrightarrow{p} I$ implies that $\forall \epsilon, \delta > 0 \exists \hat{m} \in \mathbb{N}$, such that $\forall m > \hat{m}$ and with probability larger
534 than $1 - \delta$.

$$B^l(m) = I + e_1 \quad |e_1| < \epsilon, \quad \forall l \in [2 \dots N] \quad (14)$$

535 In addition, from (13) and Thms. 4-5

$$Q^{N\top} Q^N = \frac{K}{m} I + e_2 \quad |e_2| < \epsilon$$

536 We fix m and let B^l be a shorthand for $B^l(m)$. Now

$$B^{N+1} = Q_{N-1}^\top Q^{N\top} Q^N Q_{N-1} = \frac{K}{m} Q^{N-1\top} Q^{N-1} + O(\epsilon) = \frac{K}{m} I + O(\epsilon) \quad (15)$$

537 To evaluate ΔB^l from (12), we start from

$$Q^{l\top} \Delta Q^l = \sum_{j=1}^{l-1} \left(\prod_{l-1}^1 Q_n \right)^\top \left(\prod_{l-1}^{j+1} Q_n \right) \Delta Q_j \prod_{j-1}^1 Q_n$$

538 Simplifying t_j – the j^{th} term in the sum

$$t_j = \mu \prod_1^{l-1} Q_n^\top \prod_{l-1}^{j+1} Q_n \prod_{j+1}^N Q_n^\top E_r \prod_1^{j-1} Q_n^\top \prod_{j-1}^1 Q_n = \mu B^l (B^{j+1})^{-1} Q^{N\top} E_r B^j + O(\epsilon)$$

539 The last transition is exactly true when $B^l = I$ and $B^{j+1} = I$, as shown in Lemma 3 in §A.3.
540 Substituting E_r

$$\begin{aligned} t_j &= \mu B^l (B^{j+1})^{-1} Q^{N\top} [Q^N \Sigma_{XX} - \Sigma_{YX}] B^j + O(\epsilon) \\ &= \mu B^l (B^{j+1})^{-1} [B^{N+1} \Sigma_{XX} + Q^{N\top} \Sigma_{YX}] B^j + O(\epsilon) \end{aligned}$$

541 Substituting (14) and (15)

$$t_j = \mu \left[\frac{K}{m} \Sigma_{XX} + Q^{N\top} \Sigma_{YX} \right] + O(\epsilon) \quad (16)$$

542 From (16) and Lemma 2

$$\mathbb{E}[Q^{l\top} \Delta Q^l] = \sum_{j=1}^{l-1} \mathbb{E}[t_j] = \mu l \frac{K}{m} I + O(\epsilon)$$

543 Since $\Delta Q^{l\top} Q^l = [Q^{l\top} \Delta Q^l]^\top$, it follows from (12) that

$$\mathbb{E}[\Delta B^l] = 2\mu l \frac{K}{m} I + O(\epsilon) \quad (17)$$

544 To conclude the proof, we need to show that $\forall \epsilon', \delta' > 0 \exists \hat{m}' \in \mathbb{N}$, such that $\forall m > \hat{m}'$

$$P(|\Delta B^l| > \epsilon') < \delta'$$

545 Since (17) is true with probability $(1 - \delta) \forall \epsilon, \delta$ and $\forall m > \hat{m}$, we choose ϵ and \hat{m}' such that

$$|\mathbb{E}[\Delta B^l]| < \frac{\epsilon'}{2} \quad \forall m > \hat{m}' \quad (18)$$

546

$$P(|\Delta B^l| > \varepsilon') \leq (1 - \delta)P\left(|\Delta B^l - \mathbb{E}(\Delta B^l)| > \frac{\varepsilon'}{2}\right) < \frac{4\text{var}(\Delta B^l)}{\varepsilon'^2}(1 - \delta)$$

547 $\text{var}(\Delta B^l) = O\left(\frac{1}{m}\right)$ implies that $\exists \hat{m}'' \in \mathbb{N}, \delta > 0$, such that $\forall m > \hat{m}''$

$$\frac{4\text{var}(\Delta B^l)}{\varepsilon'^2}(1 - \delta) < \delta'$$

548 It now follows that $\Delta B^l(m) \xrightarrow{p} 0$.

549 To analyze the variance, we assume that all the moments of the distribution functions used to sample
550 Q_n are bounded. Thus, from (16), the variance of $t_j \forall j$ remains $O\left(\frac{1}{m}\right)$. Likewise, since ΔB^l is a sum
551 of matrices, each with variance $O\left(\frac{1}{m}\right)$ thus bounding the covariance by $O\left(\frac{1}{m}\right)$, we can deduce that
552 $\text{var}(\Delta B^l) = O\left(\frac{1}{m}\right)$.

553 □554 **Theorem 7.** For sequence $A^l(m)$ defined as above, if

$$A^l(m) \xrightarrow{p} I \quad \text{and} \quad \text{var}[A^l(m)] = O\left(\frac{1}{m}\right)$$

555 then

$$\Delta A^l(m) \xrightarrow{p} 0 \quad \text{and} \quad \text{var}[\Delta A^l(m)] = O\left(\frac{1}{m}\right)$$

556 The proof is mostly similar to Thm 6, though we additionally need to show the following in order to
557 replace (13):

$$\mathbb{E}[\mathbf{Q}^0 \Sigma_{xx} \mathbf{Q}^{0\top}] = \left[\frac{q}{m} + O\left(\frac{1}{m}\right)\right]I$$

558 This, in turn, can be proved in a similar manner to the proof of Thm 3, when taking into account the
559 initialization scheme defined in Def. 4.

560 **Note about convergence rate.** In Thm 6, convergence to 0 when $m \rightarrow \infty$ is governed by $O\left(\frac{K}{m}\right)$.
561 In Thm 7, convergence is governed by $O\left(\frac{q}{m}\right)$.

562 **A.3 Some Useful Lemmas**563 **Lemma 1.** Given function $G(W) = \frac{1}{2}\|UWVX - Y\|_F^2$, its derivative is the following

$$\frac{dG(W)}{dW} = U^\top U W V X (V X)^\top - U^\top Y (V X)^\top = U^\top [U W V \Sigma_{xx} - \Sigma_{yx}] V^\top \quad (19)$$

564 **Lemma 2.** Given $\mathbf{Q} = \prod_{n=N}^1 Q_n$, where $Q_n \in \mathbb{R}^{m_n \times m_{n-1}}$ denotes a random matrix whose
565 elements are sampled iid from a distribution with mean 0 and variance $\sigma_n^2, \forall i, j$.

$$\mathbb{E}[\mathbf{Q}_{ij}] = 0 \quad \text{var}[\mathbf{Q}_{ij}] = \frac{1}{m_N} \prod_{n=1}^N m_n \sigma_n^2 \quad (20)$$

566 *Proof.* By induction on N . Clearly for $N = 1$:

$$\mathbb{E}[\mathbf{Q}_{ij}] = \mathbb{E}[(Q_1)_{ij}] = 0 \quad \text{var}[\mathbf{Q}_{ij}] = \text{var}[(Q_1)_{ij}] = \sigma_1^2$$

567 Assume that (20) holds for $N - 1$. Let $V = \prod_{n=N-1}^1 Q_n, U = Q_N$. It follows that

$$\mathbb{E}[\mathbf{Q}_{ij}] = \mathbb{E}[(UV)_{ij}] = \sum_k \mathbb{E}[U_{ik} V_{kj}] = \sum_k \mathbb{E}[U_{ik}] \mathbb{E}[V_{kj}] = 0$$

568 where the last transition follows from the independence of U and V . In a similar manner

$$\begin{aligned} \text{var}[\mathbf{Q}_{ij}] &= \mathbb{E}[\mathbf{Q}_{ij}^2] = \mathbb{E}\left[\left(\sum_k U_{ik} V_{kj}\right)^2\right] = \mathbb{E}\left[\sum_k U_{ik} V_{kj} \sum_l U_{il} V_{lj}\right] = \sum_k \mathbb{E}[U_{ik}^2] \mathbb{E}[V_{kj}]^2 \\ &= m_{N-1} \sigma_N^2 \frac{1}{m_{N-1}} \prod_{n=1}^{N-1} m_n \cdot \sigma_n^2 = \frac{1}{m_N} \prod_{n=1}^N m_n \cdot \sigma_n^2 \end{aligned}$$

569 With the initialization scheme defined in Def. 4, $\text{var}(\mathbf{Q}_{ij}) = O(\frac{1}{mn})$. \square

570 **Lemma 3.** Consider matrix multiplication CD where $C \in \mathbb{R}^{k \times m}$, $D \in \mathbb{R}^{m \times k}$, $k \ll m$ and
 571 $\text{rank}(CD) = k$. Define $\Delta_1 \in \mathbb{R}^{m \times k}$, $\Delta_2 \in \mathbb{R}^{k \times m}$. Then

$$C\Delta_1 = \Delta_2\Delta_1 = I \implies CD = C\Delta_1\Delta_2D$$

572 *Proof.* Since $C = \Delta_1^+$ and $\Delta_2 = \Delta_1^+$

$$C = \Delta_1^+\Delta_1C = C\Delta_1\Delta_1^+ = C\Delta_1\Delta_2$$

573 \square

574 B Supplementary Proofs and Additional Models

575 B.1 Deep Linear networks

576 Here we prove Thm 1 as defined in Section 2.1.

577 **Theorem 1.** The compact matrix representation \mathbf{W} obeys the following dynamics

$$\mathbf{W}^{s+1} = \mathbf{W}^s - \mu \sum_{l=1}^L A_l^s \cdot Er^s \cdot B_l^s + O(\mu^2)$$

578 where the gradient scale matrices A_l^s , B_l^s are defined in (3)

$$A_l^s := \left(\prod_{j=L}^{l+1} W_j^s \right) \left(\prod_{j=L}^{l+1} W_j^s \right)^\top \in \mathbb{R}^{K \times K} \quad B_l^s := \left(\prod_{j=l-1}^1 W_j^s \right)^\top \left(\prod_{j=l-1}^1 W_j^s \right) \in \mathbb{R}^{q \times q}$$

579

580 *Proof.* At time s , the gradient step ΔW_l^s of layer l is defined by differentiating $L(\mathbb{X})$ with respect to
 581 W_l^s . Henceforth we omit index s for clarity. First, we rewrite $L(\mathbb{X})$ as follows:

$$L(\mathbb{X}; W_l) = \frac{1}{2} \left\| \left(\prod_{j=L}^{l+1} W_j \right) W_l \left(\prod_{j=l-1}^1 W_j \right) X - Y \right\|_F^2$$

582 Differentiating $L(\mathbb{X}; W_l)$ to obtain the gradient $\Delta W_l = \frac{\partial L(\mathbb{X}; W_l)}{\partial W_l}$, using Lemma 1 above, we get

$$\Delta W_l = \left(\prod_{j=L}^{l+1} W_j \right)^\top [\mathbf{W}\Sigma_{XX} - \Sigma_{YX}] \left(\prod_{j=l-1}^1 W_j \right)^\top \quad (21)$$

583 Finally

$$\Delta \mathbf{W} = \prod_{l=L}^1 (W_l - \mu \Delta W_l) - \prod_{l=L}^1 W_l = -\mu \sum_{l=1}^L \left(\prod_{n=L}^{l+1} W_n \right) \Delta W_l \left(\prod_{n=l-1}^1 W_n \right) + O(\mu^2)$$

584 Substituting ΔW_l and Er (as defined in Def. 3) into the above completes the proof.

585 \square

586 B.2 Adding Non-Linear ReLU Activation

587 The results shown in Fig. 2b pertain to a relatively simple non-linear model analyzed by Arora et al.
 588 (2019), here adapted to classification rather than regression. Specifically, it is a two-layer model with
 589 ReLU activation, where only the weights of the first layer are being learned. Similarly to (1), the loss
 590 is defined as

$$L(\mathbb{X}) = \frac{1}{2} \sum_{i=1}^n \|f(\mathbf{x}_i) - \mathbf{y}_i\|^2 \quad f(\mathbf{x}_i) = \mathbf{a}^\top \cdot \sigma(W\mathbf{x}_i), \quad \mathbf{a} \in \mathbb{R}^m, \quad W \in \mathbb{R}^{m \times d}$$

591 m denotes the number of neurons in the hidden layer. We consider a binary classification problem
 592 with 2 classes, where $y_i = 1$ for $\mathbf{x}_i \in C_1$, and $y_i = -1$ for $\mathbf{x}_i \in C_2$. $\sigma(\cdot)$ denotes the ReLU
 593 activation function applied element-wise to vectors, where $\sigma(u) = u$ if $u \geq 0$, and 0 otherwise.

594 At time s , each gradient step is defined by differentiating $L(\mathbb{X})$ with respect to W . Due to the
 595 non-linear nature of the activation function $\sigma(\cdot)$, we separately³ differentiate each row of W , denoted
 596 \mathbf{w}_r where $r \in [m]$, as follows:

$$\begin{aligned} \mathbf{w}_r^{s+1} - \mathbf{w}_r^s &= -\mu \frac{\partial L(\mathbb{X})}{\partial \mathbf{w}_r} \Big|_{\mathbf{w}_r = \mathbf{w}_r^s} = -\mu \sum_{i=1}^n \left[\mathbf{a}^\top \cdot \sigma(W^s \mathbf{x}_i) - y_i \right] \frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{w}_r} \Big|_{\mathbf{w}_r = \mathbf{w}_r^s} \\ &= -\mu \sum_{i=1}^n \left[\sum_{j=1}^m a_j \sigma(\mathbf{w}_j^s \cdot \mathbf{x}_i) - y_i \right] a_r \mathbf{x}_i^\top \mathbb{1}_{\mathbf{w}_r^s}(\mathbf{x}_i) \\ &= -\mu a_r \sum_{i=1}^n \mathbb{1}_{\mathbf{w}_r^s}(\mathbf{x}_i) \left[\Psi^s(\mathbf{x}_i) \cdot \mathbf{x}_i - y_i \right] \mathbf{x}_i^\top \quad \text{where } \Psi^s(\mathbf{x}_i) = \sum_{j=1}^m a_j \mathbf{w}_j^s \mathbb{1}_{\mathbf{w}_j^s}(\mathbf{x}_i) \end{aligned}$$

597 Above $\mathbb{1}_{\mathbf{w}_r^s}(\mathbf{x}_i)$ denotes the indicator function that equals 1 when $\mathbf{w}_r^s \cdot \mathbf{x}_i \geq 0$, and 0 otherwise.

598 In order to proceed, we make two assumptions:

- 599 1. The distribution of the data is symmetric where $P(\mathbf{x}_i) = P(-\mathbf{x}_i)$.
- 600 2. W and \mathbf{a} are initialized so that $\mathbf{w}_{2i}^0 = -\mathbf{w}_{2i-1}^0$ and $a_{2i} = -a_{2i-1} \forall i \in [\frac{m}{2}]$.

601 It follows from Assumption 2 that at the beginning of training $\mathbb{1}_{\mathbf{w}_{2j}^0}(\mathbf{x}_i) + \mathbb{1}_{\mathbf{w}_{2j-1}^0}(\mathbf{x}_i) = 1, \forall \mathbf{x}_i$
 602 such that $\mathbf{w}_{2j-1} \mathbf{x}_i \neq \mathbf{w}_{2j} \mathbf{x}_i \neq 0$, and $\forall j \in [\frac{m}{2}]$. Consequently

$$\Psi^0(\mathbf{x}_i) = \sum_{j=1}^m a_j \mathbf{w}_j^0 \mathbb{1}_{\mathbf{w}_j^0}(\mathbf{x}_i) = \frac{1}{2} \sum_{j=1}^m a_j \mathbf{w}_j^0 = \frac{1}{2} \mathbf{a}^\top W^0$$

603 $\forall \mathbf{x}_i$ such that $\mathbf{w}_{2j-1} \mathbf{x}_i \neq \mathbf{w}_{2j} \mathbf{x}_i \neq 0$. Finally

$$\mathbf{w}_r^1 - \mathbf{w}_r^0 = -\mu a_r \left[\frac{1}{2} \mathbf{a}^\top W^0 \sum_{\substack{i=1 \\ \mathbf{w}_r^0 \mathbf{x}_i \geq 0}}^n \mathbf{x}_i \mathbf{x}_i^\top - \sum_{\substack{i=1 \\ \mathbf{w}_r^0 \mathbf{x}_i \geq 0}}^n y_i \mathbf{x}_i^\top \right]$$

604 Next, we note that Assumption 1 implies

$$\mathbb{E} \left[\sum_{\substack{i=1 \\ \mathbf{w} \cdot \mathbf{x}_i \geq 0}}^n \mathbf{x}_i \mathbf{x}_i^\top \right] = \frac{1}{2} \mathbb{E} \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right] = \frac{1}{2} \mathbb{E}[\Sigma_{XX}]$$

605 for any vector \mathbf{w} . Thus, if the sample-size n is large enough, at the beginning of training we expect
 606 to see

$$\mathbf{w}_r^{s+1} - \mathbf{w}_r^s \approx -\mu \frac{a_r}{2} \left[\mathbf{a}^\top W^s \Sigma_{XX} - \tilde{\mathbf{m}}_r^s \right] \quad \forall r$$

607 where row vector $\tilde{\mathbf{m}}_r^s$ denotes the vector difference between the centroids of classes C_1 and C_2 ,
 608 computed in the half-space defined by $\mathbf{w}_r^s \cdot \mathbf{x} \geq 0$. Finally (for small s)

$$W^{s+1} - W^s \approx -\mu \frac{1}{2} \left[(\mathbf{a} \mathbf{a}^\top) W^s \Sigma_{XX} - \tilde{M}^s \right]$$

609 where \tilde{M}^s denotes the matrix whose r -th row is $a_r \tilde{\mathbf{m}}_r^s$. This equation is reminiscent of the single
 610 layer linear model dynamics $\mathbf{W}^{s+1} = \mathbf{W}^s - \mu E r^s$, and we may conclude that when it holds and
 611 using the principal coordinate system, the rate of convergence of the j -th column of W^s is governed
 612 by the singular value d_j .

³Since the ReLU function is not everywhere differentiable, the following may be considered the definition of the update rule.

613 **B.3 Weight Evolution**

614 To analyze the weight dynamics, we first shift to the principal coordinate system defined in Def 1.
 615 In this representation $Er^s = W^s D - M$, where $D = \text{diag}(\{d_j\}_{j=1}^q)$ is a diagonal matrix. Based
 616 on Thm 1 and the subsequent discussion of convergence rate, assuming that the width of the hidden
 617 layers is very large, we can readily substitute $B_l^s \approx I \forall l$ in (2), to obtain

$$\mathbf{W}^{s+1} = \mathbf{W}^s - \mu \sum_{l=1}^L A_l^s E r^s + O(\mu^2) \quad (22)$$

618 Let $\mathbf{w}_j \in \mathbb{R}^K$ denote the j -th column of \mathbf{W} , \mathbf{m}_j denote the j -th column of M . From (22) we have

$$\mathbf{w}_j^{s+1} = \mathbf{w}_j^s - \mu \sum_{l=1}^L A_l^s (d_j \mathbf{w}_j^s - \mathbf{m}_j) \quad j \in [K]$$

619 This is a telescoping series; denoting $A^s = \sum_{l=1}^L A_l^s$,

$$\begin{aligned} \mathbf{w}_j^{s+1} &= \mathbf{w}_j^s - \mu A^s (d_j \mathbf{w}_j^s - \mathbf{m}_j) = (I - \mu d_j A^s) \mathbf{w}_j^s + \mu A^s \mathbf{m}_j = \dots \\ &= \prod_{\nu=1}^s (I - \mu d_j A^\nu) \mathbf{w}_j^0 + \mu \left[\sum_{\nu=1}^s \prod_{\rho=\nu+1}^s (I - \mu d_j A^\rho) A^\nu \right] \mathbf{m}_j \end{aligned} \quad (23)$$

620 The only difference between individual columns lies in d_j , which governs the rate of convergence of
 621 the first term to 0, and the rate of convergence of the second term to the optimal value of $\frac{1}{d_j} \mathbf{m}_j$.

622 In the discussion following the proof of Thm 2, we noted that the approximation $A_l^s \approx I$ breaks down
 623 before $B_l^s \approx I$. Nevertheless, while it is still valid, (23) further simplifies to the following

$$\begin{aligned} \mathbf{w}_j^{s+1} &= (1 - \mu d_j L)^s \mathbf{w}_j^0 + \mu \left[\sum_{\nu=1}^s (1 - \mu d_j L)^{(s-\nu)} L I \right] \mathbf{m}_j = \lambda_j^s \mathbf{w}_j^0 + \mu L \left[\sum_{k=0}^{s-1} \lambda_j^k \right] \mathbf{m}_j \\ &= \lambda_j^s \mathbf{w}_j^0 + (1 - \lambda_j^s) \frac{\mathbf{m}_j}{d_j} \quad \lambda_j = 1 - \mu d_j L \end{aligned}$$

624 **C Additional Empirical Results**

625 **C.1 Weight Initialization**

626 We evaluate empirically the weight initialization scheme from Def. 4. When compared to Glorot
 627 uniform initialization (Glorot & Bengio, 2010), the only difference between the two schemes lies
 628 in how the first and last layers are scaled. Thus, in order to highlight the difference between the
 629 methods, we analyze a fully connected linear network with a single hidden layer, whose dimension
 630 (the number of hidden neurons) is much larger than the input and output dimensions. We trained
 631 $N=10$ such networks on a binary classification problem, once with the initialization suggested in
 632 Def. 4, and again with Glorot uniform initialization. While both initialization schemes achieve the
 633 same final accuracy upon convergence, our proposed initialization variant converges faster on both
 634 train and test datasets (see Fig. 10).

635 **C.2 Spectral Bias**

636 The *spectral bias*, discussed in Section 4.3, can also induce similar learning order in different
 637 networks. To support the discussion in Section 4.3, in §C.2.2) we analyze the relation between the
 638 *spectral bias* and *accessibility*, in order to clarify its relation to the *Learning Order Constancy* and
 639 the *PC-bias*. First, however, we expand the scope of the empirical evidence for this effect to the
 640 classification scenario and real image data (§C.2.1).

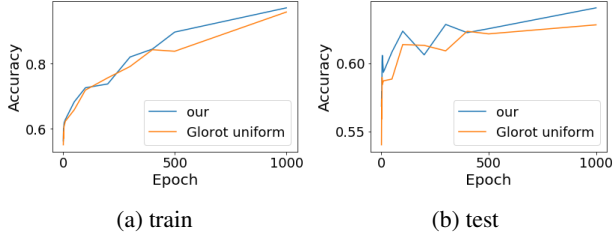


Figure 10: Learning curves of a fully connected linear network with one hidden layer, trained on the dogs and cats dataset, and initialized by Glorot uniform initialization (orange) and the initialization proposed in Def. 4 (blue).

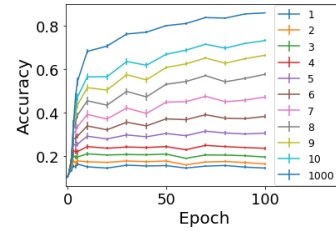


Figure 11: Evaluations on test-sets projected to the first P principal components, for different values of P (see legend) of 10 VGG-19 models trained on CIFAR-10

641 C.2.1 Spectral Bias in Classification

642 Rahaman et al. (2019) showed that when regressing a 2D function by a neural network, the model
 643 seems to approximate the lower frequencies of the function before its higher frequencies. Here
 644 we extend this empirical observation to the classification framework. Thus, given frequencies
 645 $\kappa = (\kappa_1, \kappa_2, \dots, \kappa_m)$ with corresponding phases $\phi = (\varphi_1, \varphi_2, \dots, \varphi_m)$, we consider the mapping
 646 $\lambda : [-1, 1] \rightarrow \mathbb{R}$ given by

$$\lambda(z) = \sum_{i=1}^m \sin(2\pi\kappa_i z + \varphi_i) := \sum_{i=1}^m \text{freq}_i(z) \quad (24)$$

647 Above κ is strictly monotonically increasing, while ϕ is sampled uniformly.

648 The classification rule is defined by $\lambda(z) \leq 0$. We created a binary dataset whose points are fully
 649 separated by $\lambda(z)$, henceforth called the *frequency dataset* (see visualization in Fig. 13 and details in
 650 §D.4). When training on this dataset, we observe that the frequency of the corresponding separator
 651 increases as learning proceeds, in agreement with the results of Rahaman et al. (2019).

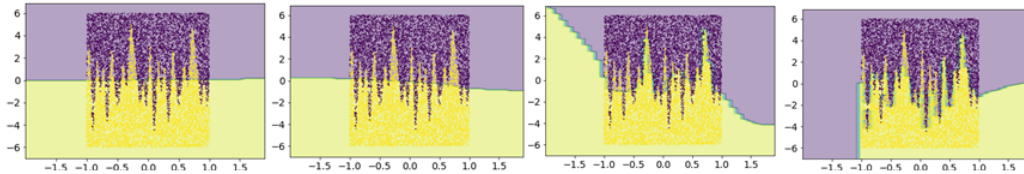


Figure 12: Visualization of the separator learned by st-VGG when trained on the frequency dataset, as captured in advancing epochs (from left to right): 1, 100, 1000, 10000. Each point represents a training example (yellow for one class and purple for the other). The background color represents the classification that the network predicts for points in that region.

652 To visualize the decision boundary of an st-VGG network trained on this dataset as it evolves with
 653 time, we trained $N=100$ st-VGG networks. Since the data lies in \mathbb{R}^2 , we can visualize it and the
 654 corresponding network’s inter-class boundary at each epoch as shown in Fig. 12. We can see that the
 655 decision boundary incorporates low frequencies at the beginning of the learning, adding the higher
 656 frequencies only later on. The same qualitative results are achieved with other instances of st-VGG
 657 as well. We note that while the decision functions are very similar in the region where the training
 658 data is, at points outside of the data they differ drastically across networks.

659 C.2.2 Spectral Bias: Relation to Accessibility

660 In order to connect between the learning order, which is defined over examples, and the Fourier
 661 analysis of a separator, we define for each example its *critical frequency*, which characterizes the
 662 smallest number of frequencies needed to correctly classify the example. To illustrate, consider the
 663 *frequency dataset* defined above. Here, the *critical frequency* is defined as the smallest $j \in [m]$ such
 664 that $\lambda_j(z) = \sum_{i=1}^j \text{freq}_i(z)$ classifies the example correctly (see Figs. 14a,14b).

665 In this binary classification task, we observe a strong connection between the order of learning and
 666 the *critical frequency*. Specifically, we trained $N=100$ st-VGG networks on the *frequency dataset*,

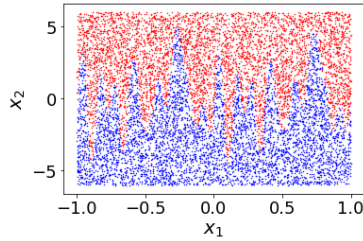


Figure 13: Visualization of the classification dataset used to extend Rahaman et al. (2019) to a classification framework.

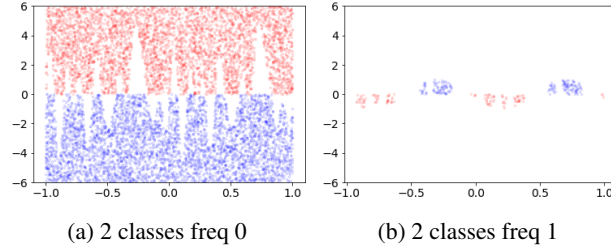
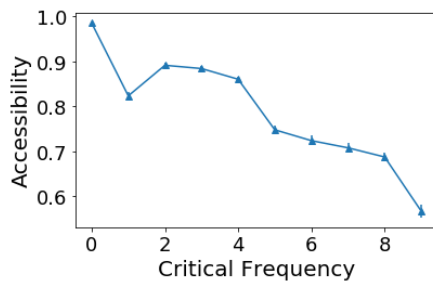
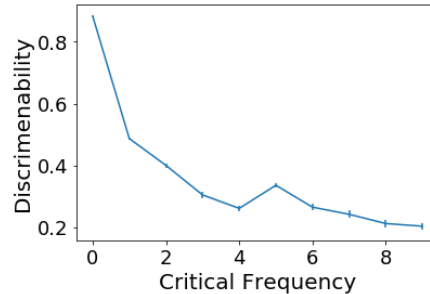


Figure 14: Visualization of the *critical frequency*, showing all the points in the 2D-frequency dataset with *critical frequency* of (a) 0, and (b) 1.



(a)



(b)

Figure 15: (a) Correlation between *critical frequency* and *accessibility* score in the 2D-frequency dataset. (b) Correlation between *discriminability* and *critical frequency* in the 2D-frequency dataset.

667 and correlated the *accessibility* scores with the *critical frequency* of the examples (see Fig. 15a). We
 668 see a strong negative correlation ($r = -0.93$, $p < 10^{-2}$), suggesting that examples whose *critical*
 669 *frequency* is high are learned last by the networks.

670 In order to see the effect of the *spectral bias* in real classification task and extend the above analysis
 671 to natural images, we need to define a score that captures the notion of *critical frequency*. To this
 672 end, we define the *discriminability* measure of an example - the percentage out of its k neighbors that
 673 share the same class as the example. Intuitively, an example has a low *discriminability* score when
 674 it is surrounded by examples from other classes, which forces the learned boundary to incorporate
 675 high frequencies. In Fig. 15b we plot the correlation between the *discriminability* and the *critical*
 676 *frequency* for the 2D frequency dataset. The high correlation ($r = -0.8$, $p < 10^{-2}$) indicates that
 677 *discriminability* indeed captures the notion of *critical frequency*.

678 C.3 Projection to higher PC's

679 In Section 3.3 we described an evaluation methodology, based on the creation of a modified *test-set*
 680 by projecting each test example on the span of the first P principal components. We repeat this
 681 experiment with VGG-19 networks on CIFAR-10, and plot the results in Fig. 11.

682 D Methodology

683 D.1 Implementation details and hyper parameters

684 The results reported in Section 5 represent the mean performance of 100 st-VGG and linear st-VGG
 685 networks, trained on the small mammals dataset. The results reported in Section 5 represent the
 686 mean performance of 10 2-layers fully connected linear networks trained over the cats and dogs
 687 dataset. The results in Fig. 9 represent the mean performance of 100 st-VGG network trained on the
 688 small mammals dataset. In every experimental setup the network's hyper-parameters were coarsely
 689 grid-searched to achieve good performance over the validation set, for a fair comparison. Other
 690 hyper-parameters exhibit similar results.

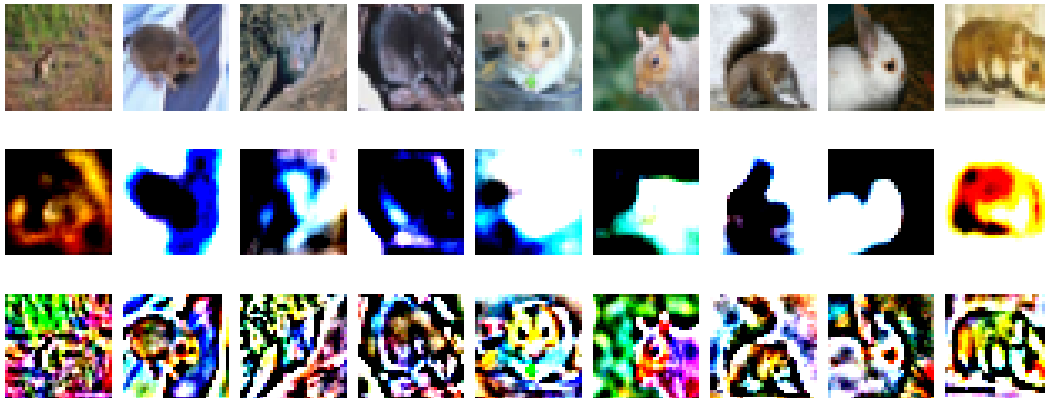


Figure 16: Visualization of the small mammals dataset, with amplification of 1.5% of its principal components by a factor of 10. Top: original data; middle: data amplified along the first principal components; bottom: data amplified along the last principal components

691 D.2 Generalization Gap

692 In Section 5 we discuss the evaluation of networks on datasets with amplified principal components.
 693 Examples of these images are shown in Fig. 16: the top row shows examples of the original
 694 images, the middle row shows what happens to each image when its 1.5% most significant principal
 695 components are amplified, and the bottom row shows what happens when its 1.5% least significant
 696 principal components are amplified. Amplification was done by a factor of 10, which is significantly
 697 smaller than the ratio between the values of the first and last principal components of the data. After
 698 amplification, all the images were re-normalized to have 0 mean and std 1 in every channel as
 699 customary.

700 D.3 Architectures

701 **st-VGG.** A stripped version of VGG which we used in many of the experiments. It is a convolutional
 702 neural network, containing 8 convolutional layers with 32, 32, 64, 64, 128, 128, 256, 256 filters
 703 respectively. The first 6 layers have filters of size 3×3 , and the last 2 layers have filters of size
 704 2×2 . Every other layer is followed by a 2×2 max-pooling layer and a 0.25 dropout layer. After
 705 the convolutional layers, the units are flattened, and there is a fully-connected layer with 512 units
 706 followed by 0.5 dropout. The batch size we used was 100. The output layer is a fully-connected layer
 707 with output units matching the number of classes in the dataset, followed by a softmax layer. We
 708 trained the network using the SGD optimizer, with cross-entropy loss. When training st-VGG, we
 709 used a learning rate of 0.05.

710 **Linear st-VGG.** A linear version of the st-VGG network. In linear st-VGG, we change the activation
 711 function to the identity function, and replace max-pooling by average pooling with a similar stride.

712 **Linear fully connected network.** An L -layered fully connected network. Each layer contains 1024
 713 weights, initialized with Glorot uniform initialization. 0.5 dropout is used before the output layer.
 714 Networks are trained with an SGD optimizer, without momentum or L_2 regularization.

715 D.4 Datasets

716 In all the experiments and all the datasets, the data was always normalized to have 0 mean and std 1,
 717 in each channel separately.

718 **Small Mammals.** The small-mammals dataset used in our experiments is the relevant super-class
 719 of the CIFAR-100 dataset. It contains 2500 train images divided into 5 classes equally, and 500 test
 720 images. Each image is of size $32 \times 32 \times 3$. This dataset was chosen due to its small size.

721 **Cats and Dogs.** The cats and dogs dataset is a subset of CIFAR-10. It uses only the 2 relevant classes,
 722 to create a binary problem. Each image is of size $32 \times 32 \times 3$. The dataset is divided to 20000 train

723 images (10000 per class) and 2000 test images (1000 per class). This dataset is used when a binary
724 problem is required.

725 **ImageNet-20.** The ImageNet-20 dataset is a subset of ImageNet containing 20 classes. This data
726 resembles ImageNet in terms of image resolution and data variability, but contains a smaller number
727 of examples in order to reduce computation time. The dataset contains 26000 train images (1300 per
728 class) and 1000 test images (50 per class). The choice of the 20 classes was arbitrary, and contained
729 the following classes: boa constrictor, jellyfish, American lobster, little blue heron, Shih-Tzu, scotch
730 terrier, Chesapeake Bay retriever, komondor, snow leopard, tiger, long-horned beetle, warthog, cab,
731 holster, remote control, toilet seat, pretzel, fig, burrito and toilet tissue.

732 **Frequency dataset** A binary 2D dataset, used in Section 4.3, to examine the effects of spectral bias
733 in classification. The data is define by the mapping $\lambda : [-1, 1] \rightarrow \mathbb{R}$ given in (24) by

$$\lambda(z) = \sum_{i=1}^m \sin(2\pi\kappa_i z + \varphi_i) := \sum_{i=1}^m \text{freq}_i(z)$$

734 with frequencies $\kappa = (\kappa_1, \kappa_2, \dots, \kappa_m)$ and corresponding phases $\phi = (\varphi_1, \varphi_2, \dots, \varphi_m)$. The classifi-
735 cation rule is defined by $\lambda(z) \leq 0$.

736 In our experiments, we chose $m = 10$, with frequencies $\kappa_1 = 0, \kappa_2 = 1, \kappa_3 = 2, \dots, \kappa_{10} = 9$. Other
737 choices of m yielded similar qualitative results. The phases were chosen randomly between 0 and 2π ,
738 and were set to be: $\varphi_1 = 0, \varphi_2 = 3.46, \varphi_3 = 5.08, \varphi_4 = 0.45, \varphi_5 = 2.10, \varphi_6 = 1.4, \varphi_7 = 5.36,$
739 $\varphi_8 = 0.85, \varphi_9 = 5.9, \varphi_{10} = 5.16$. As the first frequency is $\kappa_1 = 0$, the choice of φ_0 does not matter,
740 and is set to 0. The dataset contained 10000 training points, and 1000 test points, all uniformly
741 distributed in the first dimension between -1 and 1 and in the second dimension between -2π and
742 2π . The labels were set to be either 0 or 1, in order to achieve perfect separation with the classification
743 rule $\lambda(z)$.