

## 440 A Rotation Invariance w.r.t. the Initialized Weights

441 In this paper, we analyze neural networks trained on high-dimensional data that lies on a low dimensional  
 442 linear subspace denoted by  $P$ . We assume that the dimension of  $P$  is  $d - \ell$ . Throughout the paper  
 443 it will be more convenient to analyze data which lies on the subspace  $M = \text{span}(\{e_1, \dots, e_{d-\ell}\})$ ,  
 444 because then the ‘‘off manifold’’ directions correspond exactly to certain coordinates of the input. In  
 445 this section we show that we can essentially analyze the data as if it is rotated to lie on  $M$ , and it  
 446 would imply the same consequences as the original data from  $P$ .

447 **Theorem A.1.** *Let  $P \subseteq \mathbb{R}^d$  be a subspace of dimension  $d - \ell$ , and let  $M = \text{span}\{e_1, \dots, e_{d-\ell}\}$ .  
 448 Let  $R$  be an orthogonal matrix such that  $R \cdot P = M$ , let  $X \subseteq P$  be a training dataset and let  
 449  $X_R = \{R \cdot x : x \in X\}$ . Assume we train a neural network  $N(x) = \sum_{i=1}^m u_i \sigma(w_i^\top x)$  as explained  
 450 in Section 3, and denote by  $N^X$  and  $N^{X_R}$  the network trained on  $X$  and  $X_R$  respectively for the  
 451 same number of iterations. Let  $x_0 \in P$ , then we have:*

- 452 1. *W.p.  $p$  (over the initialization) we have  $\left\| \Pi_{P^\perp} \left( \frac{\partial N^X(x_0)}{\partial x} \right) \right\| \geq c$  (resp.  $\leq c$ ) for some  $c \in \mathbb{R}$ ,*  
 453 *iff w.p.  $p$  also  $\left\| \Pi_{M^\perp} \left( \frac{\partial N^{X_R}(Rx_0)}{\partial x} \right) \right\| \geq c$  (resp.  $\leq c$ ).*
- 454 2. *For any  $c, p \geq 0$ , w.p.  $p$  (over the initialization) there exists  $z \in P^\perp$  with  $\|z\| = c$  such that*  
 455  *$\text{sign}(N^X(x_0 + z)) \neq \text{sign}(N^X(x_0))$ , iff w.p.  $p$  there exists  $z' \in M^\perp$  with  $\|z'\| = c$  such*  
 456 *that  $\text{sign}(N^{X_R}(Rx_0 + z')) \neq \text{sign}(N^{X_R}(Rx_0))$ .*

457 *Proof.* Denote by  $\mathbf{w}_{1:m} := (w_1, \dots, w_m)$  and by  $R\mathbf{w}_{1:m} = (Rw_1, \dots, Rw_m)$ . Let  $\mathbf{w}_{1:m}^{(t)}$  the  
 458 weights of the network trained on the dataset  $X$  where  $\mathbf{w}_{1:m}^{(0)}$  is some initialization, and  $\tilde{\mathbf{w}}_{1:m}^{(t)} =$   
 459  $(\tilde{w}_1^{(t)}, \dots, \tilde{w}_m^{(t)})$  the weights of the network trained on  $X_R$  and initialized at  $R\mathbf{w}_{1:m}^{(0)}$ . In the proof,  
 460 when taking derivatives w.r.t. the  $w_i$ ’s we will explicitly write  $N(x, \mathbf{w}_{1:m})$ .

461 We first show by induction on the number of training steps that  $\tilde{\mathbf{w}}_{1:m}^{(t)} = R\mathbf{w}_{1:m}^{(t)}$ . For  $t = 0$  it is clear  
 462 by the assumption on the initialization. Assume it is true for  $t$ , then we have for some  $x \in X$ :

$$\begin{aligned} \frac{\partial N(Rx, \tilde{\mathbf{w}}_{1:m}^{(t)})}{\partial w_i} &= u_i \sigma'(\langle \tilde{w}_i^{(t)}, Rx \rangle) Rx \\ &= u_i \sigma'(\langle Rw_i^{(t)}, Rx \rangle) Rx \\ &= u_i \sigma'(\langle w_i^{(t)}, x \rangle) Rx \\ &= R \cdot \frac{\partial N(x, \mathbf{w}_{1:m}^{(t)})}{\partial w_i}. \end{aligned}$$

463 This is true for every  $i \in [m]$  and for every  $x \in X$ . Also note that by our induction assumption we  
 464 have:

$$N(x, \mathbf{w}_{1:m}^{(t)}) = \sum_{i=1}^m u_i \sigma(\langle w_i^{(t)}, x \rangle) = \sum_{i=1}^m u_i \sigma(\langle Rw_i^{(t)}, Rx \rangle) = N(Rx, \tilde{\mathbf{w}}_{1:m}^{(t)}). \quad (1)$$

465 Finally, the derivative of the loss on a single data point  $x \in X$  with label  $y$  can be written as:

$$\frac{\partial L(N(x, \mathbf{w}_{1:m}^{(t)}) \cdot y)}{\partial w_i} = L'(N(x, \mathbf{w}_{1:m}^{(t)}) \cdot y) \cdot \frac{\partial N(x, \mathbf{w}_{1:m}^{(t)})}{\partial w_i},$$

466 where the first term depends only on the value of  $N(x, \mathbf{w}_{1:m}^{(t)})$ . Hence, taking a single gradient step of  
 467  $N$  with weights  $\mathbf{w}_{1:m}^{(t)}$  and dataset  $X$  will change the weights by the same term up to multiplication  
 468 by  $R$  as if taking a gradient step with weights  $\tilde{\mathbf{w}}_{1:m}^{(t)}$  and dataset  $X_R$ . This finishes the induction.

469 Let  $\mathbf{w}_{1:m}^{(0)}$  be an initialization for the training of  $N^X$ , where there exists  $z \in P^\perp$  with  $\|z\| = c$  such  
 470 that  $\text{sign}(N^X(x_0 + z)) \neq \text{sign}(N^X(x_0))$ . Then, by Eq. (1) the initialization  $R\mathbf{w}_{1:m}^{(0)}$  for the training  
 471 of  $N^{X_R}$  is such that for  $z' = Rz$  we have  $\|z'\| = c$  and  $\text{sign}(N^{X_R}(Rx_0 + z')) \neq \text{sign}(N^{X_R}(Rx_0))$ .  
 472 This argument holds also in the opposite direction. Let  $A \subseteq \{\mathbf{w}_{1:m} \in \mathbb{R}^{d-m}\}$  be the set of all

473 initializations to  $N^X$  where there exists  $z \in P^\perp$  with  $\|z\| = c$  such that  $\text{sign}(N^X(x_0 + z)) \neq$   
474  $\text{sign}(N^X(x_0))$ , then by the above the set  $R \cdot A = \{R\mathbf{w}_{1:m} : \mathbf{w}_{1:m} \in A\}$  are exactly all the  
475 initializations to  $N^{XR}$  where there exists  $z' \in M^\perp$  with  $\|z'\| = c$  such that  $\text{sign}(N^{XR}(Rx_0 + z')) \neq$   
476  $\text{sign}(N^{XR}(Rx_0))$ . Since we initialize the  $w_i$ 's using a Gaussian initialization which is spherically  
477 symmetric, we have that  $\Pr(A) = \Pr(RA)$ . This proves item (2). Item (1) follows from similar  
478 arguments (which we do not repeat for conciseness).  $\square$

479 Under the assumption that the data lies on  $M = \text{span}\{e_1, \dots, e_{d-\ell}\}$ , and no regularization is used,  
480 we can show that the weights of the first layer projected on  $M^\perp$  do not change during training. This  
481 is an essential part of the proofs, as it allows us to analyze those weights as random Gaussian vectors,  
482 and apply concentration bounds on them.

483 **Theorem A.2.** *Let  $M = \text{span}\{e_1, \dots, e_{d-\ell}\}$ . Assume we train a neural network  $N(x, \mathbf{w}_{1:m}) :=$   
484  $\sum_{i=1}^m u_i \sigma(w_i^\top x)$  as explained in Section 3 (where  $\mathbf{w}_{1:m} = (w_1, \dots, w_m)$ ). Denote by  $\hat{w} :=$   
485  $\Pi_{M^\perp}(w)$  for  $w \in \mathbb{R}^d$ , then after training, for each  $i \in [m]$ ,  $\hat{w}_i$  did not change from their initial value.*

486 *Proof.* Note that for each  $i \in [m]$  and  $x \in M$  we have:

$$\Pi_{M^\perp} \left( \frac{\partial N(x, \mathbf{w}_{1:m})}{\partial w_i} \right) = \Pi_{M^\perp} (u_i \sigma'(w_i^\top x) x) = u_i \sigma'(w_i^\top x) \hat{x} = \mathbf{0}.$$

487 Taking the derivative of the loss we have:

$$\begin{aligned} \Pi_{M^\perp} \left( \frac{\partial L(N(x, \mathbf{w}_{1:m}) \cdot y)}{\partial w_i} \right) &= \Pi_{M^\perp} \left( L'(N(x, \mathbf{w}_{1:m}) \cdot y) \cdot \frac{\partial N(x, \mathbf{w}_{1:m})}{\partial w_i} \right) \\ &= L'(N(x, \mathbf{w}_{1:m}) \cdot y) \cdot \Pi_{M^\perp} \left( \frac{\partial N(x, \mathbf{w}_{1:m})}{\partial w_i} \right) = \mathbf{0}. \end{aligned}$$

488 The above calculation did not depend on the specific value of the  $w_i$ 's. Hence, the value of the  $\hat{w}_i$ 's  
489 for every  $i \in [m]$  did not change during training from their initial value.  $\square$

## 490 B Proofs from Section 4

491 Before proving the main theorem, we will first need the next two lemmas about the concentration of  
492 Gaussian random variables:

493 **Lemma B.1.** *Let  $w \in \mathbb{R}^n$  such that  $w \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_n)$ . Then:*

$$\mathbb{P} \left[ \|w\|^2 \leq \frac{1}{2} \sigma^2 n \right] \leq e^{-\frac{n}{16}}.$$

494 *Proof.* Note that  $\left\| \frac{w}{\sigma} \right\|^2$  has the Chi-squared distribution. A concentration bound by Laurent and  
495 Massart [Laurent and Massart, 2000, Lemma 1] implies that for all  $t > 0$  we have

$$\Pr \left[ n - \left\| \frac{w}{\sigma} \right\|^2 \geq 2\sqrt{nt} \right] \leq e^{-t}.$$

496 Plugging-in  $t = \frac{n}{16}$ , we get

$$\Pr \left[ n - \left\| \frac{w}{\sigma} \right\|^2 \geq \frac{1}{2} n \right] = \Pr \left[ \left\| \frac{w}{\sigma} \right\|^2 \leq \frac{1}{2} n \right] \leq e^{-n/16}.$$

497 Thus, we have

$$\Pr \left[ \|w\| \leq \sigma \sqrt{\frac{n}{2}} \right] \leq e^{-n/16}.$$

498  $\square$

499 **Lemma B.2.** *Let  $w_1, \dots, w_m \in \mathbb{R}^n$  such that for all  $i \in [m]$ ,  $w_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_n)$ , then we have:*

$$\mathbb{P} \left[ \left\| \sum_{i=1}^m w_i \right\|^2 \leq \frac{1}{2} m \sigma^2 n \right] \leq e^{-\frac{n}{16}}.$$

500 *Proof.* We denote the  $j$ -th coordinate of the vector  $w_i \in \mathbb{R}^n$  by  $w_{i,j}$ . Note, for any  $i \in \{1, \dots, m\}$   
501 and  $j \in \{1, \dots, n\}$  we have  $w_{i,j} \sim \mathcal{N}(0, \sigma^2)$ . We denote by  $s$  the sum vector  $s := \sum_{i=1}^m w_i$ , and by  
502  $s_j$  the  $j$ -th coordinate of  $s$ . By this definition,  $s_j = \sum_{i=1}^m w_{i,j}$  is a sum of  $m$  independent Gaussian  
503 variables and therefore also a Gaussian variable. Particularly,  $s \sim \mathcal{N}(\mathbf{0}, m\sigma^2 I_n)$ . We use Lemma B.1  
504 with variance  $m\sigma^2$  and get that:

$$\mathbb{P} \left[ \left\| \sum_{i=1}^m w_i \right\|^2 \leq \frac{1}{2} m\sigma^2 n \right] \leq e^{-\frac{n}{16}} .$$

505 □

506 We are now ready to prove the main theorem of this section:

507 *Proof of Theorem 4.1.* Let  $M = \text{span}\{e_1, \dots, e_{d-\ell}\}$ . By Theorem A.1(1), given a training dataset  
508  $X \subseteq P$ , it is enough to consider a training set  $X_R = \{Rx : x \in X\}$ , where  $R$  is an orthogonal  
509 matrix such that  $R \cdot P = M$ , and training is done over  $X_R$ . From now on, we assume that the training  
510 data, as well as  $x_0$  lie on  $M$ , and the consequences of this proof would also imply for a dataset  $X$   
511 and  $x_0 \in P$ .

512 The projection of the gradient on  $M^\perp$  is equal to:

$$\Pi_{M^\perp} \left( \frac{\partial N(x_0)}{\partial x} \right) = \Pi_{M^\perp} \left( \sum_{i=1}^m u_i w_i \mathbb{1}_{\langle w_i, x_0 \rangle \geq 0} \right) = \sum_{i=1}^m \Pi_{M^\perp} (u_i w_i) \mathbb{1}_{i \in S} = \sum_{i \in S} \Pi_{M^\perp} (u_i w_i) .$$

513 Denote by  $\hat{w}_i = (w_i)_{d-\ell+1:d}$ , the last  $\ell$  coordinates of  $w_i$ . By Theorem A.2 we get that for every  
514  $i \in [m]$ ,  $\hat{w}_i$  did not change from their initial value during training.

515 Recall that we initialized  $\hat{w}_i \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_\ell)$ . Note that the set  $S$  is independent of the value of the  
516  $\hat{w}_i$ 's. This is because  $\hat{w}_i$  does not effect the training, hence will not effect  $w_i - \Pi_{M^\perp}(w_i)$ . Also, after  
517 choosing  $x_0$  we have  $\langle \hat{w}_i, \hat{x}_0 \rangle = 0$ , since  $\hat{x}_0 = \mathbf{0}$ , which means that the choice of  $S$  is independent  
518 of the  $\hat{w}_i$ 's. We can conclude that the random variables  $\hat{w}_i$  for  $i \in S$  are sampled independently.

519 Note, since for all  $i \in \{1, \dots, m\}$ ,  $|u_i| = \frac{1}{\sqrt{m}}$  and they are not trained, we get that  $u_i \hat{w}_i$  are also  
520 Gaussian random variables with the same mean, and variance multiplied by  $\frac{1}{m}$ . Therefore, from  
521 Lemma B.2 we get that w.p.  $\geq 1 - e^{-\ell/16}$ :

$$\left\| \sum_{i \in S} u_i \hat{w}_i \right\| \geq \sqrt{\frac{1}{2}} \sqrt{\frac{kl}{dm}} .$$

522 Combining the above, we get:

$$\left\| \Pi_{M^\perp} \left( \frac{\partial N(x_0)}{\partial x} \right) \right\| \geq \sqrt{\frac{1}{2}} \sqrt{\frac{kl}{dm}} .$$

523 □

## 524 C Proofs from Section 5

525 Before proving the main theorem, we prove a few lemmas about concentration of Gaussian random  
526 variables:

527 **Lemma C.1.** *Let  $w \in \mathbb{R}^n$  with  $w \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_n)$ . Then:*

$$\Pr \left[ \|w\|^2 \geq 2\sigma^2 n \right] \leq e^{-\frac{n}{16}} .$$

528 *Proof.* Note that  $\left\|\frac{w}{\sigma}\right\|^2$  has the Chi-squared distribution. A concentration bound by Laurent and  
 529 Massart [Laurent and Massart, 2000, Lemma 1] implies that for all  $t > 0$  we have

$$\Pr \left[ \left\| \frac{w}{\sigma} \right\|^2 - n \geq 2\sqrt{nt} + 2t \right] \leq e^{-t}.$$

530 Plugging-in  $t = \frac{n}{16}$ , we get

$$\Pr \left[ \left\| \frac{w}{\sigma} \right\|^2 \geq 2n \right] \leq \Pr \left[ \left\| \frac{w}{\sigma} \right\|^2 - n \geq n/2 + n/8 \right] \leq e^{-n/16}.$$

531 Thus, we have

$$\Pr \left[ \|w\| \geq \sigma\sqrt{2n} \right] \leq e^{-n/16}.$$

532

□

533 **Lemma C.2.** Let  $u \in \mathbb{R}^n$ , and  $v \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_n)$ . Then, for every  $t > 0$  we have

$$\Pr [ |\langle u, v \rangle| \geq \|u\| t ] \leq 2 \exp \left( -\frac{t^2}{2\sigma^2} \right).$$

534 *Proof.* We first consider  $\langle \frac{u}{\|u\|}, v \rangle$ . As the distribution  $\mathcal{N}(\mathbf{0}, \sigma^2 I_n)$  is rotation invariant, one can rotate

535  $u$  and  $v$  to get  $\tilde{u}$  and  $\tilde{v}$  such that  $\frac{\tilde{u}}{\|\tilde{u}\|} = e_1$ , the first standard basis vector and  $\langle \frac{u}{\|u\|}, v \rangle = \langle \frac{\tilde{u}}{\|\tilde{u}\|}, \tilde{v} \rangle$ .

536 Note,  $v$  and  $\tilde{v}$  have the same distribution. We can see that  $\langle \frac{\tilde{u}}{\|\tilde{u}\|}, \tilde{v} \rangle \sim \mathcal{N}(0, \sigma^2)$  since it is the first

537 coordinate of  $\tilde{v}$ . By a standard tail bound, we get that for  $t > 0$ :

$$\Pr \left[ \left| \langle \frac{u}{\|u\|}, v \rangle \right| \geq t \right] = \Pr \left[ \left| \langle \frac{\tilde{u}}{\|\tilde{u}\|}, \tilde{v} \rangle \right| \geq t \right] = \Pr [ |\tilde{v}_1| \geq t ] \leq 2 \exp \left( -\frac{t^2}{2\sigma^2} \right).$$

538 Therefore

$$\Pr [ |\langle u, v \rangle| \geq \|u\| t ] \leq 2 \exp \left( -\frac{t^2}{2\sigma^2} \right).$$

539

□

540 **Lemma C.3.** Let  $u \sim \mathcal{N}(\mathbf{0}, \sigma_1^2 I_n)$ , and  $v \sim \mathcal{N}(\mathbf{0}, \sigma_2^2 I_n)$ . Then, for every  $t > 0$  we have

$$\Pr \left[ |\langle u, v \rangle| \geq \sigma_1 \sqrt{2nt} \right] \leq e^{-n/16} + 2e^{-t^2/2\sigma_2^2}.$$

541 *Proof.* Using Lemma C.1 we get that w.p.  $\leq e^{-n/16}$  we have  $\|u\| \geq \sigma_1 \sqrt{2n}$ . Moreover, by

542 Lemma C.2, w.p.  $\leq 2 \exp \left( -\frac{t^2}{2\sigma_2^2} \right)$  we have  $|\langle u, v \rangle| \geq \|u\| t$ . By the union bound, we get

$$\Pr \left[ |\langle u, v \rangle| \geq \sigma_1 \sqrt{2nt} \right] \leq \Pr \left[ \|u\| \geq \sigma_1 \sqrt{2n} \right] + \Pr [ |\langle u, v \rangle| \geq \|u\| t ] \leq e^{-n/16} + 2 \exp \left( -\frac{t^2}{2\sigma_2^2} \right).$$

543

□

544 We are now ready to prove the main theorem of this section:

545 **Theorem 5.1.** By Theorem A.1(2), we can assume w.l.o.g. that  $P = M = \text{span}\{e_1, \dots, e_{d-\ell}\}$ .

546 We also assume w.l.o.g. that  $y_0 = 1$ , the case  $y_0 = -1$  is proved in a similar manner. Denote by

547  $\bar{w} := (w)_{d-\ell+1:d}$ , the last  $\ell$  coordinates of  $w$ . By Theorem A.2 we have that  $\bar{w}_i$  have not changed

548 after training from their initial value.

549 We can write  $N(x_0 + z)$  as:

$$\begin{aligned}
N(x_0 + z) &= \sum_{i=1}^m u_i \sigma(\langle w_i, x_0 \rangle + \langle w_i, z \rangle) \\
&= \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle + \langle w_i, z \rangle) + \sum_{i \in I_+} u_i \sigma(\langle w_i, x_0 \rangle + \langle w_i, z \rangle) \\
&= \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle + \langle \bar{w}_i, \bar{z} \rangle) + \sum_{i \in I_+} u_i \sigma(\langle w_i, x_0 \rangle + \langle \bar{w}_i, \bar{z} \rangle) \tag{2}
\end{aligned}$$

550 where the last equality is since  $(z)_{1:d-\ell} = \mathbf{0}$ , hence  $\langle w, z \rangle = \langle \bar{w}, \bar{z} \rangle$  for every  $w \in \mathbb{R}^d$ . We will  
551 bound each term of the above separately.

552 For the first term in Eq. (2), where  $i \in I_-$  we can write:

$$\langle \bar{w}_i, \bar{z} \rangle = \alpha \|\bar{w}_i\|^2 + \alpha \langle \bar{w}_i, \sum_{j \neq i} \text{sign}(u_j) \bar{w}_j \rangle.$$

553 By our assumptions,  $\bar{w}_i \sim \mathcal{N}(\mathbf{0}, \frac{1}{d} I_\ell)$  and  $\sum_{j \neq i} \text{sign}(u_j) \bar{w}_j \sim \mathcal{N}(\mathbf{0}, \frac{m-1}{d} I_\ell)$ , since it is a sum of  
554  $m-1$  i.i.d. Gaussian random variables, which are also symmetric hence multiplying them by  $-1$   
555 does not change their distribution. From Lemma B.1 we get w.p.  $\geq 1 - e^{-\ell/16}$  that

$$\alpha \cdot \|\bar{w}_i\|^2 \geq \alpha \cdot \frac{\ell}{2d}.$$

556 From Lemma C.3, and using  $t = \sqrt{\frac{(m-1) \log(dm^2)}{d}}$  we get w.p.  $\geq 1 - e^{-\ell/16} + 2e^{-t^2 d/2(m-1)} =$   
557  $1 - e^{-\ell/16} + 2m^{-1} d^{-1/2}$  that

$$\begin{aligned}
\langle \bar{w}_i, \sum_{j \neq i} \text{sign}(u_j) \bar{w}_j \rangle &\leq \frac{1}{\sqrt{d}} t \sqrt{2\ell} \\
&= \frac{1}{d} \cdot \sqrt{2\ell(m-1) \log(m^2 d)}. \tag{3}
\end{aligned}$$

558 Applying union bound over the above two events, and for every  $i \in I_-$ , we get w.p.  $\geq 1 -$   
559  $2(me^{-\ell/16} + d^{-1/2})$  that:

$$\langle \bar{w}_i, \bar{z} \rangle \geq \frac{\alpha \ell}{2d} - \frac{\alpha}{d} \sqrt{2\ell(m-1) \log(m^2 d)}.$$

560 For the second term in Eq. (2), where  $i \in I_+$  we can write in a similar way:

$$\langle \bar{w}_i, \bar{z} \rangle = -\alpha \|\bar{w}_i\|^2 + \alpha \langle \bar{w}_i, \sum_{j \neq i} \text{sign}(u_j) \bar{w}_j \rangle.$$

561 Using the same argument as above, we get w.p.  $\geq 1 - 2(me^{-\ell/16} + d^{-1/2})$  that:

$$\langle \bar{w}_i, \bar{z} \rangle \leq -\frac{\alpha \ell}{2d} + \frac{\alpha}{d} \sqrt{2\ell(m-1) \log(m^2 d)}.$$

562 By assuming that  $\ell \geq 8(m-1) \log(m^2 d)$  we get that  $\langle \bar{w}_i, z \rangle \leq 0$ . Denote  $C := \frac{\alpha \ell}{2d} -$   
563  $\frac{\alpha}{d} \sqrt{2\ell(m-1) \log(m^2 d)}$ , then going back to Eq. (2), using the above bounds and applying union  
564 bound, we get w.p.  $\geq 1 - 4(me^{-\ell/16} + d^{-1/2})$  that:

$$\begin{aligned}
N(x_0 + z) &\leq \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle + C) + \sum_{i \in I_+} u_i \sigma(\langle w_i, x_0 \rangle) \\
&= \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle + C) + \sum_{i \in I_+} u_i \sigma(\langle w_i, x_0 \rangle) + \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle) - \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle) \\
&= \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle + C) - \sum_{i \in I_-} u_i \sigma(\langle w_i, x_0 \rangle) + N(x_0) \\
&= \sum_{i \in I_-} u_i (\sigma(\langle w_i, x_0 \rangle + C) - \sigma(\langle w_i, x_0 \rangle)) + N(x_0).
\end{aligned}$$

565 Define  $F_- := \{i \in I_- : \langle w_i, x_0 \rangle \geq 0\}$ , and  $k_- = |F_-|$ . We have that:

$$\begin{aligned} \sum_{i \in I_-} u_i (\sigma(\langle w_i, x_0 \rangle + C) - \sigma(\langle w_i, x_0 \rangle)) &\leq \sum_{i \in F_-} u_i (\sigma(\langle w_i, x_0 \rangle + C) - \sigma(\langle w_i, x_0 \rangle)) \\ &= \sum_{i \in F_-} u_i C = -\frac{k_- C}{\sqrt{m}}, \end{aligned}$$

566 where the first inequality is since we only sum over negative terms, and the second inequality is since  
567 both  $\langle w_i, x_0 \rangle \geq 0$  (because  $i \in F_-$ ) and  $C \geq 0$  (because  $\ell \geq 32(m-1) \log(m^2 d)$ ). Combining all  
568 of the above, we get that:

$$N(x_0 + z) \leq -\frac{k_- C}{\sqrt{m}} + N(x_0). \quad (4)$$

569 By our assumption that  $\ell \geq 32(m-1) \log(m^2 d)$  we have that

$$\begin{aligned} C &= \alpha \left( \frac{1}{2} \frac{\ell}{d} - \sqrt{2} \sqrt{m-1} \frac{\sqrt{\ell}}{d} \sqrt{\log(dm^2)} \right) \\ &= \frac{\alpha \sqrt{\ell}}{d} \left( \frac{\sqrt{\ell}}{2} - \sqrt{2(m-1) \log(m^2 d)} \right) \\ &\geq \frac{\alpha \ell}{4d}. \end{aligned}$$

570 Plugging in  $C$  and  $\alpha = \frac{8\sqrt{md}N(x_0)}{k_- \ell}$  to Eq. (4) we get that:

$$\begin{aligned} N(x_0 + z) &\leq -\frac{k_- C}{\sqrt{m}} + N(x_0) \\ &\leq -\frac{k_-}{\sqrt{m}} \cdot \frac{\ell}{4d} \cdot \frac{8\sqrt{md}N(x_0)}{k_- \ell} + N(x_0) = -N(x_0) < 0, \end{aligned}$$

571 and in particular  $\text{sign}(N(x_0)) \neq \text{sign}(N(x_0 + z))$ .

572 We are left with calculating the norm of  $z$ :

$$\begin{aligned} \|z\| &= \alpha \cdot \left\| \sum_{i \in I_-} \Pi_{M^\perp}(w_i) - \sum_{i \in I_+} \Pi_{M^\perp}(w_i) \right\| \\ &= \alpha \cdot \left\| \sum_{i=1}^m -\text{sign}(u_i) \Pi_{M^\perp}(w_i) \right\| \\ &= \alpha \cdot \left\| \sum_{i=1}^m -\text{sign}(u_i) \bar{w}_i \right\|. \end{aligned}$$

573 Since for each  $i \in [m]$ ,  $\bar{w}_i \sim \mathcal{N}(\mathbf{0}, \frac{1}{d} I_\ell)$ , then  $-\text{sign}(u_i) \bar{w}_i$  also have the same distribution, because  
574 this is a symmetric distribution. Hence,  $\sum_{i=1}^m -\text{sign}(u_i) \bar{w}_i \sim \mathcal{N}(\mathbf{0}, \frac{m}{d} I_\ell)$  as a sum of Gaussian  
575 random variables. Using Lemma C.1 we get w.p  $\geq 1 - e^{-\ell/16}$  that  $\|\sum_{i=1}^m -\text{sign}(u_i) \bar{w}_i\|^2 \leq \frac{2m\ell}{d}$ .  
576 Plugging in  $\alpha$  we get that:

$$\|z\| \leq \sqrt{\frac{2m\ell}{d}} \cdot \frac{8\sqrt{md}N(x_0)}{k_- \ell} = 8\sqrt{2}N(x_0) \cdot \frac{m}{k_-} \cdot \sqrt{\frac{d}{\ell}}.$$

577

□

578 **D Proofs for Section 6**

579 For proving the main theorem, we will use the following lemma that upper bounds the norm of a sum  
580 of Gaussian random variables:

581 **Lemma D.1.** *Let  $w_1, \dots, w_m \in \mathbb{R}^n$  such that for all  $i \in [m]$ ,  $w_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_n)$ , then we have:*

$$\mathbb{P} \left[ \left\| \sum_{i=1}^m w_i \right\|^2 \geq 2m\sigma^2 n \right] \leq e^{-\frac{n}{16}}$$

582 *Proof.* We denote the  $j$ -th coordinate of the vector  $w_i \in \mathbb{R}^n$  by  $w_{i,j}$ . Note, for any  $i \in [m]$  and  
583  $j \in [n]$  we have  $w_{i,j} \sim \mathcal{N}(0, \sigma^2)$ . We denote by  $s$  the sum vector  $s := \sum_{i=1}^m w_i$ , and by  $s_j$  the  $j$ -th

584 coordinate of  $s$ . By this definition,  $s_j = \sum_{i=1}^m w_{i,j}$  is a sum of  $m$  independent Gaussian variables and  
585 therefore also a Gaussian variable. Therefore,  $s \sim \mathcal{N}(\mathbf{0}, m\sigma^2 I_n)$ . We use Lemma C.1 with variance  
586  $m\sigma^2$  and get that:

$$\mathbb{P} \left[ \left\| \sum_{i=1}^m w_i \right\|^2 \geq 2m\sigma^2 n \right] \leq e^{-\frac{n}{16}}.$$

587 □

588 We now prove the main theorem of this section:

589 *Proof of Theorem 6.1.* Similar to the lower bound of the norm, let  $M = \text{span}\{e_1, \dots, e_{d-\ell}\}$ . By  
590 Theorem A.1(1), given a training dataset  $X \subseteq P$ , it is enough to consider a training set  $X_R = \{Rx : x \in X\}$ ,  
591 where  $R$  is an orthogonal matrix such that  $R \cdot P = M$ , and training is done over  $X_R$ . From  
592 now on, we assume that the training data, as well as  $x_0$  lie on  $M$ , and the consequences of this proof  
593 would also imply for a dataset  $X$  and  $x_0 \in P$ .

594 The projection of the gradient on  $M^\perp$  is equal to:

$$\Pi_{M^\perp} \left( \frac{\partial N(x_0)}{\partial x} \right) = \Pi_{M^\perp} \left( \sum_{i=1}^m u_i w_i \mathbb{1}_{\langle w_i, x_0 \rangle \geq 0} \right) = \sum_{i=1}^m \Pi_{M^\perp} (u_i w_i) \mathbb{1}_{i \in S} = \sum_{i \in S} \Pi_{M^\perp} (u_i w_i).$$

595 Denote by  $\hat{w}_i = (w_i)_{d-\ell+1:d}$ , the last  $\ell$  coordinates of  $w_i$ . By Theorem A.2 we get that for every  
596  $i \in [m]$ ,  $\hat{w}_i$  did not change from their initial value during training.

597 Recall that we initialized  $\hat{w}_i \sim \mathcal{N}(\mathbf{0}, \beta^2 I_\ell)$ . Note that the set  $S$  is independent of the value of the  
598  $\hat{w}_i$ 's. This is because  $\hat{w}_i$  does not effect the training, hence will not effect  $w_i - \Pi_{M^\perp}(w_i)$ . Also, after  
599 choosing  $x_0$  we have  $\langle \hat{w}_i, \hat{x}_0 \rangle = 0$ , since  $\hat{x}_0 = \mathbf{0}$ , which means that the choice of  $S$  is independent  
600 of the  $\hat{w}_i$ 's. We can conclude that the random variables  $\hat{w}_i$  for  $i \in S$  are sampled independently.

601 Therefore, from Lemma B.2 we get that w.p.  $\geq 1 - e^{-\ell/16}$ :

$$\left\| \sum_{i \in S} \hat{w}_i \right\| \leq \beta \sqrt{2k\ell}.$$

602 Note, since for all  $i \in [m]$ ,  $|u_i| = \frac{1}{\sqrt{m}}$  and they are not trained, we get w.p.  $\geq 1 - e^{-\ell/16}$  that:

$$\left\| \Pi_{M^\perp} \left( \frac{\partial N(x_0)}{\partial x} \right) \right\| \leq \beta \sqrt{\frac{2k\ell}{m}}.$$

603 □

604 **D.1 Explicit  $L_2$  regularization**

605 *Proof of Theorem 6.2.* As before, for this proof we rotate the data subspace  $P$  to lie on  $M =$   
 606  $\text{span}\{e_1, \dots, e_{d-\ell}\}$  and rotate the model's weights accordingly. For a dataset  $(x_1, y_1), \dots, (x_r, y_r),$   
 607 we train over the following objective:

$$\sum_{j=1}^r L(y_j \cdot N(x_j, \mathbf{w}_{1:m})) + \frac{1}{2} \lambda \|\mathbf{w}_{1:m}\|^2$$

608 In Theorem A.2, we showed for all  $(x_j, y_j)$  that if we train the model using the loss  $L$  we get:

$$\Pi_{M^\perp} \left( \frac{\partial L(N(x_j, \mathbf{w}_{1:m}) \cdot y_j)}{\partial w_i} \right) = 0$$

609 Now, we analyze the training process using the new loss which includes the regularization term. We  
 610 denote by  $w_i^{(t)}$  the weight vector  $w_i$  after  $t$  training steps, and by  $\hat{w}_i^{(t)} := \Pi_{M^\perp} \left( w_i^{(t)} \right)$  its projection  
 611 on the subspace orthogonal to  $M$ . We look at the projected gradient of  $w_i^{(t)}$  w.r.t. the loss:

$$\begin{aligned} & \Pi_{M^\perp} \left( \frac{\partial \sum_{j=1}^r L(N(x_j, \mathbf{w}_{1:m}^{(t)}) \cdot y_j)}{\partial w_i} + \frac{\partial \frac{1}{2} \lambda \|w_i^{(t)}\|^2}{\partial w_i} \right) = \\ & = \sum_{j=1}^r \Pi_{M^\perp} \left( \frac{\partial L(N(x_j, \mathbf{w}_{1:m}^{(t)}) \cdot y_j)}{\partial w_i} \right) + \Pi_{M^\perp} \left( \frac{\partial \frac{1}{2} \lambda \|w_i^{(t)}\|^2}{\partial w_i} \right) \\ & = \Pi_{M^\perp} \left( \frac{\partial \frac{1}{2} \lambda \|w_i^{(t)}\|^2}{\partial w_i} \right) \\ & = \Pi_{M^\perp} \left( \lambda w_i^{(t)} \right) \\ & = \lambda \hat{w}_i^{(t)}. \end{aligned}$$

612 For a training step of size  $\eta$ , using gradient descent we get that:

$$\hat{w}_i^{(t+1)} = \hat{w}_i^{(t)} - \eta \lambda \hat{w}_i^{(t)}.$$

613 Thus, after a total of  $T$  iteration of training we get that:

$$\hat{w}_i^{(T)} = (1 - \eta \lambda)^T \hat{w}_i^{(0)}.$$

614 Therefore, the projection of gradients after training onto  $P^\perp$  will be the same as if they were initialized  
 615 to  $\sim \mathcal{N} \left( 0, \frac{(1-\eta\lambda)^{2T}}{d} I_d \right)$  and trained using logistic loss without regularization. The rest of the proof  
 616 is the same as Theorem 6.1 for  $\beta = \frac{(1-\eta\lambda)^T}{\sqrt{d}}$ .  $\square$

617 **E Further Experiments and Experimental Details**

618 **E.1 Further Experiments**

619 In Figure 3 we present the boundary of a two-layer ReLU network trained over a 25-point dataset  
 620 on a two-dimensional linear subspace, similar to Figure 2. We train the networks until reaching a  
 621 constant positive margin. The difference between the figures is that in Figure 3 we initialize the  
 622 weights using the default PyTorch initialization, while in Figure 2 we initialized using a smaller scale

623 for the robustness effect to be smaller, and visualized more easily. The experiment in Figure 3 is  
 624 demonstrating an extreme robustness effect, occurring when using the standard settings.

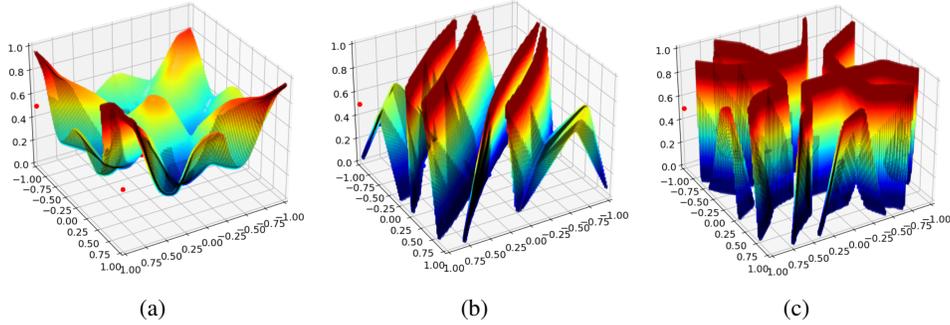


Figure 3: **Experiments on two-dimensional dataset.** We plot the dataset points and the decision boundary in 3 settings: (a) Vanilla trained network, (b) The network’s weights are initialized from a smaller variance distribution, and (c) Training with regularization. Colors are used to emphasise the values in the  $z$  axis.

625 In Figure 4 we go beyond the theory discussed in this paper, and present similar phenomena in all  
 626 three settings for a five-layer ReLU network. In Figure 4a we can see the boundary of the regularly  
 627 trained network within a small distance in  $P^\perp$  from the data points. In Figure 4b we use small  
 628 initialization for all five layers, and present a boundary almost orthogonal to the data manifold. In  
 629 Figure 4c, the boundary of a regularized trained network is in a similar form. This experiment  
 630 suggests that our theoretical results might be extended also to deeper networks, where all layers are  
 631 trained.

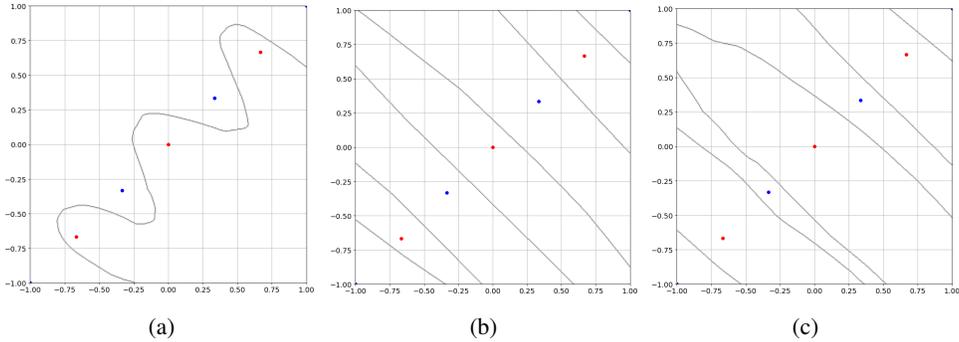


Figure 4: **Experiments on one-dimensional dataset with deep network.** We plot the dataset points and the decision boundary in 3 settings: (a) Vanilla trained network, (b) The network’s weights are initialized from a smaller variance distribution, and (c) Training with regularization.

632 **E.2 One-dimensional dataset experiment - 2 layer network (Figure 1)**

633 **Dataset** For all the three experiments we used a 7-point data set, spread equally on the two  
 634 dimensional line  $y = x$  from  $(-1, -1)$  to  $(1, 1)$ .

635 **Network** For all the three experiments we used two-layer ReLU network of width 100 with biases  
 636 in both layers. The weights of both layers were initialized using (1+3) default PyTorch initialization  
 637 for linear layers, (2) default initialization divided by 3.

638 **Training** We used train step of size 0.02 for (1+3) and 0.04 for (2). We trained both layers until the  
 639 margin reached 0.3. The losses we used were (1+2) Logistic loss, (3) Logistic loss with 0.005  $L_2$   
 640 regularization.

641 **E.3 Two-dimensional dataset experiment - smaller effect (Figure 2)**

642 **Dataset** For all the three experiments we used a 25-point data set, spread equally on a grid which  
643 lies on the  $z = 0.5$  axis.

644 **Network** For all the three experiments we used two-layer ReLU network of width 4000 with biases  
645 in both layers. The weights in the first layer were initialized in (1+3) from  $\mathcal{N}(\mathbf{0}, 1/3I_3)$ , and in (2)  
646 from  $\mathcal{N}(\mathbf{0}, 1/36I_3)$ . The weight of the output layer were initialized to the uniform distribution over  
647 the set  $\{-1, 1\}$ .

648 **Training** For all the experiments we trained both layers until the margin reached 0.3 and we used  
649 train step of size 0.002. The losses we used were (1+2) Logistic loss, (3) Logistic loss with 0.8  $L_2$   
650 regularization on the weights of the first layer.

651 **E.4 Two-dimensional dataset experiment (Figure 3)**

652 **Dataset** For all the three experiments we used a 25-point data set, spread equally on a grid which  
653 lies on the  $x - y$  axis.

654 **Network** For all the three experiments we used two-layer ReLU network of width 400 with biases  
655 in both layers. The weights in any layer were initialized using (1+3) default PyTorch initialization for  
656 linear layers, (2) default initialization divided by 3.

657 **Training** For (1) experiments we used train step of size 0.005, and for (2+3) we used step of size  
658 0.05. We trained both layers until the margin reached 0.1. The losses we used were (1+2) Logistic  
659 loss, (3) Logistic loss with 0.005  $L_2$  regularization.

660 **E.5 One-dimensional dataset experiment - 5 layer network (Figure 4)**

661 **Dataset** For all the three experiments we used a 7-point data set, spread equally on the two  
662 dimensional line  $y = x$  from  $(-1, -1)$  to  $(1, 1)$ .

663 **Network** For all the three experiments we used 5-layer ReLU network of width 100 with biases in  
664 all layers. The weights in any layer were initialized using (1+3) default PyTorch initialization for  
665 linear layers, (2) default initialization divided by 3.

666 **Training** For (1+3) experiments we used train step of size 0.02, and for (2) we used step of size  
667 0.06. we trained all layers until the margin reached 0.3. The losses we used were (1+2) Logistic loss,  
668 (3) Logistic loss with 0.01  $L_2$  regularization.